

## CONTENTS

<b>1</b>	<b>Math and Machine Learning Basics</b>	<b>7</b>
1.1	Linear Algebra (Quick Review) (Ch. 2) . . . . .	8
1.1.1	Example: Principal Component Analysis . . . . .	10
1.2	Probability & Information Theory (Quick Review) (Ch. 3) . . . . .	12
1.3	Numerical Computation (Ch. 4) . . . . .	14
1.4	Machine Learning Basics (Ch. 5) . . . . .	17
1.4.1	Estimators, Bias and Variance (5.4) . . . . .	17
1.4.2	Maximum Likelihood Estimation (5.5) . . . . .	19
1.4.3	Bayesian Statistics (5.6) . . . . .	20
1.4.4	Supervised Learning Algorithms (5.7) . . . . .	22
<b>2</b>	<b>Deep Networks: Modern Practices</b>	<b>23</b>
2.1	Deep Feedforward Networks (Ch. 6) . . . . .	24
2.1.1	Back-Propagation (6.5) . . . . .	25
2.2	Regularization for Deep Learning (Ch. 7) . . . . .	26
2.3	Optimization for Training Deep Models (Ch. 8) . . . . .	28
2.4	Convolutional Neural Networks (Ch. 9) . . . . .	32
2.5	Sequence Modeling (RNNs) (Ch. 10) . . . . .	35
2.5.1	Review: The Basics of RNNs . . . . .	35
2.5.2	RNNs as Directed Graphical Models . . . . .	40
2.5.3	Challenge of Long-Term Deps. (10.7) . . . . .	42
2.5.4	LSTMs and Other Gated RNNs (10.10) . . . . .	43
2.6	Applications (Ch. 12) . . . . .	44
2.6.1	Natural Language Processing (12.4) . . . . .	44
2.6.2	Neural Language Models (12.4.2) . . . . .	45
<b>3</b>	<b>Deep Learning Research</b>	<b>46</b>
3.1	Linear Factor Models (Ch. 13) . . . . .	47
3.2	Autoencoders (Ch. 14) . . . . .	50
3.3	Representation Learning (Ch. 15) . . . . .	51
3.4	Structured Probabilistic Models for DL (Ch. 16) . . . . .	52
3.4.1	Sampling from Graphical Models . . . . .	54
3.4.2	Inference and Approximate Inference . . . . .	54
3.5	Monte Carlo Methods (Ch. 17) . . . . .	56
3.6	Confronting the Partition Function (Ch. 18) . . . . .	58
3.7	Approximate Inference (Ch. 19) . . . . .	59
3.8	Deep Generative Models (Ch. 20) . . . . .	61
<b>4</b>	<b>Papers and Tutorials</b>	<b>66</b>
4.1	WaveNet . . . . .	69

4.2	Neural Style . . . . .	72
4.3	Neural Conversation Model . . . . .	74
4.4	NMT By Jointly Learning to Align & Translate . . . . .	76
4.4.1	Detailed Model Architecture . . . . .	77
4.5	Effective Approaches to Attention-Based NMT . . . . .	79
4.6	Using Large Vocabularies for NMT . . . . .	81
4.7	Candidate Sampling – TensorFlow . . . . .	84
4.8	Attention Terminology . . . . .	86
4.9	TextRank . . . . .	88
4.9.1	Keyword Extraction . . . . .	90
4.9.2	Sentence Extraction . . . . .	91
4.10	Simple Baseline for Sentence Embeddings . . . . .	92
4.11	Survey of Text Clustering Algorithms . . . . .	94
4.11.1	Distance-based Clustering Algorithms . . . . .	97
4.11.2	Probabilistic Document Clustering and Topic Models . . . . .	98
4.11.3	Online Clustering with Text Streams . . . . .	100
4.12	Deep Sentence Embedding Using LSTMs . . . . .	102
4.13	Clustering Massive Text Streams . . . . .	105
4.14	Supervised Universal Sentence Representations (InferSent) . . . . .	107
4.15	Dist. Rep. of Sentences from Unlabeled Data (FastSent) . . . . .	108
4.16	Latent Dirichlet Allocation . . . . .	110
4.17	Conditional Random Fields . . . . .	113
4.18	Attention Is All You Need . . . . .	116
4.19	Hierarchical Attention Networks . . . . .	120
4.20	Joint Event Extraction via RNNs . . . . .	123
4.21	Event Extraction via Bidi-LSTM Tensor NNs . . . . .	125
4.22	Reasoning with Neural Tensor Networks . . . . .	127
4.23	Language to Logical Form with Neural Attention . . . . .	128
4.24	Seq2SQL: Generating Structured Queries from NL using RL . . . . .	130
4.25	SLING: A Framework for Frame Semantic Parsing . . . . .	133
4.26	Poincaré Embeddings for Learning Hierarchical Representations . . . . .	135
4.27	Enriching Word Vectors with Subword Information (FastText) . . . . .	137
4.28	DeepWalk: Online Learning of Social Representations . . . . .	139
4.29	Review of Relational Machine Learning for Knowledge Graphs . . . . .	141
4.30	Fast Top-K Search in Knowledge Graphs . . . . .	144
4.31	Dynamic Recurrent Acyclic Graphical Neural Networks (DRAGNN) . . . . .	146
4.31.1	More Detail: Arc-Standard Transition System . . . . .	149
4.32	Neural Architecture Search with Reinforcement Learning . . . . .	150
4.33	Joint Extraction of Events and Entities within a Document Context . . . . .	152
4.34	Globally Normalized Transition-Based Neural Networks . . . . .	155
4.35	An Introduction to Conditional Random Fields . . . . .	158
4.35.1	Inference (Sec. 4) . . . . .	162
4.35.2	Parameter Estimation (Sec. 5) . . . . .	165
4.35.3	Related Work and Future Directions (Sec. 6) . . . . .	168

4.36	Co-sampling: Training Robust Networks for Extremely Noisy Supervision . . . . .	169
4.37	Hidden-Unit Conditional Random Fields . . . . .	170
4.37.1	Detailed Derivations . . . . .	172
4.38	Pre-training of Hidden-Unit CRFs . . . . .	177
4.39	Structured Attention Networks . . . . .	179
4.40	Neural Conditional Random Fields . . . . .	181
4.41	Bidirectional LSTM-CRF Models for Sequence Tagging . . . . .	183
4.42	Relation Extraction: A Survey . . . . .	184
4.43	Neural Relation Extraction with Selective Attention over Instances . . . . .	187
4.44	On Herding and the Perceptron Cycling Theorem . . . . .	189
4.45	Non-Convex Optimization for Machine Learning . . . . .	191
4.45.1	Non-Convex Projected Gradient Descent (3) . . . . .	194
4.46	Improving Language Understanding by Generative Pre-Training . . . . .	195
4.47	Deep Contextualized Word Representations . . . . .	196
4.48	Exploring the Limits of Language Modeling . . . . .	198
4.49	Connectionist Temporal Classification . . . . .	200
4.50	BERT . . . . .	202
4.51	Wasserstein is all you need . . . . .	204
4.52	Noise Contrastive Estimation . . . . .	206
4.52.1	Self-Normalized NCE . . . . .	208
4.53	Neural Ordinary Differential Equations . . . . .	210
4.54	On the Dimensionality of Word Embedding . . . . .	212
4.55	Generative Adversarial Nets . . . . .	213
4.56	A Framework for Intelligence and Cortical Function . . . . .	216
4.57	Large-Scale Study of Curiosity Driven Learning . . . . .	217
4.58	Universal Language Model Fine-Tuning for Text Classification . . . . .	218
4.59	The Marginal Value of Adaptive Gradient Methods in Machine Learning . . . . .	220
4.60	A Theoretically Grounded Application of Dropout in Recurrent Neural Networks . . . . .	221
4.61	Improving Neural Language Models with a Continuous Cache . . . . .	222
4.62	Protection Against Reconstruction and Its Applications in Private Federated Learning . . . . .	223
4.63	Context Dependent RNN Language Model . . . . .	225
4.64	Strategies for Training Large Vocabulary Neural Language Models . . . . .	226
4.65	Product quantization for nearest neighbor search . . . . .	228
4.66	Large Memory Layers with Product Keys . . . . .	229
4.67	Show, Ask, Attend, and Answer . . . . .	231
4.68	Did the Model Understand the Question? . . . . .	233
4.69	XLNet . . . . .	234
4.70	Transformer-XL . . . . .	236
4.71	Efficient Softmax Approximation for GPUs . . . . .	237
4.72	Adaptive Input Representations for Neural Language Modeling . . . . .	238
4.73	Neural Module Networks . . . . .	239
4.74	Learning to Compose Neural Networks for QA . . . . .	241
4.75	End-to-End Module Networks for VQA . . . . .	243
4.76	Fast Multi-language LSTM-based Online Handwriting Recognition . . . . .	245

4.77	Multi-Language Online Handwriting Recognition . . . . .	246
4.78	Modular Generative Adversarial Networks . . . . .	248
4.79	Transfer Learning from Speaker Verification to TTS . . . . .	250
<b>5</b>	<b>NLP with Deep Learning</b>	<b>251</b>
5.1	Word Vector Representations (Lec 2) . . . . .	252
5.2	GloVe (Lec 3) . . . . .	255
<b>6</b>	<b>Speech and Language Processing</b>	<b>258</b>
6.1	Introduction (Ch. 1 2nd Ed.) . . . . .	259
6.2	Morphology (Ch. 3 2nd Ed.) . . . . .	260
6.3	N-Grams (Ch. 6 2nd Ed.) . . . . .	261
6.4	Naive Bayes and Sentiment (Ch. 6 3rd Ed.) . . . . .	263
6.5	Hidden Markov Models (Ch. 9 3rd Ed.) . . . . .	265
6.6	POS Tagging (Ch. 10 3rd Ed.) . . . . .	268
6.7	Formal Grammars (Ch. 11 3rd Ed.) . . . . .	271
6.8	Vector Semantics (Ch. 15) . . . . .	272
6.9	Semantics with Dense Vectors (Ch. 16) . . . . .	275
6.10	Information Extraction (Ch. 21 3rd Ed) . . . . .	278
<b>7</b>	<b>Probabilistic Graphical Models</b>	<b>281</b>
7.1	Foundations (Ch. 2) . . . . .	282
7.1.1	Appendix . . . . .	284
7.1.2	L-BFGS . . . . .	287
7.1.3	Exercises . . . . .	290
7.2	The Bayesian Network Representation (Ch. 3) . . . . .	292
7.3	Undirected Graphical Models (Ch. 4) . . . . .	296
7.3.1	Exercises . . . . .	303
7.4	Local Probabilistic Models (Ch. 5) . . . . .	305
7.5	Template-Based Representations (Ch. 6) . . . . .	306
7.6	Gaussian Network Models (Ch. 7) . . . . .	309
7.7	Variable Elimination (Ch. 9) . . . . .	310
7.8	Clique Trees (Ch. 10) . . . . .	314
7.9	Inference as Optimization (Ch. 11) . . . . .	319
7.10	Parameter Estimation (Ch. 17) . . . . .	320
7.11	Partially Observed Data (Ch. 19) . . . . .	322
<b>8</b>	<b>Information Theory, Inference, and Learning Algorithms</b>	<b>324</b>
8.1	Introduction to Information Theory (Ch. 1) . . . . .	325
8.2	Probability, Entropy, and Inference (Ch. 2) . . . . .	327
8.2.1	More About Inference (Ch. 3 Summary) . . . . .	329
8.3	The Source Coding Theorem (Ch. 4) . . . . .	331

8.3.1	Data Compression and Typicality . . . . .	333
8.3.2	Further Analysis and Q&A . . . . .	335
8.4	Monte Carlo Methods (Ch. 29) . . . . .	337
8.5	Variational Methods (Ch. 33) . . . . .	339
<b>9</b>	<b>Machine Learning: A Probabilistic Perspective</b>	<b>341</b>
9.1	Probability (Ch. 2) . . . . .	342
9.1.1	Exercises . . . . .	343
9.2	Generative Models for Discrete Data (Ch. 3) . . . . .	344
9.2.1	Exercises . . . . .	348
9.3	Gaussian Models (Ch. 4) . . . . .	349
9.4	Bayesian Statistics (Ch. 5) . . . . .	353
9.5	Linear Regression (Ch. 7) . . . . .	355
9.6	Generalized Linear Models and the Exponential Family (Ch. 9) . . . . .	358
9.7	Mixture Models and the EM Algorithm (Ch. 11) . . . . .	361
9.8	Latent Linear Models (Ch. 12) . . . . .	364
9.9	Markov and Hidden Markov Models (Ch. 17) . . . . .	366
9.10	Undirected Graphical Models (Ch. 19) . . . . .	368
<b>10</b>	<b>Convex Optimization</b>	<b>369</b>
10.1	Convex Sets (Ch. 2) . . . . .	370
<b>11</b>	<b>Bayesian Data Analysis</b>	<b>372</b>
11.1	Probability and Inference (Ch. 1) . . . . .	373
11.2	Single-Parameter Models (Ch. 2) . . . . .	375
11.3	Asymptotics and Connections to Non-Bayesian Approaches (Ch. 4) . . . . .	378
11.4	Gaussian Process Models (Ch. 21) . . . . .	381
<b>12</b>	<b>Gaussian Processes for Machine Learning</b>	<b>383</b>
12.1	Regression (Ch. 2) . . . . .	384
<b>13</b>	<b>Blogs</b>	<b>387</b>
13.1	Conv Nets: A Modular Perspective . . . . .	388
13.2	Understanding Convolutions . . . . .	389
13.3	Deep Reinforcement Learning . . . . .	391
13.4	Deep Learning for Chatbots (WildML) . . . . .	393
13.5	Attentional Interfaces – Neural Perspective . . . . .	395
<b>14</b>	<b>Appendix</b>	<b>396</b>
14.1	Common Distributions and Models . . . . .	397
14.2	Math . . . . .	399
14.3	Matrix Cookbook . . . . .	408

14.4	Main Tasks in NLP . . . . .	409
14.5	Misc. Topics . . . . .	412
14.5.1	BLEU Score . . . . .	412
14.5.2	Connectionist Temporal Classification (CTC) . . . . .	413
14.5.3	Perplexity . . . . .	415
14.5.4	Byte Pair Encoding . . . . .	417
14.5.5	Grammars . . . . .	417
14.5.6	Bloom Filter . . . . .	418
14.5.7	Distributed Training . . . . .	418
14.5.8	Traditional Language Modeling . . . . .	419

# MATH AND MACHINE LEARNING BASICS

## CONTENTS

1.1	Linear Algebra (Quick Review) (Ch. 2)	8
1.1.1	Example: Principal Component Analysis	10
1.2	Probability & Information Theory (Quick Review) (Ch. 3)	12
1.3	Numerical Computation (Ch. 4)	14
1.4	Machine Learning Basics (Ch. 5)	17
1.4.1	Estimators, Bias and Variance (5.4)	17
1.4.2	Maximum Likelihood Estimation (5.5)	19
1.4.3	Bayesian Statistics (5.6)	20
1.4.4	Supervised Learning Algorithms (5.7)	22

## Linear Algebra (Quick Review) (Ch. 2)

Table of Contents Local

Written by Brandon McKinzie

- For  $A^{-1}$  to exist,  $Ax = b$  must have exactly one solution for every value of  $b$ . Determining whether a solution exists  $\forall b \in \mathbb{R}^m$  means requiring that the **column space** (range) of  $A$  be all of  $\mathbb{R}^m$ . It is helpful to see  $Ax$  expanded out explicitly in this way:

$$Ax = \sum_i x_i A_{:,i} = x_1 \begin{pmatrix} A_{1,1} \\ \vdots \\ A_{m,1} \end{pmatrix} + \cdots + x_m \begin{pmatrix} A_{1,n} \\ \vdots \\ A_{m,n} \end{pmatrix} \quad (2.27)$$

- Necessary:  $A$  must have at least  $m$  columns ( $n \geq m$ ). (“wide”).
- Necessary *and* sufficient: matrix must contain at least one set of  $m$  linearly independent columns.
- Invertibility: In addition to above, need matrix to be *square* (re: at most  $m$  columns  $\wedge$  at least  $m$  columns).
- A square matrix with linearly dependent columns is known as **singular**. A (necessarily square) matrix is singular if and only if one or more eigenvalues are zero.
- A **norm** is any function  $f$  that satisfies the following properties:  $\|x\|_\infty = \max_i |x_i|$

$$f(\mathbf{x}) = 0 \Rightarrow \mathbf{x} = \mathbf{0} \quad (1)$$

$$f(\mathbf{x} + \mathbf{y}) \leq f(\mathbf{x}) + f(\mathbf{y}) \quad (2)$$

$$\forall \alpha \in \mathbb{R}, f(\alpha \mathbf{x}) = |\alpha| f(\mathbf{x}) \quad (3)$$

- An **orthogonal matrix** is a square matrix whose rows are mutually orthonormal and whose columns are mutually orthonormal:

$$A^T A = A A^T = I \quad (2.37)$$

$$A^{-1} = A^T \quad (2.38)$$

- Suppose square matrix  $A \in \mathbb{R}^{n \times n}$  has  $n$  linearly independent eigenvectors  $\{\mathbf{v}^{(1)}, \dots, \mathbf{v}^{(n)}\}$ . The **eigendecomposition** of  $A$  is then given by<sup>1</sup>

$$A = V \text{diag}(\boldsymbol{\lambda}) V^{-1} \quad (2.40)$$

In the special case where  $A$  is real-symmetric,  $A = Q \Lambda Q^T$ . **Interpretation:**  $Ax$  can be decomposed into the following three steps:

<sup>1</sup>This appear to imply that unless the columns of  $V$  are also normalized, can't guarantee that its inverse equals its transpose? (since that is the only difference between it and an orthogonal matrix)

Unless stated otherwise,  
assume  $A \in \mathbb{R}^{m \times n}$

Note that orthonorm cols  
implies orthonorm rows  
(if square). To prove,  
consider the relationship  
between  $A^T A$  and  $A A^T$

All real-symmetric  $A$   
have an  
eigendecomposition, but  
it might not be  
unique!

- 1) **Change of basis:** The vector  $(\mathbf{Q}^T \mathbf{x})$  can be thought of as how  $\mathbf{x}$  would appear in the basis of eigenvectors of  $\mathbf{A}$ .
- 2) **Scale:** Next, we scale each component  $(\mathbf{Q}^T \mathbf{x})_i$  by an amount  $\lambda_i$ , yielding the new vector  $(\Lambda(\mathbf{Q}^T \mathbf{x}))$ .
- 3) **Change of basis:** Finally, we rotate this new vector back from the eigen-basis into its original basis, yielding the transformed result of  $\mathbf{Q}\Lambda\mathbf{Q}^T \mathbf{x}$ .

A common convention to sort the entries of  $\Lambda$  in descending order.

- **Positive definite:** all  $\lambda$  are positive; **positive semidefinite:** all  $\lambda$  are positive or zero.

$$\begin{aligned} \rightarrow \text{PSD: } & \forall \mathbf{x}, \quad \mathbf{x}^T \mathbf{A} \mathbf{x} \geq 0 \\ \rightarrow \text{PD: } & \mathbf{x}^T \mathbf{A} \mathbf{x} = 0 \Rightarrow \mathbf{x} = \mathbf{0}.^2 \end{aligned}$$

- Any real matrix  $\mathbf{A} \in \mathbb{R}^{m \times n}$  has a **singular value decomposition** of the form,

$$\mathbf{A} = \mathbf{U} \mathbf{D} \mathbf{V}^T \quad (10)$$

$$\mathbf{U} \in \mathbb{R}^{m \times m} \quad (7)$$

$$\mathbf{D} \in \mathbb{R}^{m \times n} \quad (8)$$

$$\mathbf{V} \in \mathbb{R}^{n \times n} \quad (9)$$

where both  $\mathbf{U}$  and  $\mathbf{V}$  are orthogonal matrices, and  $\mathbf{D}$  is diagonal.

- The **singular values** are the diagonal entries  $\mathbf{D}_{ii}$ .
- The **left(right)-singular vectors** are the columns of  $\mathbf{U}(\mathbf{V})$ .
- Eigenvectors of  $\mathbf{A}\mathbf{A}^T$  are the L-S vectors. Eigenvectors of  $\mathbf{A}^T\mathbf{A}$  are the R-S vectors.
- The eigenvalues of both  $\mathbf{A}\mathbf{A}^T$  and  $\mathbf{A}^T\mathbf{A}$  are given by the singular values squared.

- The Moore-Penrose **pseudoinverse**, denoted  $\mathbf{A}^+$ , enables us to find an “inverse” of sorts for a (possibly) non-square matrix  $\mathbf{A}$ . Most algorithms compute  $\mathbf{A}^+$  via

$$\mathbf{A}^+ = \mathbf{V} \mathbf{D}^+ \mathbf{U}^T \quad (11)$$

- The **determinant** of a matrix is  $\det(\mathbf{A}) = \prod_i \lambda_i$ . Conceptually,  $|\det(\mathbf{A})|$  tells how much [multiplication by]  $\mathbf{A}$  expands/contracts space. If  $\det(\mathbf{A}) = 1$ , the transformation preserves volume.

$\mathbf{A}^+$  is useful, e.g., when we want to solve  $\mathbf{A}\mathbf{x} = \mathbf{y}$  by left-multiplying each side to obtain  $\mathbf{x} = \mathbf{B}\mathbf{y}$ . It is far more likely for solution(s) to exist when  $\mathbf{A}$  is wider than it is tall.

---

<sup>2</sup>I proved this and it made me happy inside. Check it out. Let  $\mathbf{A}$  be positive definite. Then

$$\mathbf{x}^T \mathbf{A} \mathbf{x} = \mathbf{x}^T \mathbf{Q} \Lambda \mathbf{Q}^T \mathbf{x} \quad (4)$$

$$= \sum_i (\mathbf{Q}^T \mathbf{x})_i \lambda_i (\mathbf{Q}^T \mathbf{x})_i \quad (5)$$

$$= \sum_i \lambda_i (\mathbf{Q}^T \mathbf{x})_i^2 \quad (6)$$

Since all terms in the summation are non-negative and all  $\lambda_i > 0$ , we have that  $\mathbf{x}^T \mathbf{A} \mathbf{x} = 0$  if and only if  $(\mathbf{Q}^T \mathbf{x})_i = 0 = \mathbf{q}^{(i)} \cdot \mathbf{x}$  for all  $i$ . Since the set of eigenvectors  $\{\mathbf{q}^{(i)}\}$  form an orthonormal basis, we have that  $\mathbf{x}$  must be the zero vector.

---

### 1.1.1 EXAMPLE: PRINCIPAL COMPONENT ANALYSIS

---

**Task.** Say we want to apply lossy compression (less memory, but may lose precision) to a collection of  $m$  points  $\{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m)}\}$ . We will do this by converting each  $\mathbf{x}^{(i)} \in \mathbb{R}^n$  to some  $\mathbf{c}^{(i)} \in \mathbb{R}^l$  ( $l < n$ ), i.e. finding functions  $f$  and  $g$  such that:

$$f(\mathbf{x}) = \mathbf{c} \quad \text{and} \quad \mathbf{x} \approx g(f(\mathbf{x})) \quad (12)$$

**Decoding function ( $g$ ).** As is, we still have a rather general task to solve. *PCA* is defined by choosing  $g(\mathbf{c}) = \mathbf{D}\mathbf{c}$ , with  $\mathbf{D} \in \mathbb{R}^{n \times l}$ , where all columns of  $\mathbf{D}$  are both (1) orthogonal and (2) unit norm.

**Encoding function ( $f$ ).** Now we need a way of mapping  $\mathbf{x}$  to  $\mathbf{c}$  such that  $g(\mathbf{c})$  will give us back a vector optimally close to  $\mathbf{x}$ . We've already defined  $g$ , so this amounts to finding the optimal  $\mathbf{c}^*$  such that:

$$\mathbf{c}^* = \arg \min_{\mathbf{c}} \|\mathbf{x} - g(\mathbf{c})\|_2^2 \quad (13)$$

$$(\mathbf{x} - g(\mathbf{c}))^T (\mathbf{x} - g(\mathbf{c})) = \mathbf{x}^T \mathbf{x} - 2\mathbf{x}^T g(\mathbf{c}) + g(\mathbf{c})^T g(\mathbf{c}) \quad (14)$$

$$\mathbf{c}^* = \arg \min_{\mathbf{c}} [-2\mathbf{x}^T \mathbf{D}\mathbf{c} + \mathbf{c}^T \mathbf{c}] \quad (15)$$

$$= \mathbf{D}^T \mathbf{x} = f(\mathbf{x}) \quad (16)$$

which means the PCA *reconstruction operation* is defined as  $r(\mathbf{x}) = \mathbf{D}\mathbf{D}^T \mathbf{x}$ .

**Optimal  $\mathbf{D}$ .** It is important to notice that we've been able to determine e.g. the optimal  $\mathbf{c}^*$  for some  $\mathbf{x}$  because each  $\mathbf{x}$  has a (allowably) different  $\mathbf{c}^*$ . However, we use *the same* matrix  $\mathbf{D}$  for all our samples  $\mathbf{x}^{(i)}$ , and thus must optimize it over all points in our collection. With that out of the way, we just do what we always do: minimize over the  $L^2$  distance between points and their reconstruction. Formally, we minimize the Frobenius norm of the matrix of errors:

$$\mathbf{D}^* = \arg \min_{\mathbf{D}} \sqrt{\sum_{i,j} \left( \mathbf{x}_j^{(i)} - r(\mathbf{x}^{(i)})_j \right)^2} \quad \text{s.t.} \quad \mathbf{D}^T \mathbf{D} = \mathbf{I} \quad (17)$$

Consider the case of  $l = 1$  which means  $\mathbf{D} = \mathbf{d} \in \mathbb{R}^n$ . In this case, after [insert math here], we obtain

$$\mathbf{d}^* = \arg \max_{\mathbf{d}} \text{Tr} (\mathbf{d}^T \mathbf{X}^T \mathbf{X} \mathbf{d}) \quad s.t. \quad \mathbf{d}^T \mathbf{d} = 1 \quad (18)$$

where, as usual,  $\mathbf{X} \in \mathbb{R}^{m,n}$ . It should be clear that the optimal  $\mathbf{d}$  is just the largest eigenvector of  $\mathbf{X}^T \mathbf{X}$ .

## Probability &amp; Information Theory (Quick Review) (Ch. 3)

Table of Contents Local

Written by Brandon McKinzie

**Expectation.** For some function  $f(x)$ ,  $\mathbb{E}_{x \sim P}[f(x)]$  is the mean value that  $f$  takes on when  $x$  is drawn from  $P$ . The formula for discrete and continuous variables, respectively is as follows:

$$\mathbb{E}_{x \sim P}[f(x)] = \sum_x P(x)f(x) \quad (3.9)$$

$$\mathbb{E}_{x \sim P}[f(x)] = \int p(x)f(x)dx \quad (3.10)$$

**Variance.** A measure of how much the values of a function of a random variable  $x$  vary as we sample different values of  $x$  from its distribution.

$$\text{Var}[f(x)] = \mathbb{E}[(f(x) - \mathbb{E}[f(x)])^2] \quad (3.11)$$

**Covariance.** Gives some sense of how much two values are *linearly* related to each other, as well as the *scale* of these variables.

$$\text{Cov}[f(x), g(x)] = \mathbb{E}[(f(x) - \mathbb{E}[f(x)])(g(x) - \mathbb{E}[g(x)])] \quad (3.13)$$

- Large  $|\text{Cov}[f, g]|$  means the function values change a lot and both functions are far from their means at the same time.
- **Correlation** normalizes the contribution of each variable in order to measure only how much the variables are related.

**Covariance Matrix** of a random vector  $\mathbf{x} \in \mathbb{R}^n$  is an  $n \times n$  matrix, such that

$$\text{Cov}[\mathbf{x}]_{i,j} = \text{Cov}[x_i, x_j] \quad (3.14)$$

and if we want the “sample” covariance matrix taken over  $m$  data point samples, then

$$\Sigma := \frac{1}{m} \sum_{k=1}^m (x_k - \bar{x})(x_k - \bar{x})^T \quad (19)$$

where  $m$  is the number of data points.

## Measure Theory.

- A set of points that is negligibly small is said to have **measure zero**. In practical terms, think of such a set as occupying no volume in the space we are measuring (interested in).
- A property that holds **almost everywhere** holds throughout all space except for on a set of measure zero.

In  $\mathbb{R}^2$ , a line has measure zero.

## Functions of RVs.

- **Common mistake:** Suppose  $\mathbf{y} = g(\mathbf{x})$ , and  $g$  is invertible/continuous/differentiable. It is NOT true that  $p_y(\mathbf{y}) = p_x(g^{-1}(\mathbf{y}))$ . This fails to account for the distortion of [probability] space introduced by  $g$ . Rather,

$$p_x(\mathbf{x}) = p_y(g(\mathbf{x})) \left| \frac{\partial g(\mathbf{x})}{\partial \mathbf{x}} \right| \quad (3.47)$$

**Information Theory.** Denote the **self-information** of an event  $\mathbf{x} = x$  to be

$$I(x) \triangleq -\log P(x) \quad (20)$$

where  $\log$  is always assumed to be the natural logarithm. We can quantify the amount of uncertainty in an entire probability distribution using the **Shannon entropy**,

$$H(\mathbf{x}) = \mathbb{E}_{\mathbf{x} \sim P} [I(x)] = -\mathbb{E}_{\mathbf{x} \sim P} [\log P(x)] \quad (21)$$

which gives the expected amount of information in an event drawn from that distribution. Taking it a step further, say we have two separate probability distributions  $P(\mathbf{x})$  and  $Q(\mathbf{x})$ . We can measure how different these distributions are with the **Kullback-Leibler (KL) divergence**:

$$D_{KL}(P||Q) \triangleq \mathbb{E}_{\mathbf{x} \sim P} \left[ \log \frac{P(x)}{Q(x)} \right] = \mathbb{E}_{\mathbf{x} \sim P} [\log P(x) - \log Q(x)] \quad (22)$$

Note that the expectation is taken over  $P$ , thus making  $D_{KL}$  not symmetric (and thus not a true distance measure), since  $D_{KL}(P||Q) \neq D_{KL}(Q||P)$ . Finally, a closely related quantity is the **cross-entropy**,  $H(P, Q)$ , defined as:

$$H(P, Q) \triangleq H(P) + D_{KL}(P||Q) \quad (23)$$

$$= -\mathbb{E}_{\mathbf{x} \sim P} [\log Q(x)] \quad (24)$$

## Numerical Computation (Ch. 4)

Table of Contents Local

Written by Brandon McKinzie

**Some terminology.** **Underflow** is when numbers near zero are rounded to zero. Similarly, **overflow** is when large [magnitude] numbers are approximated as  $\pm\infty$ . **Conditioning** refers to how rapidly a function changes w.r.t. small changes in its inputs. Consider the function  $f(\mathbf{x}) = \mathbf{A}^{-1}\mathbf{x}$ . When  $\mathbf{A}$  has an eigenvalue decomposition, its *condition number* is

$$\max_{i,j} \left| \frac{\lambda_i}{\lambda_j} \right| \quad (4.2)$$

which is the ratio of the magnitude of the largest and smallest eigenvalue. When this is large, matrix inversion is sensitive to error in the input [of  $f(\mathbf{x})$ ].

**Gradient-based optimization.** Recall from basic calculus that the **directional derivative** of  $f(\mathbf{x})$  in direction  $\hat{\mathbf{u}}$  (a unit vector) is defined as the slope of the function  $f$  in direction  $\hat{\mathbf{u}}$ . By definition of the derivative, this is given by (with  $\mathbf{v} := \mathbf{x} + \alpha\hat{\mathbf{u}}$ )

$$\lim_{\alpha \rightarrow 0} \frac{f(\mathbf{x} + \alpha\hat{\mathbf{u}}) - f(\mathbf{x})}{\alpha} = \frac{\partial f(\mathbf{x} + \alpha\hat{\mathbf{u}})}{\partial \alpha} \Big|_{\alpha=0} \quad (25)$$

$$= \sum_i \frac{\partial f(\mathbf{v})}{\partial v_i} \frac{\partial v_i}{\partial \alpha} \Big|_{\alpha=0} \quad (26)$$

$$= \sum_i (\nabla_{\mathbf{v}} f(\mathbf{v}))_i u_i \Big|_{\alpha=0} \quad (27)$$

$$= \hat{\mathbf{u}}^T \nabla_{\mathbf{v}} f(\mathbf{v}) \Big|_{\alpha=0} \quad (28)$$

$$= \hat{\mathbf{u}}^T \nabla_{\mathbf{x}} f(\mathbf{x}) \quad (29)$$

where it's important to recognize the distinction between  $\lim_{\alpha \rightarrow 0}$  and *setting*  $\alpha$  to zero, which is denoted by  $|_{\alpha=0}$ . If we want to *find* the direction  $\hat{\mathbf{u}}$  such that this directional derivative is a minimum, i.e.

$$\hat{\mathbf{u}}^* = \arg \min_{\hat{\mathbf{u}}, \hat{\mathbf{u}}^T \hat{\mathbf{u}}=1} \hat{\mathbf{u}}^T \nabla_{\mathbf{x}} f(\mathbf{x}) \quad (30)$$

$$= \arg \min_{\hat{\mathbf{u}}, \hat{\mathbf{u}}^T \hat{\mathbf{u}}=1} \|\hat{\mathbf{u}}\|_2 \|\nabla_{\mathbf{x}} f(\mathbf{x})\|_2 \cos(\theta) \quad (31)$$

$$= \cos(\theta) \quad (32)$$

and we see that  $\hat{\mathbf{u}}$  points in the opposite direction as the gradient.

**Jacobian and Hessian Matrices.** For when we want partial derivatives of some function  $f$  whose input and output are both vectors. The **Jacobian matrix** contains all such partial derivatives. Sometimes we want to know about second derivatives too, since this tells us whether a gradient step will cause as much of an improvement as we would expect based on the gradient alone. The **Hessian matrix**  $\mathbf{H}(f)(\mathbf{x})$  is defined such that

$$\mathbf{H}(f)(\mathbf{x})_{i,j} = \frac{\partial^2}{\partial x_i \partial x_j} f(\mathbf{x}) \quad (4.6)$$

$$f : \mathbb{R}^m \rightarrow \mathbb{R}^n$$

$$\mathbf{J} \in \mathbb{R}^{n \times m} \text{ where}$$

$$J_{i,j} = \frac{\partial}{\partial x_j} f(\mathbf{x})_i$$

The Hessian is the Jacobian of the gradient.

The second derivative in a specific direction  $\hat{\mathbf{d}}$  is given by  $\hat{\mathbf{d}}^T \mathbf{H} \hat{\mathbf{d}}$ <sup>3</sup>. It tells us how well we can expect a gradient descent step to perform. How so? Well, it shows up in the second-order approximation to the function  $f(\mathbf{x})$  about our current spot, which we can denote  $\mathbf{x}^{(0)}$ . The standard gradient descent step will move us from  $\mathbf{x}^{(0)} \rightarrow \mathbf{x}^{(0)} - \epsilon \mathbf{g}$ , where  $\mathbf{g}$  is the gradient evaluated at  $\mathbf{x}^{(0)}$ . Plugging this in to the 2nd order approximation shows us how  $\mathbf{H}$  can give information related to how “good” of a step that really was. Mathematically,

$$f(\mathbf{x}) \approx f(\mathbf{x}^{(0)}) + (\mathbf{x} - \mathbf{x}^{(0)})^T \mathbf{g} + \frac{1}{2} (\mathbf{x} - \mathbf{x}^{(0)})^T \mathbf{H} (\mathbf{x} - \mathbf{x}^{(0)}) \quad (4.8)$$

$$f(\mathbf{x}^{(0)} - \epsilon \mathbf{g}) \approx f(\mathbf{x}^{(0)}) - \epsilon \mathbf{g}^T \mathbf{g} + \frac{1}{2} \epsilon^2 \mathbf{g}^T \mathbf{H} \mathbf{g} \quad (4.9)$$

If  $\mathbf{g}^T \mathbf{H} \mathbf{g}$  is positive, then we can easily solve for the optimal  $\epsilon = \epsilon^*$  that decreases the Taylor series approximation as

$$\epsilon^* = \frac{\mathbf{g}^T \mathbf{g}}{\mathbf{g}^T \mathbf{H} \mathbf{g}} \quad (4.10)$$

which can be as low as  $1/\lambda_{max}$  (the worst case), and as high as  $1/\lambda_{min}$  with the  $\lambda$  being the eigenvalues of the Hessian. The best (and perhaps only) way to take what we learned about the “second derivative test” in single-variable calculus and apply it to the multidimensional case with  $\mathbf{H}$  is by using the *eigendecomposition of  $\mathbf{H}$* . Why? Because we can examine the eigenvalues of the Hessian to determine whether the critical point  $\mathbf{x}^{(0)}$  is a local maximum, local minimum, or saddle point<sup>4</sup>. If all eigenvalues are positive (remember that this is equivalent to saying that the Hessian is **positive definite!**), the point is a local minimum.

The condition number of the Hessian at a given point can give us an idea about how much the second derivatives (along different directions) differ from each other

---

<sup>3</sup> In the same manner that I derived equation 29, we can derive the second derivative in a specified direction  $\hat{\mathbf{d}}$ :

$$\frac{\partial^2}{\partial \alpha^2} f(\mathbf{x} + \alpha \hat{\mathbf{d}}) \Big|_{\alpha=0} = \frac{\partial}{\partial \alpha} \hat{\mathbf{d}}^T \nabla_{\mathbf{v}} f(\mathbf{v}) \Big|_{\alpha=0} \quad (33)$$

$$= \sum_i d_i \frac{\partial}{\partial \alpha} \frac{\partial f(\mathbf{v})}{\partial v_i} \Big|_{\alpha=0} \quad (34)$$

$$= \sum_i d_i \frac{\partial}{\partial v_i} \frac{\partial f(\mathbf{v})}{\partial \alpha} \Big|_{\alpha=0} \quad (35)$$

$$= \sum_i \sum_j d_i \frac{\partial^2 f(\mathbf{v})}{\partial v_i v_j} d_j \Big|_{\alpha=0} \quad (36)$$

$$= \hat{\mathbf{d}}^T \mathbf{H} \hat{\mathbf{d}} \quad (37)$$

<sup>4</sup>Emphasis on “values” in “eigenvalues” because it’s important not to get tripped up here about what the

**Constrained optimization:** minimizing/maximizing a function  $f(\mathbf{x})$  constrained to only values of  $\mathbf{x}$  in some set  $\mathbb{S}$ . One way of approaching such a problem is to re-design the unconstrained optimization problem such that the re-designed problem's solution satisfies the constraints. For example, to minimize  $f(\mathbf{x})$  for  $\mathbf{x} \in \mathbb{R}^2$  with constraint  $\|\mathbf{x}\|_2 = 1$ , we can minimize  $g(\theta) = f([\cos \theta, \sin \theta]^T)$  wrt  $\theta$ , then return  $[\cos \theta, \sin \theta]^T$  as the solution to the original problem.

The **Karush-Kuhn-Tucker** (KKT) approach, a generalization of **Lagrange multipliers**, provides a general approach for re-designing the optimization problem, with procedure as follows:

1. Find  $m$  functions  $g^{(i)}$  and  $n$  functions  $h^{(j)}$  such that your set of allowed values  $\mathbb{S}$  can be written

$$\mathbb{S} = \{\mathbf{x} \mid \forall i, g^{(i)}(\mathbf{x}) = 0 \text{ and } \forall j, h^{(j)}(\mathbf{x}) \leq 0\} \quad (38)$$

The equations involving  $g^{(i)}$  are called the **equality constraints** and the inequalities involving  $h^{(j)}$  are called the **inequality constraints**.

2. Introduce new variables  $\lambda_i$  (for the equality constraints) and  $\alpha_j$  (for the inequality constraints). These are called the KKT multipliers. The generalized Lagrangian is then defined as

$$\mathcal{L}(\mathbf{x}, \boldsymbol{\lambda}, \boldsymbol{\alpha}) = f(\mathbf{x}) + \sum_i \lambda_i g^{(i)}(\mathbf{x}) + \sum_j \alpha_j h^{(j)}(\mathbf{x}) \quad (39)$$

3. Solve the re-designed unconstrained optimization problem:

$$\min_{\mathbf{x}} \max_{\boldsymbol{\lambda}} \max_{\boldsymbol{\alpha}, \boldsymbol{\alpha} \geq 0} \mathcal{L}(\mathbf{x}, \boldsymbol{\lambda}, \boldsymbol{\alpha}) \quad (40)$$

which has the same optimal objective function value and set of optimal points  $\mathbf{x}$  as the original constrained problem,  $\min_{\mathbf{x} \in \mathbb{S}} f(\mathbf{x})$ . Any time the constraints are satisfied, the expression  $\max_{\boldsymbol{\lambda}} \max_{\boldsymbol{\alpha}, \boldsymbol{\alpha} \geq 0} \mathcal{L}(\mathbf{x}, \boldsymbol{\lambda}, \boldsymbol{\alpha})$  evaluates to  $f(\mathbf{x})$ , and any time a constraint is violated, the same expression evaluates to  $\infty$ .

---

eigenvectors of the Hessian mean. The reason for the decomposition is that it gives us an orthonormal basis (out of which we can get any direction) and therefore the magnitude of the second derivative along each of these directions as the eigenvalues.

## Machine Learning Basics (Ch. 5)

Table of Contents Local

Written by Brandon McKinzie

**Capacity, Overfitting, and Underfitting.** Difference between ML and optimization is that, in addition to wanting low training error, we want **generalization error** (test error) to be low as well. The ideal model is an oracle that simply knows the true probability distribution  $p(\mathbf{x}, y)$  that generates the data. The error incurred by such an oracle, due things like inherently stochastic mappings from  $\mathbf{x}$  to  $y$  or other variables, is called the **Bayes error**. The **no free lunch theorem** states that, averaged over all possible data-generating distributions, every classification algorithm has the same error rate when classifying previously unobserved points. Therefore, the goal of ML research is to understand what kinds of distributions are relevant to the “real world” that an AI agent experiences, and what kinds of ML algorithms perform well on data drawn from the relevant data-generating distributions.

## 1.4.1 ESTIMATORS, BIAS AND VARIANCE (5.4)

**Point Estimation:** attempt to provide “best” prediction of some quantity, such as some parameter or even a whole function. Formally, a point estimator or *statistic* is any function of the data:

$$\hat{\theta}_m = g(x^{(1)}, \dots, x^{(m)}) \quad (5.19)$$

where, since the data is drawn from a random process,  $\hat{\theta}$  is a random variable. **Function estimation** is identical in form, where we want to estimate some  $f(x)$  with  $\hat{f}$ , a point estimator in *function space*.

**Bias.** Defined below, where the expectation is taken over the data-generating distribution<sup>5</sup>. Bias measures the expected deviation from the true value of the func/param.

$$\text{bias}[\hat{\theta}_m] = \mathbb{E}[\hat{\theta}_m] - \theta \quad (5.20)$$

**TODO:** Figure out how to derive  $\mathbb{E}[\hat{\theta}_m^2]$  for Gaussian distribution [helpful link].

---

<sup>5</sup>May want to double-check this, but I’m fairly certain this is what the book meant when it said “data,” based on later examples.

## Bias-Variance Tradeoff.

- **Conceptual Info.** Two sources of error for an estimator are (1) bias and (2) variance, which are both defined as deviations from a certain value. Bias gives deviation from the *true* value, while variance gives the [expected] deviation from this *expected* value.
- **Summary of main formulas.**

$$\text{bias} [\hat{\theta}_m] = \mathbb{E} [\hat{\theta}_m] - \theta \quad (41)$$

$$\text{Var} [\hat{\theta}_m] = \mathbb{E} \left[ (\hat{\theta}_m - \mathbb{E} [\hat{\theta}_m])^2 \right] \quad (42)$$

- **MSE decomposition.** The MSE of the estimates is given by<sup>6</sup>

$$\text{MSE} = \mathbb{E} [(\hat{\theta}_m - \theta)^2] \quad (5.53)$$

$$= \text{Bias}(\hat{\theta})^2 + \text{Var} [\hat{\theta}_m] \quad (5.54)$$

and desirable estimators are those with low MSE.

**Consistency.** As the number of training data points increases, we want the estimators to converge to the true values. Specifically, below are the definitions for *weak* and *strong* consistency, respectively.

$$\begin{aligned} \text{plim}_{m \rightarrow \infty} \hat{\theta}_m &= 0 \\ p \left( \lim_{m \rightarrow \infty} \hat{\theta}_m = \theta \right) &= 1 \end{aligned} \quad (5.55)$$

where the symbol plim means  $P(|\hat{\theta}_m - \theta| > \epsilon) \rightarrow 0$  as  $m \rightarrow \infty$ .

---

<sup>6</sup>Derivation:

$$\text{MSE} = \mathbb{E} [\hat{\theta}^2 + \theta^2 - 2\theta\hat{\theta}] \quad (43)$$

$$= \mathbb{E} [\hat{\theta}^2] + \theta^2 - 2\theta\mathbb{E} [\hat{\theta}] \quad (44)$$

$$= (\mathbb{E} [\hat{\theta}]^2 - \mathbb{E} [\hat{\theta}]^2) + \mathbb{E} [\hat{\theta}^2] + \theta^2 - 2\theta\mathbb{E} [\hat{\theta}] \quad (45)$$

$$= (\mathbb{E} [\hat{\theta}]^2 + \theta^2 - 2\theta\mathbb{E} [\hat{\theta}]) + (\mathbb{E} [\hat{\theta}^2] - \mathbb{E} [\hat{\theta}]^2) \quad (46)$$

$$= \text{Bias}(\hat{\theta})^2 + \text{Var} [\hat{\theta}_m] \quad (47)$$

---

#### 1.4.2 MAXIMUM LIKELIHOOD ESTIMATION (5.5)

Consider set of  $m$  examples  $\mathbb{X} = \{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m)}\}$  drawn independently from the true (but unknown)  $p_{data}(\mathbf{x})$ . Let  $p_{model}(\mathbf{x}; \boldsymbol{\theta})$  be parametric family of probability distributions over the same space indexed by  $\boldsymbol{\theta}$ . The maximum likelihood estimator for  $\boldsymbol{\theta}$  can be expressed as

$$\boldsymbol{\theta}_{ML} = \arg \max_{\boldsymbol{\theta}} \mathbb{E}_{\mathbf{x} \sim \hat{p}_{data}} [\log p_{model}(\mathbf{x}; \boldsymbol{\theta})] \quad (5.59)$$

where we've chosen to express with log for underflow/gradient reasons. One interpretation of ML is to view it as minimizing the dissimilarity, as measured by the KL divergence<sup>7</sup>, between  $\hat{p}_{data}$  and  $p_{model}$ .

Any loss consisting of a negative log-likelihood is a **cross-entropy** between the  $\hat{p}_{data}$  distribution and the  $p_{model}$  distribution.

Thoughts: Let's look at  $D_{KL}$  in some more detail. First, I'll rewrite it with the explicit definition of  $\mathbb{E}_{\mathbf{x} \sim \hat{p}_{data}} [\log (\hat{p}_{data}(\mathbf{x}))]$ :

$$D_{KL}(\hat{p}_{data} || p_{model}) = \mathbb{E}_{\mathbf{x} \sim \hat{p}_{data}} [\log (\hat{p}_{data}(\mathbf{x})) - \log (p_{model}(\mathbf{x}))] \quad (48)$$

$$= \left( \frac{1}{N} \left( \sum_{i=1}^N \log (\text{Counts}(\mathbf{x}_i)) \right) - \log N \right) - \mathbb{E}_{\mathbf{x} \sim \hat{p}_{data}} [\log (p_{model}(\mathbf{x}))] \quad (49)$$

Note also that our goal is to find parameters  $\boldsymbol{\theta}$  such that  $D_{KL}$  is minimized. It is for *this* reason, that we wish to optimize over  $\boldsymbol{\theta}$ , that minimizing  $D_{KL}$  amounts to maximizing the quantity,  $\mathbb{E}_{\mathbf{x} \sim \hat{p}_{data}} [\log (p_{model}(\mathbf{x}))]$ . Sure, I can agree this is *true*, but **why is our goal to minimize  $D_{KL}$ , as opposed to minimizing  $|D_{KL}|$ ?** I'm assuming it is because optimizing w.r.t. an absolute value is challenging numerically.

**Conditional Log-Likelihood and MSE.** We can readily generalize  $\boldsymbol{\theta}_{ML}$  to estimate a conditional probability  $p(\mathbf{y} | \mathbf{x}; \boldsymbol{\theta})$  in order to predict  $\mathbf{y}$  given  $\mathbf{x}$ , since

We are assuming the examples are i.i.d. here.

$$\boldsymbol{\theta}_{ML} = \arg \max_{\boldsymbol{\theta}} \sum_{i=1}^m \log P(\mathbf{y}^{(i)} | \mathbf{x}^{(i)}; \boldsymbol{\theta}) \quad (5.63)$$

where  $\mathbf{x}^{(i)}$  are fed as *inputs* to the model; this is why we can formulate MLE as a conditional probability.

---

<sup>7</sup> The KL divergence is given by

$$D_{KL}(\hat{p}_{data} || p_{model}) = \mathbb{E}_{\mathbf{x} \sim \hat{p}_{data}} [\log \hat{p}_{data}(\mathbf{x}) - \log p_{model}(\mathbf{x})] \quad (5.60)$$

---

### 1.4.3 BAYESIAN STATISTICS (5.6)

Distinction between frequentist and bayesian approach:

- **Frequentist:** Estimate  $\boldsymbol{\theta}$  —> make predictions thereafter based on this estimate.
- **Bayesian:** Consider all possible values of  $\boldsymbol{\theta}$  when making predictions.

**The prior.** Before observing the data, we represent our knowledge of  $\boldsymbol{\theta}$  using the **prior probability distribution**  $p(\boldsymbol{\theta})$ . Unlike maximum likelihood, which makes predictions using a *point estimate* of  $\boldsymbol{\theta}$  (a single value), the Bayesian approach uses Bayes' rule to make predictions using the *full distribution* over  $\boldsymbol{\theta}$ . In other words, rather than focusing on the most accurate value estimate of  $\boldsymbol{\theta}$ , we instead focus on pinning down a range of possible  $\boldsymbol{\theta}$  values and how likely we believe each of these values to be.

It is common to choose a high-entropy prior, e.g. uniform.

So what happens to  $\boldsymbol{\theta}$  after we observe the data? We update it using Bayes' rule<sup>8</sup>:

$$p(\boldsymbol{\theta} \mid x^{(1)}, \dots, x^{(m)}) = \frac{p(x^{(1)}, \dots, x^{(m)} \mid \boldsymbol{\theta})p(\boldsymbol{\theta})}{p(x^{(1)}, \dots, x^{(m)})} \quad (50)$$

Note that we still haven't mentioned how to actually make *predictions*. Since we no longer have just one value for  $\boldsymbol{\theta}$ , but rather we have a posterior *distribution*  $p(\boldsymbol{\theta} \mid x^{(1)}, \dots, x^{(m)})$ , we must integrate over this to get the predicted likelihood of the next sample  $x^{(m+1)}$ :

$$p(x^{(m+1)} \mid x^{(1)}, \dots, x^{(m)}) = \int p(x^{(m+1)} \mid \boldsymbol{\theta})p(\boldsymbol{\theta} \mid x^{(1)}, \dots, x^{(m)})d\boldsymbol{\theta} \quad (51)$$

$$= \mathbb{E}_{\boldsymbol{\theta} \sim p(\boldsymbol{\theta} \mid x^{(1)}, \dots, x^{(m)})} [p(x^{(m+1)} \mid \boldsymbol{\theta})] \quad (52)$$

**Linear Regression: MLE vs. Bayesian.** Both want to model the conditional distribution  $p(y \mid \mathbf{x})$  (the conditional likelihood). To derive the standard linear regression algorithm, we *define*

$$p(y \mid \mathbf{x}) = \mathcal{N}(y; \hat{y}(\mathbf{x}; \mathbf{w}), \sigma^2) \quad (53)$$

$$\hat{y}(\mathbf{x}; \mathbf{w}) = \mathbf{w}^T \mathbf{x} \quad (54)$$

Assume  $\sigma^2$  is some fixed constant chosen by the user.

---

<sup>8</sup>In practice, we typically compute the denominator by simply normalizing the probability distribution, i.e. it is effectively the partition function.

- **Maximum Likelihood Approach:** We can use the definition above (and the i.i.d. assumption) to evaluate the conditional log-likelihood as

$$\sum_{i=1}^m \log p(y^{(i)} | \mathbf{x}^{(i)}; \boldsymbol{\theta}) = -m \log \sigma - \frac{m}{2} \log(2\pi) - \sum_{i=1}^m \frac{\|\hat{y}^{(i)} - y^{(i)}\|^2}{2\sigma^2} \quad (5.65)$$

where only the last term has any dependence on  $\mathbf{w}$ . Therefore, to obtain  $\mathbf{w}_{ML}$  we take the derivative of the last term w.r.t.  $\mathbf{w}$ , set that to zero, and solve for  $\mathbf{w}$ . We see that finding the  $\mathbf{w}$  that maximizes the conditional log-likelihood is equivalent to finding the  $\mathbf{w}$  that minimizes the training MSE.

Recall that the training MSE is  $\frac{1}{m} \sum_{i=1}^m \|\hat{y}^{(i)} - y^{(i)}\|^2$

- **Bayesian Approach:** Our conditional likelihood is already given in equation 53. Next, we must define a prior distribution over  $\mathbf{w}$ . As is common, we choose a Gaussian prior to express our high degree of uncertainty about  $\boldsymbol{\theta}$  (implying we'll choose a relatively large variance):

$$p(\mathbf{w}) := \mathcal{N}(\mathbf{w}; \boldsymbol{\mu}_0, \boldsymbol{\Lambda}_0) \quad (55)$$

Typically assume  $\boldsymbol{\Lambda}_0 = \text{diag}(\boldsymbol{\lambda}_0)$

We can then compute [the unnormalized]  $p(\mathbf{w} | \mathbf{X}, \mathbf{y}) \propto p(\mathbf{y} | \mathbf{X}, \mathbf{w})p(\mathbf{w})$  [and then normalize it].

**Maximum A Posteriori (MAP) Estimation.** Often we either prefer a point estimate for  $\boldsymbol{\theta}$ , or we find out that computing the posterior distribution is intractable and a point estimate offers a tractable estimation. The obvious way of obtaining this while still taking the Bayesian route is to just argmax the posterior and use that as your point estimate:

$$\boldsymbol{\theta}_{MAP} = \arg \max_{\boldsymbol{\theta}} p(\boldsymbol{\theta} | \mathbf{x}) = \arg \max_{\boldsymbol{\theta}} \log p(\mathbf{x} | \boldsymbol{\theta}) + \log p(\boldsymbol{\theta}) \quad (56)$$

where the second form shows how this is basically maximum likelihood with incorporation of the prior. We don't want just any  $\boldsymbol{\theta}$  that maximizes the likelihood of our data if there is virtually no chance of that value of  $\boldsymbol{\theta}$  in the first place.

---

#### 1.4.4 SUPERVISED LEARNING ALGORITHMS (5.7)

---

**Logistic Regression.** We've already seen that linear regression corresponds to the family

$$p(y | \mathbf{x}) = \mathcal{N}(y; \boldsymbol{\theta}^T \mathbf{x}, \mathbf{I}) \quad (5.80)$$

which we can generalize to the binary **classification** scenario by interpreting as the probability of class 1. One way of doing this while ensuring the output is between 0 and 1 is to use the logistic sigmoid function:

$$p(y = 1 \mid \mathbf{x}; \boldsymbol{\theta}) = \sigma(\boldsymbol{\theta}^T \mathbf{x}) \quad (5.81)$$

Equation 5.81 is the definition of logistic regression

Unfortunately, there is no closed-form solution for  $\boldsymbol{\theta}$ , so we must search via maximizing the log-likelihood.

**Support Vector Machines.** Driving by a linear function  $\mathbf{w}^T \mathbf{x} + b$  like logistic regression, but instead of outputting probabilities it outputs a class identity, which depends on the sign of  $\mathbf{w}^T \mathbf{x} + b$ . SVMs make use of the **kernel trick**, the “trick” being that we can rewrite  $\mathbf{w}^T \mathbf{x} + b$  completely in terms of dot products between examples. The general form of our prediction function becomes

$$f(\mathbf{x}) = b + \sum_i \alpha_i k(\mathbf{x}, \mathbf{x}^{(i)}) \quad (5.83)$$

If our kernel function is just  $k(x, x^{(i)}) = x^T x^{(i)}$  then we've just rewritten  $\mathbf{w}$  in the form  $\mathbf{w} \rightarrow \mathbf{X}^T \boldsymbol{\alpha}$

where the *kernel* [function] takes the general form  $k(\mathbf{x}, \mathbf{x}^{(i)}) = \phi(\mathbf{x}) \cdot \phi(\mathbf{x}^{(i)})$ . A major drawback to kernel machines (methods) in general is that the cost of evaluating the decision function  $f(\mathbf{x})$  is linear in the number of training examples. SVMs, however, are able to mitigate this by learning an  $\alpha$  with mostly zeros. The training examples with *nonzero*  $\alpha_i$  are known as **support vectors**.

# DEEP NETWORKS: MODERN PRACTICES

## CONTENTS

2.1	Deep Feedforward Networks (Ch. 6) . . . . .	24
2.1.1	Back-Propagation (6.5) . . . . .	25
2.2	Regularization for Deep Learning (Ch. 7) . . . . .	26
2.3	Optimization for Training Deep Models (Ch. 8) . . . . .	28
2.4	Convolutional Neural Networks (Ch. 9) . . . . .	32
2.5	Sequence Modeling (RNNs) (Ch. 10) . . . . .	35
2.5.1	Review: The Basics of RNNs . . . . .	35
2.5.2	RNNs as Directed Graphical Models . . . . .	40
2.5.3	Challenge of Long-Term Deps. (10.7) . . . . .	42
2.5.4	LSTMs and Other Gated RNNs (10.10) . . . . .	43
2.6	Applications (Ch. 12) . . . . .	44
2.6.1	Natural Language Processing (12.4) . . . . .	44
2.6.2	Neural Language Models (12.4.2) . . . . .	45

## Deep Feedforward Networks (Ch. 6)

Table of Contents Local

Written by Brandon McKinzie

The strategy/purpose of [feedforward] deep learning is to *learn the set of features/representation describing  $\mathbf{x}$*  with a mapping  $\phi$  before applying a linear model. In this approach, we have a model

$$y = f(\mathbf{x}; \boldsymbol{\theta}, \mathbf{w}) = \phi(\mathbf{x}; \boldsymbol{\theta})^T \mathbf{w}$$

with  $\phi$  defining a hidden layer.

**ReLUs and their generalizations.** Some nice properties of ReLUs are...

- Derivatives through a ReLU remain large and consistent whenever the unit is active.
- Second derivative is 0 a.e. and the derivative is 1 everywhere the unit is active, meaning the gradient direction is more useful for learning than it would be with activation functions that introduce 2nd-order effects (see equation 4.9)

Recall the ReLU activation function:  
 $g(z) = \max\{0, z\}$  almost everywhere”

**Generalizing to aid gradients when  $z < 0$ .** Three such generalizations are based on using a nonzero slope  $\alpha_i$  when  $z_i < 0$ :

$$h_i = g(\mathbf{z}, \boldsymbol{\alpha})_i = \max(0, z_i) + \alpha_i \min(0, z_i) \quad (57)$$

- Absolute value rectification: fix  $\alpha_i = -1$  to obtain  $g(z) = |z|$ .
- Leaky ReLU: fix  $\alpha_i$  to a small value like 0.01.
- Parametric ReLU (PReLU): treats  $\alpha_i$  like a learnable parameter.

**Logistic sigmoid and hyperbolic tangent.** Sigmoid activations on hidden units is a bad idea, since they’re only sensitive to their inputs near zero, with small gradients everywhere else. If sigmoid activations must be used, tanh is probably a better substitute, since it resembles the identity (i.e. a linear function) near zero.

---

### 2.1.1 BACK-PROPAGATION (6.5)

---

**The chain rule.** Suppose  $z = f(\mathbf{y})$  where  $\mathbf{y} = g(\mathbf{x})$  (see margin for dimensions). Then<sup>9</sup>,

$$\frac{\partial z}{\partial x_i} = (\nabla_{\mathbf{x}} z)_i = \sum_{j=1}^n \frac{\partial z}{\partial y_j} \frac{\partial y_j}{\partial x_i} = \sum_{j=1}^n (\nabla_{\mathbf{y}} z)_j \frac{\partial y_j}{\partial x_i} = \sum_{j=1}^n (\nabla_{\mathbf{y}} z)_j (\nabla_{\mathbf{x}} y_j)_i \quad (6.45)$$

$$\rightarrow \nabla_{\mathbf{x}} z = \left( \frac{\partial \mathbf{y}}{\partial \mathbf{x}} \right)^T \nabla_{\mathbf{y}} z = \mathbf{J}_{y=g(\mathbf{x})}^T \nabla_{\mathbf{y}} z \quad (6.46)$$

$$\begin{aligned} \mathbf{x} &\in \mathbb{R}^m \\ \mathbf{y} &\in \mathbb{R}^n \\ z &: \mathbb{R}^n \rightarrow \mathbb{R} \\ g &: \mathbb{R}^m \rightarrow \mathbb{R}^n \end{aligned}$$

From this we see that the gradient of a variable  $x$  can be obtained by multiplying a Jacobian matrix  $\frac{\partial \mathbf{y}}{\partial \mathbf{x}}$  by a gradient  $\nabla_{\mathbf{y}} z$ .

---

<sup>9</sup>Note that we can view  $z = f(\mathbf{y})$  as a multi-variable function of the dimensions of  $\mathbf{y}$ ,

$$z = f(y_1, y_2, \dots, y_n)$$

## Regularization for Deep Learning (Ch. 7)

Table of Contents Local

Written by Brandon McKinzie

Recall the definition of regularization: “any modification we make to a learning algorithm that is intended to reduce its generalization error but not its training error.”

**Limiting Model Capacity.** Recall that **Capacity** [of a model] is the ability to fit a wide variety of functions. Low cap models may struggle to fit training set, while high cap models may overfit by simply memorizing the training set. We can limit model capacity by adding a parameter norm penalty  $\Omega(\theta)$  to the objective function  $J$ :

$$\tilde{J}(\theta; X, y) = J(\theta; X, y) + \alpha\Omega(\theta) \quad \text{where } \alpha \in [0, \infty) \quad (7.1)$$

where we typically choose  $\Omega$  to only penalize the *weights* and leave biases unregularized.

**L2-Regularization.** Defined as setting  $\Omega(\theta) = \frac{1}{2}\|w\|_2^2$ . Assume that  $J(w)$  is quadratic, with minimum at  $w^*$ . Since quadratic, we can approximate  $J$  with a second-order expansion about  $w^*$ .

$$\hat{J}(w) = J(w^*) + \frac{1}{2}(w - w^*)^T H(w - w^*) \quad (7.6)$$

$$\nabla_w \hat{J}(w) = H(w - w^*) \quad (7.7)$$

where  $H_{ij} = \frac{\partial^2 J}{\partial w_i \partial w_j}|_{w^*}$ . If we add in the [derivative of] the weight decay and set to zero, we obtain the solution

$$\tilde{w} = (H + \alpha I)^{-1} H w^* \quad (7.10)$$

$$= Q(\Lambda + \alpha I)^{-1} \Lambda Q^T w^* \quad (7.13)$$

which shows that the effect of regularization is to rescale the  $i$  eigenvectors of  $H$  by  $\frac{\lambda_i}{\lambda_i + \alpha}$ . This means that eigenvectors with  $\lambda_i >> \alpha$  are relatively unchanged, but the eigenvectors with  $\lambda_i << \alpha$  are shrunk to nearly zero. In other words, only directions along which the parameters contribute significantly to reducing the objective function are preserved relatively intact.

**Sparse Representations.** Weight decay acts by placing a penalty directly on the model parameters. Another strategy is to place a penalty on the *activations* of the units, encouraging their activations to be sparse. It's important to distinguish the difference between sparse parameters and sparse *representations*. In the former, if we take the example of some  $\mathbf{y} = \mathbf{B}\mathbf{h}$ , there are many zero entries in some parameter matrix  $\mathbf{B}$  while, in the latter, there are many zero entries in the representation vector  $\mathbf{h}$ . The modification to the loss function, analogous to 7.1, takes the form

$$\tilde{J}(\boldsymbol{\theta}; \mathbf{X}, \mathbf{y}) = J(\boldsymbol{\theta}; \mathbf{X}, \mathbf{y}) + \alpha\Omega(\mathbf{h}) \quad \text{where } \alpha \in [0, \infty) \quad (7.48)$$

**Adversarial Training** Even for networks that perform at human level accuracy have a nearly 100 percent error rate on examples that are intentionally constructed to search for an input  $\mathbf{x}'$  near a data point  $\mathbf{x}$  such that the model output for  $\mathbf{x}'$  is very different than the output for  $\mathbf{x}$ .

$$\mathbf{x}' \leftarrow \mathbf{x} + \epsilon \cdot \text{sign}(\nabla_{\mathbf{x}} J(\boldsymbol{\theta}; \mathbf{x}, \mathbf{y})) \quad (58)$$

In many cases,  $\mathbf{x}'$  can be so similar to  $\mathbf{x}$  that a human cannot tell the difference!

In the context of regularization, one can reduce the error rate on the original i.i.d. test set via **adversarial training** – training on adversarially perturbed training examples.

## Optimization for Training Deep Models (Ch. 8)

Table of Contents Local

Written by Brandon McKinzie

**Empirical Risk Minimization.** The ultimate goal of any machine learning algorithm is to reduce the expected generalization error, also called the **risk**:

$$J^*(\boldsymbol{\theta}) = \mathbb{E}_{(\mathbf{x},y) \sim p_{data}} [L(f(\mathbf{x}; \boldsymbol{\theta}), y)] \quad (59) \quad \text{risk definitions}$$

with emphasis that the risk is over the *true* underlying data distribution  $p_{data}$ . If we knew  $p_{data}$ , this would be an optimization problem. Since we don't, and only have a set of training samples, it is a machine learning problem. However, we can still just minimize the **empirical risk**, replacing  $p_{data}$  in the equation above with  $\hat{p}_{data}$ <sup>10</sup>.

So, how is minimizing the empirical risk any different than familiar gradient descent approaches? Aren't they designed to do just that? Well, sort of, but it's technically not the same. When we say "minimize the empirical risk" in the context of optimization, we mean this very literally. Gradient descent methods emphatically do *not* just go and *set* the weights to values such that the empirical risk reaches its lowest possible value – that's not machine learning. Furthermore, many useful loss function such as 0-1 loss<sup>11</sup> do not have useful derivatives.

ERM  $\neq$  GD

**Surrogate Loss Functions and Early Stopping.** In cases such as 0-1 loss, where minimization is intractable, one typically optimizes a **surrogate loss function** instead, such as the negative log-likelihood of the correct class. Also, an important difference between pure optimization and our training algorithms is that the latter usually don't halt at a local minimum. Instead, we usually must define some early stopping condition to terminate training before overfitting begins to occur.

*We want to minimize the risk, but we don't have access to  $p_{data}$ , so ...*

*We want to minimize the empirical risk, but it's prone to overfitting and our loss function's derivative may be zero/undefined, so ...*

*We minimize a surrogate loss function iteratively over minibatches until early stopping is triggered.*

---

<sup>10</sup>This amounts to a simple average over the loss function at each training point.

<sup>11</sup>The 0-1 loss function is defined as

$$L(\hat{y}, y) = I(\hat{y} \neq y) \quad (60)$$

**Batch and Minibatch Algorithms.** Computing  $\nabla_{\boldsymbol{\theta}} J(\boldsymbol{\theta})$  as an expectation over the entire training set is expensive, so we typically compute the expectation over a small subset of the examples. Recall that the standard deviation, or standard error  $SE(\mu_m)$ , of the mean taken over some subset of  $m \leq n$  samples,  $\mu_m = \frac{1}{m} \sum_{i \sim \text{Rand}(0, n, \text{size}=m)} x^{(i)}$ , is given by  $\sigma/\sqrt{m}$ , where  $\sigma$  is the true [sample] standard deviation of the full  $n$  data samples. In other words, to improve such a gradient by a factor of 10 requires 100 times more samples-per-batch (and thus 100 times more computation). For this reason, most optimization algorithms actually *converge* much faster if they can rapidly compute approximate estimates of the gradient (re: smaller batches) rather than slowly computing the exact gradient.

The key points to consider when choosing your batch size:

1. Larger batches = more accurate estimates of the gradient, but with less than linear returns.
2. If examples in the batch are processed in parallel (as is typical), then memory roughly scales with batch size.
3. Small batches can offer a regularizing effect. Generalization error is often best for a batch size of 1. However, this requires a low learning rate to maintain stability and thus a longer overall training runtime.

Also, note that online SGD, where we never reuse data points, but simply update parameters as new data comes in, gives an unbiased estimator of the exact gradient of the generalization error (the risk). Once data samples are reused (e.g. when training with multiple epochs), the gradient estimates become biased. The interesting point here is that the availability of increasingly massive datasets is making single-epoch<sup>12</sup> training more common. In such cases, *overfitting is no longer an issue*, but rather underfitting and computational efficiency.

**Ill-conditioning** of the Hessian matrix  $\mathbf{H}$  can cause SGD to get “stuck” in the sense that even very small steps increase the cost function. Recall that a second-order Taylor series expansion of the cost function predicts that an SGD step of  $-\epsilon \mathbf{g}$  will add

$$\frac{1}{2} \epsilon^2 \mathbf{g}^T \mathbf{H} \mathbf{g} - \epsilon \mathbf{g}^T \mathbf{g} \quad (61)$$

to the cost. If  $\mathbf{H}$  has a large condition number (re: if  $\mathbf{H}$  is ill-conditioned), then the range of possible values,  $[1/\lambda_{\max}, 1/\lambda_{\min}]$ , for  $\mathbf{g}^T \mathbf{H} \mathbf{g}$  can become very large. In particular, if  $\mathbf{g}^T \mathbf{H} \mathbf{g}$  exceeds  $\epsilon \mathbf{g}^T \mathbf{g}$ , then the SGD step will *increase* the cost!

---

<sup>12</sup>Or even less, i.e. not using all of the training data.

**Training algorithms.** Below, I list some popular training algorithms and their update equations.

- **SGD.**

$$\mathbf{g} \leftarrow \frac{1}{m} \nabla_{\boldsymbol{\theta}} \sum_i^m L(f(\mathbf{x}^{(i)}; \boldsymbol{\theta}), \mathbf{y}^{(i)}) \quad (62)$$

$$\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} - \epsilon \mathbf{g} \quad (63)$$

- **Momentum.**

$$\mathbf{g} \leftarrow \frac{1}{m} \nabla_{\boldsymbol{\theta}} \sum_i^m L(f(\mathbf{x}^{(i)}; \boldsymbol{\theta}), \mathbf{y}^{(i)}) \quad (64)$$

$$\mathbf{v} \leftarrow \alpha \mathbf{v} - \epsilon \mathbf{g} \quad (65)$$

$$\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} + \mathbf{v} \quad (66)$$

- **Nesterov Momentum.** Gradient computations instead evaluated after current velocity is applied.

$$\mathbf{g} \leftarrow \frac{1}{m} \nabla_{\boldsymbol{\theta}} \sum_i^m L(f(\mathbf{x}^{(i)}; \boldsymbol{\theta} + \alpha \mathbf{v}), \mathbf{y}^{(i)}) \quad (67)$$

$$\mathbf{v} \leftarrow \alpha \mathbf{v} - \epsilon \mathbf{g} \quad (68)$$

$$\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} + \mathbf{v} \quad (69)$$

- **AdaGrad.** Different learning rate for each model parameter. Individually adapts the learning rates of all model parameters by scaling them inversely proportional to the square root of the sum of all historical squared values of the gradient. Empirically, can result in premature and excessive decrease in the effective learning rate.

$$\mathbf{g} \leftarrow \frac{1}{m} \nabla_{\boldsymbol{\theta}} \sum_i^m L(f(\mathbf{x}^{(i)}; \boldsymbol{\theta}), \mathbf{y}^{(i)}) \quad (70)$$

$$\mathbf{r} \leftarrow \mathbf{r} + \mathbf{g} \odot \mathbf{g} \quad (71)$$

$$\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} - \frac{\epsilon}{\delta + \sqrt{\mathbf{r}}} \odot \mathbf{g} \quad (72)$$

where the gradient accumulation variable  $\mathbf{r}$  is initialized to the zero vector, and the fraction and square root in the last equation is applied element-wise.

- **RMSProp.** Modifies AdaGrad by changing the gradient accumulation into an exponentially weighted moving average.

$$\mathbf{g} \leftarrow \frac{1}{m} \nabla_{\boldsymbol{\theta}} \sum_i^m L(f(\mathbf{x}^{(i)}; \boldsymbol{\theta}), \mathbf{y}^{(i)}) \quad (73)$$

$$\mathbf{r} \leftarrow \rho \mathbf{r} + (1 - \rho) \mathbf{g} \odot \mathbf{g} \quad (74)$$

$$\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} - \frac{\epsilon}{\delta + \sqrt{\mathbf{r}}} \odot \mathbf{g} \quad (75)$$

It is also common to modify RMSProp to use Nesterov momentum.

- **Adam.** So-named to mean “adaptive moments.” We now call  $\mathbf{r}$  the 2nd moment (variance) variable, and introduce  $\mathbf{s}$  as the 1st moment (mean) variable, where the moments are for the [true] gradient; the new variables act as estimates of the moments [since we estimate the gradient with a simple average over a minibatch]. Note that these moments are uncentered.

$$\mathbf{g} \leftarrow \frac{1}{m} \nabla_{\boldsymbol{\theta}} \sum_i^m L(f(\mathbf{x}^{(i)}; \boldsymbol{\theta}), \mathbf{y}^{(i)}) \quad (76)$$

$$\mathbf{s} \leftarrow \frac{1}{1 - \rho_1^{t-1}} [\rho_1 \mathbf{s} + (1 - \rho_1) \mathbf{g}] \quad (77)$$

$$\mathbf{r} \leftarrow \frac{1}{1 - \rho_2^{t-1}} [\rho_2 \mathbf{r} + (1 - \rho_2) \mathbf{g} \odot \mathbf{g}] \quad (78)$$

$$\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} - \frac{\epsilon}{\delta + \sqrt{\mathbf{r}}} \odot \mathbf{s} \quad (79)$$

where the factors proportional to the  $\rho$  values serve to correct for bias in the moment estimators.

## Convolutional Neural Networks (Ch. 9)

Table of Contents Local

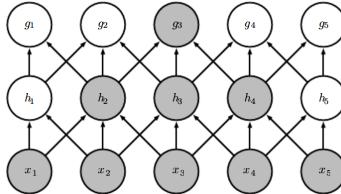
Written by Brandon McKinzie

We use a 2-D image  $I$  as our input (and therefore require a 2-D kernel  $K$ ). Note that most neural networks do not technically implement convolution<sup>13</sup>, but instead implement a related function called the *cross-correlation*, defined as

$$S(i, j) = (I * K)(i, j) = \sum_m \sum_n I(i + m, j + n)K(m, n) \quad (9.6)$$

Convolution leverages the following three important ideas:

- **Sparse interactions[/connectivity/weights].** Individual input units only interact/connect with a subset of the output units. Accomplished by making the kernel smaller than the input. It's important to recognize that the receptive field of the units in the deeper layers of a convolutional network is *larger* than the receptive field of the units in the shallow layers, as seen below.



- **Parameter sharing.**
- **Equivariance** (to translation). Changes in inputs [to a function] cause output to change in the same way. Specifically,  $f$  is equivariant to  $g$  if  $f(g(x)) = g(f(x))$ . For convolution,  $g$  would be some function that translates the input.

---

<sup>13</sup>Technically the convolution output is defined as

$$S(i, j) = (I * K)(i, j) = \sum_m \sum_n I(m, n)K(i - m, j - n) \quad (9.4)$$

$$= (K * I)(i, j) = \sum_m \sum_n I(i - m, j - n)K(m, n) \quad (9.5)$$

where 9.5 can be asserted due to commutativity of convolution.

**Pooling.** Helps make the representation approximately **invariant** to small translations of the input. The use of pooling can be viewed as adding an infinitely strong prior<sup>14</sup> that the function the layer learns must be invariant to small translations.

### Additional common tricks<sup>15</sup>.

- **Local Response Normalization (LRN)**<sup>16</sup>. Purpose is to aid generalization ability.

Let  $a_{x,y}^i$  denote the activity of a neuron computed by applying kernel  $i$  at position  $(x, y)$  and then applying the ReLU nonlinearity. The response-normalized activity  $b_{x,y}^i$  is given by the expression

$$b_{x,y}^i = \frac{a_{x,y}^i}{\left( k + \alpha \sum_j (a_{x,y}^j)^2 \right)^\beta} \quad (80)$$

where  $j$  runs from  $[i - n/2]_+$  to  $\min(N - 1, i + n/2)$ , and  $N$  is the total number of kernels in the given layer<sup>17</sup>. Authors used  $k = 2$ ,  $n = 5$ ,  $\alpha = 10^{-4}$ ,  $\beta = 0.75$ .

- **Batch Normalization**<sup>18</sup>. BN Allows us to use much higher learning rates and be less careful about initialization. Algorithm defined in image below, where each element of the batch,  $x_i \equiv x_i^{(k)}$  (where we drop the  $k$  for notational simplicity), represents the  $k$ th activation output from the previous layer [for the  $i$ th sample in the batch] and about to be fed as input to the current layer.

<b>Input:</b>	Values of $x$ over a mini-batch: $\mathcal{B} = \{x_1 \dots m\}$ ;
	Parameters to be learned: $\gamma, \beta$
<b>Output:</b>	$\{y_i = \text{BN}_{\gamma, \beta}(x_i)\}$
<hr/>	
$\mu_{\mathcal{B}} \leftarrow \frac{1}{m} \sum_{i=1}^m x_i$	// mini-batch mean
$\sigma_{\mathcal{B}}^2 \leftarrow \frac{1}{m} \sum_{i=1}^m (x_i - \mu_{\mathcal{B}})^2$	// mini-batch variance
$\hat{x}_i \leftarrow \frac{x_i - \mu_{\mathcal{B}}}{\sqrt{\sigma_{\mathcal{B}}^2 + \epsilon}}$	// normalize
$y_i \leftarrow \gamma \hat{x}_i + \beta \equiv \text{BN}_{\gamma, \beta}(x_i)$	// scale and shift

Note that one can model each layer's activations as arising from some distribution. When we feed data to a network, we model the data as coming from some data-generating

---

<sup>14</sup>Where the distribution of this prior is over all possible functions learnable by the model.

<sup>15</sup>Collected on my own. In other words, not from the deep learning book, but rather a bunch of disconnected resources over time.

<sup>16</sup>From section 3.3 of Krizhevsky et al. (2012). AlexNet paper.

<sup>17</sup>In other words, the summation is over the adjacent kernel maps, with [total] window size  $n$  (manually chosen). The min/max just says  $n/2$  to the left (right) unless that would be past the leftmost (rightmost) kernel map.

<sup>18</sup>From “Batch Normalization: Accelerating Deep Network Training by Reducing Internal Covariate Shift” by Ioffe et al.

distribution. Similarly, we can model the activations that occur when feeding the data through as coming from some activation-layer-generating distribution. The problem with this model is that the process of updating our weights during training changes the distribution of activations for each layer, which can make the learning task more difficult. Batch normalization's goal is to reduce this *internal covariate shift*, and is motivated by the practice of normalizing our data to have zero mean and unit variance.

### Intuition of some math.

- **Q:** How to intuitively understand the commutativity of convolution?
  - **A:** You must first realize that, *independent of which formula we're thinking of*, as one index into either the image or kernel increases (decreases), the other decreases (increases); they increment in opposite directions. The difference between the two formulas (9.4 and 9.5 in textbook), is just that we start at different ends of the summations (when you actually substitute in the numbers). Note that this doesn't require any symmetry on the boundaries of the summation about zero; it truly is a property of the convolution.
- **Q:** What do the authors mean by "we have flipped the kernel"?
  - **A:** Not much, and it's poor wording. They didn't *do* anything, that is just part of the definition of the convolution. They literally just mean that the convolution has the property that as one index increases, the other decreases (see previous answer). The cross-correlation, however, has the property that as one index increases, the other increases, too.

## Sequence Modeling (RNNs) (Ch. 10)

Table of Contents Local

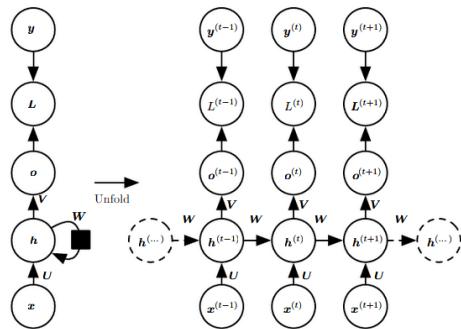
Written by Brandon McKinzie

## 2.5.1 REVIEW: THE BASICS OF RNNs

## Notation/Architecture Used.

- **U:** input  $\rightarrow$  hidden.
- **W:** hidden  $\rightarrow$  hidden.
- **V:** hidden  $\rightarrow$  output.
- **Activations:**  $\tanh$  [hidden] and softmax [after output].
- **Misc. Details:**  $\mathbf{x}^{(t)}$  is a vector of inputs fed at time  $t$ . Recall that RNNs can be unfolded for any desired number of steps  $\tau$ . For example, if  $\tau = 3$ , the general functional representation output of an RNN is  $\mathbf{s}^{(3)} = f(\mathbf{s}^{(2)}; \boldsymbol{\theta}) = f(f(\mathbf{s}^{(1)}; \boldsymbol{\theta}); \boldsymbol{\theta})$ . Typical RNNs read information out of the state  $\mathbf{h}$  to make predictions.

Shape of  $\mathbf{x}^{(t)}$  fixed, e.g.  
vocab length.



Black square on  
recurrent connection  $\equiv$   
interaction w/delay of a  
single time step.

**Forward Propagation & Loss.** Specify initial state  $\mathbf{h}^{(0)}$ . Then, for each time step from  $t = 1$  to  $t = \tau$ , feed input sequence  $\mathbf{x}^{(t)}$  and compute the output sequence  $\mathbf{o}^{(t)}$ . To determine the loss at each time-step,  $L^{(t)}$ , we compare softmax( $\mathbf{o}^{(t)}$ ) with (one-hot)  $\mathbf{y}^{(t)}$ .

$$\mathbf{h}^{(t)} = \tanh(\mathbf{a}^{(t)}) \quad \text{where} \quad \mathbf{a}^{(t)} = b + W\mathbf{h}^{(t-1)} + U\mathbf{x}^{(t)} \quad (10.9/8)$$

$$\hat{\mathbf{y}}^{(t)} = \text{softmax}(\mathbf{o}^{(t)}) \quad \text{where} \quad \mathbf{o}^{(t)} = c + V\mathbf{h}^{(t)} \quad (10.11/10)$$

Note that this is an example of an RNN that maps input seqs to output seqs of the same length<sup>19</sup>. We can then compute, e.g., the log-likelihood loss  $L = \sum_t L^{(t)}$  over all time steps as:

$$L = - \sum_t \log \left( p_{model} \left[ \mathbf{y}^{(t)} \mid \{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(t)}\} \right] \right) \quad (10.12/13/14)$$

Convince yourself this is identical to cross-entropy.

<sup>19</sup>Where “same length” is related to the number of timesteps (i.e.  $\tau$  input steps means  $\tau$  output steps), not anything about the actual shapes/sizes of each individual input/output.

where  $y^{(t)}$  is the **ground-truth** (one-hot vector) at time  $t$ , whose probability of occurring is given by the corresponding element of  $\hat{y}^{(t)}$

## Back-Propagation Through Time.

1. **Internal-Node Gradients.** In what follows, when considering what is included in the chain rule(s) for gradients with respect to a node  $N$ , just need to consider paths from it [through its **descendents**] to loss node(s).

- **Output nodes.** For any given time  $t$ , the node  $\mathbf{o}^{(t)}$  has only one direct descendant, the loss node  $L^{(t)}$ . Since no other loss nodes can be reached from  $\mathbf{o}^{(t)}$ , it is the only one we need consider in the gradient.

$$\begin{aligned}
 (\nabla_{\mathbf{o}^{(t)}} L)_i &= \frac{\partial L}{\partial \mathbf{o}_i^{(t)}} \\
 &= \frac{\partial L}{\partial L^{(t)}} \cdot \frac{\partial L^{(t)}}{\partial \mathbf{o}_i^{(t)}} \\
 &= (1) \cdot \frac{\partial L^{(t)}}{\partial \mathbf{o}_i^{(t)}} \\
 &= \frac{\partial}{\partial \mathbf{o}_i^{(t)}} \left\{ -\log(\hat{y}_{y^{(t)}}^{(t)}) \right\} \\
 &= -\frac{\partial}{\partial \mathbf{o}_i^{(t)}} \left\{ \log \left( \frac{e^{\mathbf{o}_{y^{(t)}}^{(t)}}}{\sum_j e^{\mathbf{o}_j^{(t)}}} \right) \right\} \\
 &= -\frac{\partial}{\partial \mathbf{o}_i^{(t)}} \left\{ \mathbf{o}_{y^{(t)}}^{(t)} - \log \left( \sum_j e^{\mathbf{o}_j^{(t)}} \right) \right\} \\
 &= -\left\{ \mathbf{1}_{i,y^{(t)}} - \frac{\partial}{\partial \mathbf{o}_i^{(t)}} \log \left( \sum_j e^{\mathbf{o}_j^{(t)}} \right) \right\}
 \end{aligned} \tag{81}$$

Ground-truth  $y^{(t)}$  here is a **scalar**, interpreted as the index of the correct label of output vector.

$$\mathbf{1}_{i,y^{(t)}} = \begin{cases} 1 & y^{(t)} = i \\ 0 & \text{otherwise} \end{cases}$$

$$\begin{aligned}
 &= -\left\{ \mathbf{1}_{i,y^{(t)}} - \frac{1}{\sum_j e^{\mathbf{o}_j^{(t)}}} \frac{\partial \sum_j e^{\mathbf{o}_j^{(t)}}}{\partial \mathbf{o}_i^{(t)}} \right\} \\
 &= -\left\{ \mathbf{1}_{i,y^{(t)}} - \frac{e^{\mathbf{o}_i^{(t)}}}{\sum_j e^{\mathbf{o}_j^{(t)}}} \right\} \\
 &= -\left\{ \mathbf{1}_{i,y^{(t)}} - \hat{y}_i^{(t)} \right\} \\
 &= \hat{y}_i^{(t)} - \mathbf{1}_{i,y^{(t)}}
 \end{aligned} \tag{10.18}$$

which leaves all entries of  $\mathbf{o}^{(t)}$  unchanged *except* for the entry corresponding to the true label, which will become negative in the gradient. All this means is, since we

want to increase the probability of this entry, driving this value up will *decrease* the loss (hence negative) and driving any other entries up will *increase* the loss proportional to its current estimated probability (driving up an [incorrect] entry that is already high is “worse” than driving up a small [incorrect entry]).

- **Hidden nodes.** First, consider the simplest hidden node to take the gradient of, the last one,  $\mathbf{h}^{(\tau)}$  (simplest because only one descendant [path] reaching any loss node(s)).

$$\begin{aligned}
(\nabla_{\mathbf{h}^{(\tau)}} L)_i &= \frac{\partial L}{\partial \mathbf{h}^{(\tau)}} \sum_{k=1}^{n_{out}} \frac{\partial L^{(\tau)}}{\partial \mathbf{o}_k^{(\tau)}} \frac{\partial \mathbf{o}_k^{(\tau)}}{\partial \mathbf{h}_i^{(\tau)}} \\
&= \sum_{k=1}^{n_{out}} (\nabla_{\mathbf{o}^{(\tau)}} L)_k \frac{\partial \mathbf{o}_k^{(\tau)}}{\partial \mathbf{h}_i^{(\tau)}} \\
&= \sum_{k=1}^{n_{out}} (\nabla_{\mathbf{o}^{(\tau)}} L)_k \frac{\partial}{\partial \mathbf{h}_i^{(\tau)}} \left\{ c_k + \sum_{j=1}^{n_{hid}} V_{kj} \mathbf{h}_j^{(\tau)} \right\} \\
&= \sum_{k=1}^{n_{out}} (\nabla_{\mathbf{o}^{(\tau)}} L)_k V_{ki} \\
&= \sum_{k=1}^{n_{out}} (V^T)_{ik} (\nabla_{\mathbf{o}^{(\tau)}} L)_k \\
&= (V^T \nabla_{\mathbf{o}^{(\tau)}} L)_i
\end{aligned} \tag{10.19}$$

Before proceeding, **notice the following useful pattern:** If two nodes  $a$  and  $b$ , each containing  $n_a$  and  $n_b$  neurons, are fully connected by parameter matrix  $W_{n_b \times n_a}$  and directed like  $a \rightarrow b \rightarrow L$ , then<sup>20</sup>  $\nabla_a L = W^T \nabla_b L$ . Using this result, we can then iterate and take gradients back in time from  $t = \tau - 1$  to  $t = 1$  as follows:

$$\nabla_{\mathbf{h}^{(t)}} L = \left( \frac{\partial \mathbf{h}^{(t+1)}}{\partial \mathbf{h}^{(t)}} \right)^T (\nabla_{\mathbf{h}^{(t+1)}} L) + \left( \frac{\partial \mathbf{o}^{(t)}}{\partial \mathbf{h}^{(t)}} \right)^T (\nabla_{\mathbf{o}^{(t)}} L) \tag{10.20} \quad \frac{d}{dx} \tanh(x) = 1 - \tanh^2(x)$$

$$\begin{aligned}
&= W^T (\nabla_{\mathbf{h}^{(t+1)}} L) \text{diag}(1 - \tanh^2(\mathbf{a}^{(t+1)})) + V^T (\nabla_{\mathbf{o}^{(t)}} L) \\
&= W^T (\nabla_{\mathbf{h}^{(t+1)}} L) \text{diag}(1 - (\mathbf{h}^{(t+1)})^2) + V^T (\nabla_{\mathbf{o}^{(t)}} L)
\end{aligned} \tag{10.21} \quad (\text{diag}(\mathbf{a}))_{ii} \triangleq a_i$$

2. **Parameter Gradients.** Now we can compute the gradients for the parameter matrices/vectors, where it is crucial to remember that a given parameter matrix (e.g.  $U$ ) is shared across *all* time steps  $t$ . We can treat tensor derivatives in the same form as

---

<sup>20</sup>More generally,

$$\nabla_a L = \left( \frac{\partial b}{\partial a} \right)^T \nabla_b L$$

which is a good example of how vector derivatives map into a matrix. For example, let  $\mathbf{a} \in \mathbb{R}^{n_a}$  and  $\mathbf{b} \in \mathbb{R}^{n_b}$ . Then

$$\frac{\partial \mathbf{b}}{\partial \mathbf{a}} \in \mathbb{R}^{n_b \times n_a}$$

previously done with vectors after a quick abstraction: For any tensor  $\mathbf{X}$  of arbitrary rank (e.g. if rank-4 then index like  $\mathbf{X}_{ijkl}$ ), use single variable (e.g.  $i$ ) to represent the complete tuple of indices<sup>21</sup>.

- **Bias parameters [vectors].** These are nothing new, since just vectors.

$$\begin{aligned} (\nabla_{\mathbf{c}} L) &= \sum_t \left( \frac{\partial \mathbf{o}^{(t)}}{\partial \mathbf{c}^{(t)}} \right)^T (\nabla_{\mathbf{o}^{(t)}} L) \\ &= \sum_t (\nabla_{\mathbf{o}^{(t)}} L) \end{aligned} \quad (10.22)$$

$$\begin{aligned} (\nabla_{\mathbf{c}} L) &= \sum_t \left( \frac{\partial \mathbf{h}^{(t)}}{\partial \mathbf{b}^{(t)}} \right)^T (\nabla_{\mathbf{h}^{(t)}} L) \\ &= \sum_t \text{diag} \left( 1 - (\mathbf{h}^{(t)})^2 \right) (\nabla_{\mathbf{h}^{(t)}} L) \end{aligned} \quad (10.23)$$

- **$\mathbf{V}$  ( $n_{out} \times n_{hid}$ ).**

$$\nabla_{\mathbf{V}} L = \sum_t^{\tau} \nabla_{\mathbf{V}} L^{(t)} \quad (82a)$$

$$= \sum_t^{\tau} \nabla_{\mathbf{V}} L^{(t)}(\mathbf{o}_1^{(t)}, \dots, \mathbf{o}_{n_{out}}^{(t)}) \quad (82b)$$

$$= \sum_t^{\tau} \sum_i^{n_{out}} (\nabla_{\mathbf{o}^{(t)}} L)_i \nabla_{\mathbf{V}} \mathbf{o}_i^{(t)} \quad (82c)$$

$$= \sum_t^{\tau} \sum_i^{n_{out}} (\nabla_{\mathbf{o}^{(t)}} L)_i \nabla_{\mathbf{V}} \left\{ c_i + \sum_{j=1}^{n_{hid}} V_{ij} \mathbf{h}_j^{(t)} \right\} \quad (82d)$$

$$= \sum_t^{\tau} \sum_i^{n_{out}} (\nabla_{\mathbf{o}^{(t)}} L)_i \begin{bmatrix} 0 & 0 & \dots & 0 \\ \vdots & \vdots & \dots & \vdots \\ \mathbf{h}_1^{(t)} & \mathbf{h}_2^{(t)} & \dots & \mathbf{h}_{n_{hid}}^{(t)} \\ \vdots & \vdots & \dots & \vdots \\ 0 & 0 & \dots & 0 \end{bmatrix} \quad (82e)$$

$$= \sum_t^{\tau} (\nabla_{\mathbf{o}^{(t)}} L) (\mathbf{h}^{(t)})^T \quad (82f)$$

---

<sup>21</sup>More details on tensor derivatives: Consider the chain defined by  $\mathbf{Y} = g(\mathbf{X})$ , and  $z = f(\mathbf{Y})$ , where  $z$  is some vector. Then

$$\nabla_{\mathbf{X}} z = \sum_j (\nabla_{\mathbf{X}} Y_j) \frac{\partial z}{\partial Y_j}$$

where if 82e confuses you, see the footnote<sup>22</sup>.

- **$\mathbf{W}$**  ( $n_{hid} \times n_{hid}$ ). This one is a bit odd, since  $\mathbf{W}$  is, in a sense, even more “shared” across time steps than  $\mathbf{V}$ <sup>23</sup>. The authors here define/choose, when evaluating  $\nabla_{\mathbf{W}} h_i^{(t)}$  to only concern themselves with  $\mathbf{W} := \mathbf{W}^{(t)}$ , i.e. the direct connections to  $\mathbf{h}$  at time  $t$ .

$$\nabla_{\mathbf{W}} L = \sum_t^{\tau} \nabla_{\mathbf{W}} L^{(t)} \quad (84a)$$

$$= \sum_t^{\tau} \sum_i^{n_{hid}} (\nabla_{\mathbf{h}^{(t)}} L)_i \nabla_{\mathbf{W}^{(t)}} \mathbf{h}_i^{(t)} \quad (10.25)$$

$$= \sum_t^{\tau} \sum_i^{n_{hid}} (\nabla_{\mathbf{h}^{(t)}} L)_i \left( \text{diag} \left( 1 - (\mathbf{h}^{(t)})^2 \right) \begin{bmatrix} 0 & 0 & \dots & 0 \\ \vdots & \vdots & \dots & \vdots \\ \mathbf{h}_1^{(t-1)} & \mathbf{h}_2^{(t-1)} & \dots & \mathbf{h}_{n_{hid}}^{(t-1)} \\ \vdots & \vdots & \dots & \vdots \\ 0 & 0 & \dots & 0 \end{bmatrix} \right) \quad (84b)$$

$$= \sum_t^{\tau} \text{diag} \left( 1 - (\mathbf{h}^{(t)})^2 \right) (\nabla_{\mathbf{h}^{(t)}} L) (\mathbf{h}^{(t-1)})^T \quad (10.26)$$

- **$\mathbf{U}$**  ( $n_{hid} \times n_{in}$ ). Very similar to the previous calculation.

$$\nabla_{\mathbf{U}} L = \sum_t^{\tau} \nabla_{\mathbf{U}} L^{(t)} \quad (85a)$$

$$= \sum_t^{\tau} \sum_i^{n_{hid}} (\nabla_{\mathbf{h}^{(t)}} L)_i \nabla_{\mathbf{U}^{(t)}} \mathbf{h}_i^{(t)} \quad (10.27)$$

$$= \sum_t^{\tau} \text{diag} \left( 1 - (\mathbf{h}^{(t)})^2 \right) (\nabla_{\mathbf{h}^{(t)}} L) (\mathbf{x}^{(t)})^T \quad (10.28)$$

---

<sup>22</sup> The general lesson learned here is that, for some matrix  $\mathbf{W} \in \mathbb{R}^{a \times b}$  and vector  $\mathbf{x} \in \mathbb{R}^b$ ,

$$\sum_i \nabla_{\mathbf{W}} [(\mathbf{W}\mathbf{x})_i] = \begin{bmatrix} \mathbf{x}^T \\ \mathbf{x}^T \\ \vdots \\ \mathbf{x}^T \end{bmatrix} \quad (83)$$

where, of course, the output has the same dimensions as  $\mathbf{W}$ .

<sup>23</sup>Specifically,  $\mathbf{h}^{(t)}$  is both

- An explicit function of the parameter matrix  $\mathbf{W}^{(t)}$  directly feeding into it.
- An implicit function of all other  $\mathbf{W}^{t=i}$  that came before.

This is different than before, where we had  $\mathbf{o}^{(t)}$  not implicitly depending on earlier  $\mathbf{V}^{(t=i)}$ . In other words,  $\mathbf{h}^{(t)}$  is a descendant of all earlier (and current)  $\mathbf{W}$ .

---

### 2.5.2 RNNs AS DIRECTED GRAPHICAL MODELS

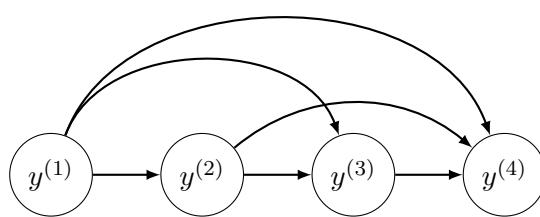
---

The advantage of RNNs is their efficient parameterization of the joint distribution over  $\mathbf{y}^{(i)}$  via parameter sharing. This introduces a built-in assumption that we can model the effect of  $y^{(i)}$  in the distant past on the current  $y^{(t)}$  *via its effect on  $h$* . We are also assuming that the conditional probability distribution over the variables at  $t + 1$  given the variables at time  $t$  is **stationary**. Next, we want to know how to draw *samples* from such a model. Specifically, how to sample from the conditional distribution  $(y^{(t)} \text{ given } y^{(t-1)})$  at each time step.

Say we want to model a sequence of scalar random variables  $\mathbb{Y} \triangleq \{y^{(1)}, \dots, y^{(\tau)}\}$  for some sequence length  $\tau$ . Without making independence assumptions just yet, we can parameterize the joint distribution  $P(\mathbb{Y})$  with basic definitions of probability:

$$P(\mathbb{Y}) \triangleq P(y^{(1)}, \dots, y^{(\tau)}) = \prod_{t=1}^{\tau} P(y^{(t)} \mid y^{(t-1)}, \dots, y^{(1)}) \quad (86)$$

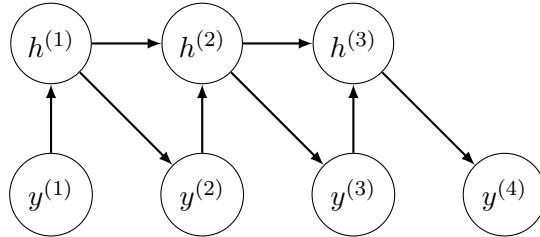
where I've drawn an example of the *complete graph* for  $\tau = 4$  below.



The complete graph can represent the direct dependencies between any pairs of  $y$  values.

If each value  $y$  could take on the same fixed set of  $k$  values, we would need to learn  $k^4$  parameters to represent the joint distribution  $P(\mathbb{Y})$ . This is clearly inefficient, since the number of parameters needed scales like  $\mathcal{O}(k^\tau)$ . If we relax the restriction that each  $y^{(i)}$  must depend *directly* on all past  $y^{(j)}$ , we can considerably reduce the number of parameters needed to compute the probability of some particular sequence.

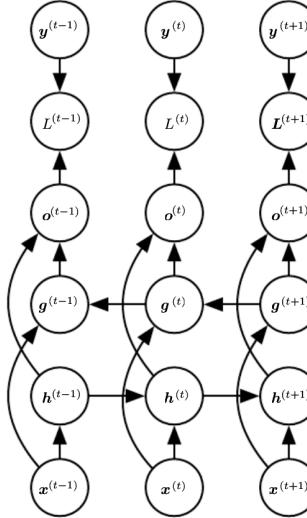
We could include latent variables  $h$  at each timestep that capture the dependencies, reminiscent of a classic RNN:



Since in the RNN case all factors  $P(h^{(t)} | h^{(t-1)})$  are deterministic, we don't need any additional parameters to compute this probability<sup>24</sup>, other than the single  $m^2$  parameters needed to convert any  $h^{(t)}$  to the next  $h^{(t+1)}$  (which is shared across all transitions). Now, the number of parameters needed as a function of sequence length is constant, and as a function of  $k$  is just  $\mathcal{O}(k)$ .

Finally, to view the RNN as a graphical model, we must describe how to sample from it, namely how to sample a sequence  $\mathbf{y}$  from  $P(\mathbb{Y})$ , if parameterized by our graphical model above. In the general case where we don't know the value of  $\tau$  for our sequence  $\mathbf{y}$ , one approach is to have a EOS symbol that, if found during sampling, means we should stop there. Also, in the typical case where we actually want to model  $P(y | x)$  for input sequence  $x$ , we can reinterpret the parameters  $\theta$  of our graphical model as a function of  $x$  the input sequence. In other words, the graphical model interpretation becomes a function of  $x$ , where  $x$  determines the exact values of the probabilities the graphical model takes on – an “instance” of the graphical model.

**Bidirectional RNNs.** In many applications, it is desirable to output a prediction of  $\mathbf{y}^{(t)}$  that may depend on *the whole sequence*. For example, in speech recognition, the interpretation of words/sentences can also depend on what is *about* to be said. Below is a typical bidirectional RNN, where the inputs  $x^{(t)}$  are fed both to a “forward” RNN ( $\mathbf{h}$ ) and a “backward” RNN ( $\mathbf{g}$ ).

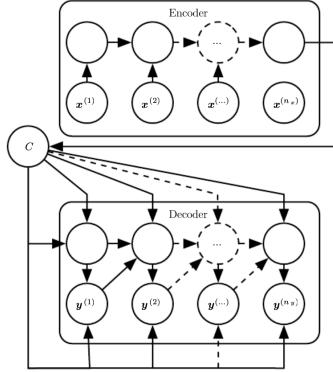


Notice how the output units  $\mathbf{o}^{(t)}$  have the nice property of depending on both the past and future while being most sensitive to input values around time  $t$ .

---

<sup>24</sup>Don't forget that, in a neural net, a variable  $y^{(t)}$  is represented by a *layer*, which itself is composed of  $k$  nodes, each associated with one of the  $k$  unique values that  $y^{(t)}$  could be.

**Encoder-Decoder Seq2Seq Architectures (10.4)** Here we discuss how an RNN can be trained to map an input sequence to output sequence which is not necessarily the same length. (Not really much of a discussion...figure below says everything.)



### 2.5.3 CHALLENGE OF LONG-TERM DEPS. (10.7)

Gradients propagated over many stages either vanish (usually) or explode. We saw how this could occur when we took parameter gradients earlier, and for weight matrices  $\mathbf{W}$  further along from the loss node, the expression for  $\nabla_{\mathbf{W}} L$  contained multiplicative Jacobian factors. Consider the (linear activation) repeated function composition of an RNN's hidden state in 10.36. We can rewrite it as a power method (10.37), and if  $\mathbf{W}$  admits an eigendecomposition (remember  $\mathbf{W}$  is necessarily square here), we can further simplify as seen in 10.38.

$$\mathbf{h}^{(t)} = \mathbf{W}^T \mathbf{h}^{(t-1)} \quad (10.36)$$

$$= (\mathbf{W}^t)^T \mathbf{h}^{(0)} \quad (10.37)$$

$$= \mathbf{Q}^T \mathbf{\Lambda}^t \mathbf{Q} \mathbf{h}^{(0)} \quad (10.38)$$

**Q:** Explain interp. of mult.  $\mathbf{h}$  by  $\mathbf{Q}$  as opposed to the usual  $\mathbf{Q}^T$  explained in the linear algebra review.

**Any component of  $\mathbf{h}^{(0)}$  that isn't aligned with the largest eigenvector will eventually be discarded.<sup>25</sup>**

If, however, we have a non-recurrent network such that the state elements are repeatedly multiplied by different  $w^{(t)}$  at each time step, the situation is different. Suppose the different  $w^{(t)}$  are i.i.d. with mean 0 and variance  $v$ . The variance of the product is easily seen to

<sup>25</sup>Make sure to think about this from the right perspective. The largest value of  $t = \tau$  in the RNNs we've seen would correspond with either (1) the largest output sequence or (2) the largest input sequence (if fixed-vector output). After we extract the output from a given forward pass, we reset the clock and either back-propagate errors (if training) or get ready to feed another sequence.

be  $\mathcal{O}(v^n)$ <sup>26</sup>. To obtain some desired variance  $v^*$  we may choose the individual weights with variance  $v = \sqrt[n]{v^*}$ .

---

#### 2.5.4 LSTMs AND OTHER GATED RNNs (10.10)

While leaky units have connection weights that are either manually chosen constants or are trainable parameters, gated RNNs generalize this to connection weights that may change *at each time step*. Furthermore, gated RNNs can learn to both accumulate and *forget*, while leaky units are designed for just accumulation<sup>27</sup>

**LSTM (10.10.1).** The idea is we want self-loops to produce paths where the gradient can flow for long durations. The self-loop weights are **gated**, meaning they are controlled by another hidden unit, interpreted as being conditioned on *context*. Listed below are the main components of the LSTM architecture.

- **Forget gate**  $f_i^{(t)} = \sigma(b_i^f + \sum_j U_{i,j}^f x_j^{(t)} + \sum_j W_{i,j}^f h_j^{(t-1)})$ .
- **Internal state**  $s_i^{(t)} = f_i^{(t)} \odot s_i^{(t-1)} + g_i^{(t)} \odot \sigma(b_i + \sum_j U_{i,j} x_j^{(t)} + \sum_j W_{i,j} h_j^{(t-1)})$ .
- **External input gate**  $g_i^{(t)} = \sigma(b_i^g + \sum_j U_{i,j}^g x_j^{(t)} + \sum_j W_{i,j}^g h_j^{(t-1)})$ .
- **Output gate**  $q_i^{(t)} = \sigma(b_i^o + \sum_j U_{i,j}^o x_j^{(t)} + \sum_j W_{i,j}^o h_j^{(t-1)})$ .

The subscript,  $i$ , identifies the cell. The superscript,  $t$ , denotes the time.

The final hidden state can then be computed via

$$h_i^{(t)} = \tanh(s_i^{(t)}) \odot q_i^{(t)} \quad (89)$$

---

<sup>26</sup>Quick sketch of (my) proof:

$$\text{Var}[w^{(i)}] = v = \mathbb{E}[(w^{(i)})^2] - \cancel{\mathbb{E}[w^{(i)}]^2} \quad (87)$$

$$\text{Var}\left[\prod_t w^{(t)}\right] = \mathbb{E}\left[\left(\prod_t w^{(t)}\right)^2\right] = \prod_t \mathbb{E}[(w^{(t)})^2] = v^n \quad (88)$$

<sup>27</sup>Q: Isn't choosing to update with higher relative weight on the present the same as forgetting? A: Sort of. It's like "soft forgetting" and will inevitably erase more/less than desired (smear). In this context, "forget" means to set the weight of a specific past cell to zero.

## Applications (Ch. 12)

Table of Contents Local

Written by Brandon McKinzie

## 2.6.1 NATURAL LANGUAGE PROCESSING (12.4)

Begins on pg. 448

**n-grams.** A **language model** defines a probability distribution over sequences of [discrete] tokens (words/characters/etc). Early models were based on the *n-gram*: a [fixed-length] sequence of  $n$  tokens. Such models define the conditional distribution for the  $n$ th token, given the  $(n - 1)$  previous tokens:

$$P(x_t \mid x_{t-(n-1)}, \dots, x_{t-1})$$

where  $x_i$  denotes the token at step/index/position  $i$  in the sequence.

To define distributions over longer sequences, we can just use Bayes rule over the shorter distributions, as usual. For example, say we want to find the [joint] distribution for some  $\tau$ -gram ( $\tau > n$ ), and we have access to an  $n$ -gram model and a [perhaps different] model for the initial sequence  $P(x_1, \dots, x_{n-1})$ . We compute the  $\tau$  distribution simply as follows:

$$P(x_1, \dots, x_\tau) = P(x_1, \dots, x_{n-1}) \prod_{t=n}^{\tau} P(x_t \mid x_{t-1}, \dots, x_{t-(n-2)}, x_{t-(n-1)}) \quad (12.5)$$

where it's important to see that each factor in the product is a distribution over a length- $n$  sequence. Since we need that initial factor, it is common to train both an  $n$ -gram model and an  $n - 1$ -gram model simultaneously.

Let's do a specific example for a trigram ( $n = 3$ ).

- **Assumptions [for this trigram model example]:**
  - For any  $n \geq 3$ ,  $P(x_n \mid x_1, \dots, x_{n-1}) = P(x_n \mid x_{n-2}, x_{n-1})$ .
  - When we get to computing the full joint distribution over some sequence of arbitrary length, we assume we have access to both  $P_3$  and  $P_2$ , the joint distributions over all subsequences of length 3 and 2, respectively.
- **Example sequence:** We want to know how to use a trigram model on the sequence ['THE', 'DOG', 'RAN', 'AWAY'].

- **Derivation:** We can use the built-in model assumption to derive the following formula.

$$\begin{aligned}
P(\text{THE DOG RAN AWAY}) &= P_3(\text{AWAY} \mid \text{THE DOG RAN}) P_3(\text{THE DOG RAN}) \\
&= P_3(\text{AWAY} \mid \text{DOG RAN}) P_3(\text{THE DOG RAN}) \\
&= \frac{P_3(\text{DOG RAN AWAY})}{P_2(\text{DOG RAN})} P_3(\text{THE DOG RAN}) \\
&= P_3(\text{THE DOG RAN}) P_3(\text{DOG RAN AWAY}) / P_2(\text{DOG RAN})
\end{aligned} \tag{12.7}$$

**Limitations of n-gram.** The last example illustrates some potential problems one may encounter that arise [if using MLE] when the full joint we seek is nonzero, but (a) some  $P_n$  factor is zero, or (b)  $P_{n-1}$  is zero. Some methods of dealing with this are as follows.

- **Smoothing:** shifting probability mass from the observed tuples to unobserved ones that are similar.
- **Back-off methods:** look up the lower-order (lower values of  $n$ )  $n$ -grams if the frequency of the context  $x_{t-1}, \dots, x_{t-(n-1)}$  is too small to use the higher-order model.

In addition,  $n$ -gram models are vulnerable to the curse of dimensionality, since most  $n$ -grams won't occur in the training set<sup>28</sup>, even for modest  $n$ .

### 2.6.2 NEURAL LANGUAGE MODELS (12.4.2)

Designed to overcome curse of dimensionality by using a distributed representation of words. Recognize that any model trained on sentences of length  $n$  and then told to generalize to new sentences [also of length  $n$ ] must deal with a space<sup>29</sup> of possible sentences that is exponential in  $n$ . Such word representations (i.e. viewing words as existing in some high-dimensional space) are often called **word embeddings**. The idea is to map the words (or sentences) from the raw high-dimensional [vocab sized] space to a smaller feature space, where similar words are closer to one another. Using distributed representations may also be used with graphical models (think Bayes' nets) in the form multiple *latent variables*.

---

<sup>28</sup>For a given vocabulary, which usually has much more than  $n$  possible words, consider how many possible sequences of length  $n$ .

<sup>29</sup>Ok I tried re-wording that from the book's confusing wording but that was also a bit confusing. Let me break it down. Say you train on a thousand sentences each of length 5. For a given vocabulary of size VOCAB\_SIZE, the number of possible sequences of length 5 is  $(\text{VOCAB\_SIZE})^5$ , which can be quite a lot more than a thousand (not to mention the possibility of duplicate training examples). To the naive model, all points in this high-dimensional space are basically the same. A neural language model, however, tries to arrange the space of possibilities in a meaningful way, so that an unforeseen sample at test time can be said "similar" as some previously seen training example. It does this by *embedding* words/sentences in a lower-dimensional feature space.

Recall that, in MLE, the  $P_n$  and  $P_{n-1}$  are usually approximated via counting occurrences in the training set

# DEEP LEARNING

# RESEARCH

## CONTENTS

3.1	Linear Factor Models (Ch. 13) . . . . .	47
3.2	Autoencoders (Ch. 14) . . . . .	50
3.3	Representation Learning (Ch. 15) . . . . .	51
3.4	Structured Probabilistic Models for DL (Ch. 16) . . . . .	52
3.4.1	Sampling from Graphical Models . . . . .	54
3.4.2	Inference and Approximate Inference . . . . .	54
3.5	Monte Carlo Methods (Ch. 17) . . . . .	56
3.6	Confronting the Partition Function (Ch. 18) . . . . .	58
3.7	Approximate Inference (Ch. 19) . . . . .	59
3.8	Deep Generative Models (Ch. 20) . . . . .	61

## Linear Factor Models (Ch. 13)

Table of Contents Local

Written by Brandon McKinzie

**Overview.** Much research is in building a *probabilistic model*<sup>30</sup> of the input,  $p_{model}(x)$ . Why? Because then we can perform *inference* to predict stuff about our environment given any of the other variables. We call the other variables **latent variables**,  $h$ , with

$$p_{model}(x) = \sum_h \Pr(h) \Pr(x | h) = \mathbb{E}_h [p_{model}(x | h)] \quad (90)$$

So what? Well, the latent variables provide another means of *data representation*, which can be useful. **Linear factor models** (LFM) are some of the simplest probabilistic models with latent variables.

A linear factor model is defined by the use of a stochastic linear decoder function that generates  $\mathbf{x}$  by adding noise to a linear transformation of  $\mathbf{h}$ .

Note that  $\mathbf{h}$  is a *vector* of arbitrary size, where we assume  $p(\mathbf{h})$  is a **factorial distribution**:  $p(\mathbf{h}) = \prod_i p(h_i)$ . This roughly means we assume the elements of  $\mathbf{h}$  are mutually independent<sup>31</sup>. The LFM describes the data-generation process as follows:

1. Sample the explanatory factors:  $\mathbf{h} \sim p(\mathbf{h})$ .
2. Sample the real-valued observable variables given the factors:

$$\mathbf{x} = \mathbf{W}\mathbf{h} + \mathbf{b} + \text{noise} \quad (91)$$

## Probabilistic PCA and Factor Analysis.

- **Factor analysis:**

$$\mathbf{h} \sim \mathcal{N}(\mathbf{h}; \mathbf{0}, \mathbf{I}) \quad (92)$$

$$\text{noise} \sim \mathcal{N}(\mathbf{0}, \boldsymbol{\psi} \equiv \text{diag}(\boldsymbol{\sigma^2})) \quad (93)$$

$$\mathbf{x} \sim \mathcal{N}(\mathbf{x}; \mathbf{b}, \mathbf{W}\mathbf{W}^T + \boldsymbol{\psi}) \quad (94)$$

where the last relation can be shown by recalling that a linear combination of Gaussian variables is itself Gaussian, and showing that  $\mathbb{E}_h [\mathbf{x}] = \mathbf{b}$ , and  $\text{Cov}(\mathbf{x}) = \mathbf{W}\mathbf{W}^T + \boldsymbol{\psi}$ .

<sup>30</sup>Whereas, before, we've been building *functions* of the input (deterministic).

<sup>31</sup>Note that, technically, this assumption isn't strictly the definition of mutual independence, which requires that every *subset* (i.e. not just the full set) of  $\{h_i \in \mathbf{h}\}$  follow this factorial property.

It is worth emphasizing the interpretation of  $\psi$  as the matrix of **conditional variances**  $\sigma_i^2$ . *Huh?* Let's take a step back. The fact that we were able to separate the distributions in the above relations for  $\mathbf{h}$  and noise is from a built-in assumption that  $\Pr(x_i|\mathbf{h}, x_{j \neq i}) = \Pr(x_i|\mathbf{h})^{32}$ .

## The Big Idea

The latent variable  $\mathbf{h}$  is a big deal because it **captures the dependencies** between the elements of  $\mathbf{x}$ . *How do I know?* Because of our assumption that the  $x_i$  are conditionally independent given  $\mathbf{h}$ . If, once we specify  $\mathbf{h}$ , all the elements of  $\mathbf{x}$  become independent, then any information about their interrelationship is hiding somewhere in  $\mathbf{h}$ .

Detailed walkthrough of Factor Analysis (a.k.a me slowly reviewing, months after taking this note):

- **Goal.** Analyze and understand the motivations behind how Factor Analysis defines the data-generation process under the framework of LFM (defined in steps 1 and 2 earlier). Assume  $\mathbf{h}$  has dimension  $n$ .
- **Prior.** Defines  $p(\mathbf{h}) := \mathcal{N}(\mathbf{h}; \mathbf{0}, \mathbf{I})$ , the unit-variance Gaussian. Explicitly,

$$p(\mathbf{h}) := \frac{1}{(2\pi)^{n/2}} e^{-\frac{1}{2} \sum_i h_i^2}$$

- **Noise.** Assumed to be drawn from a Gaussian with diagonal covariance matrix  $\psi := \text{diag}(\sigma^2)$ . Explicitly,

$$p(\text{noise} = \mathbf{a}) := \frac{1}{(2\pi)^{n/2} \prod_i \sigma_i} e^{-\frac{1}{2} \sum_i a_i^2 / \sigma_i^2}$$

- **Deriving distribution of  $\mathbf{x}$ .** We use the fact that any linear combination of Gaussians is itself Gaussian. Thus, deriving  $p(\mathbf{x})$  is reduced to computing its mean and covariance matrix.

$$\mu_x = \mathbb{E}_{\mathbf{h}} [\mathbf{W}\mathbf{h} + \mathbf{b}] \quad (95)$$

$$= \int p(\mathbf{h})(\mathbf{W}\mathbf{h} + \mathbf{b}) d\mathbf{h} \quad (96)$$

$$= \mathbf{b} + \int \frac{1}{(2\pi)^{n/2}} e^{-\frac{1}{2} \sum_i h_i^2} \mathbf{W}\mathbf{h} d\mathbf{h} \quad (97)$$

$$= \mathbf{b} \quad (98)$$

$$\text{Cov}(\mathbf{x}) = \mathbb{E} [(\mathbf{x} - \mathbb{E}[\mathbf{x}])(\mathbf{x} - \mathbb{E}[\mathbf{x}])^T] \quad (99)$$

$$= \mathbb{E} [(\mathbf{W}\mathbf{h} + \text{noise})(\mathbf{h}^T \mathbf{W}^T + \text{noise}^T)] \quad (100)$$

$$= \mathbb{E} [(\mathbf{W}\mathbf{h}\mathbf{h}^T \mathbf{W}^T)] + \psi \quad (101)$$

$$= \mathbf{W}\mathbf{W}^T + \psi \quad (102)$$

where we compute the expectation of  $\mathbf{x}$  over  $\mathbf{h}$  because  $\mathbf{x}$  is defined as a function of  $\mathbf{h}$ , and noise is always expectation zero.

---

<sup>32</sup>Due to  $\langle \text{MATH} \rangle$ , this introduces a constraint that knowing the value of some element  $x_j$  doesn't alter the probability  $\Pr(x_i = W_i \cdot h + b_i + \text{noise})$ . Given how we've defined the variable  $\mathbf{h}$ , this means that knowing  $\text{noise}_j$  provides no clues about  $\text{noise}_i$ . Mathematically, the noise must have a diagonal covariance matrix.

- **Thoughts.** Not really seeing why this is useful/noteworthy. Feels very contrived (many assumptions) and restrictive – it only applies if the dependencies between each  $x_i$  can be modeled with a random variable  $\mathbf{h}$  sampled from a unit variance Gaussian.
- **Probabilistic PCA:** Just factor analysis with  $\psi = \sigma^2 \mathbf{I}$ . So zero-mean spherical Gaussian noise. It becomes regular PCA as  $\sigma \rightarrow 0$ . Here we can use an iterative EM algorithm for estimating the parameters  $\mathbf{W}$ .

## Autoencoders (Ch. 14)

[Table of Contents](#)[Local](#)*Written by Brandon McKinzie*

**Introduction.** An autoencoder learns to copy its input to its output, via an encoder function  $\mathbf{h} = f(\mathbf{x})$  and a decoder function  $\mathbf{r} = g(\mathbf{h})$ . Modern autoencoders generalize this to allow for stochastic mappings  $p_{\text{encoder}}(\mathbf{h} \mid \mathbf{x})$  and  $p_{\text{decoder}}(\mathbf{x} \mid \mathbf{h})$ .

**Undercomplete Autoencoders.** Constrain dimension of  $\mathbf{h}$  to be smaller than that of  $\mathbf{x}$ . The learning process minimizes some  $L(\mathbf{x}, g(f(\mathbf{x})))$ , where the loss function could be e.g. mean squared error. Be careful not to have too many learnable parameters in the functions  $g$  and  $f$  (thus increasing model capacity), since that defeats the purpose of using an undercomplete autoencoder in the first place.

**Regularized Autoencoders.** We can remove the undercomplete constraint/necessity by modifying our loss function. For example, a **sparse autoencoder** one that adds a penalty  $\Omega(\mathbf{h})$  to the loss function that encourages the *activations on* (not connections to/from) the hidden layer to be sparse. One way to achieve *actual zeros* in  $\mathbf{h}$  is to use rectified linear units for the activations.

## Representation Learning (Ch. 15)

Table of Contents Local

Written by Brandon McKinzie

**Greedy Layer-Wise Unsupervised Pretraining.** Given:

- Unsupervised learning algorithm  $\mathcal{L}$  which accepts as input a training set of examples  $\mathbf{X}$ , and outputs an encoder/feature function  $f$ .
- $f^{(i)}(\tilde{\mathbf{X}})$  denotes the output of the  $i$ th layer of  $f$ , given as *immediate input* the (possibly transformed) set of examples  $\tilde{\mathbf{X}}$ .
- Let  $m$  denote the number of layers (“stages”) in the encoder function (note that each layer/stage here *must* use a representation learning algorithm for its  $\mathcal{L}$  e.g. an RBM, autoencoder, sparse coding model, etc.)

The procedure is as follows:

1. Initialize.

$$f(\cdot) \leftarrow I(\cdot) \quad (103)$$

$$\tilde{\mathbf{X}} = \mathbf{X} \quad (104)$$

2. For each layer (stage)  $i$  in range( $m$ ), do:

$$f^{(k)} = \mathcal{L}(\tilde{\mathbf{X}}) \quad (105)$$

$$f(\cdot) \leftarrow f^{(k)}(f(\cdot)) \quad (106)$$

$$\tilde{\mathbf{X}} \leftarrow f^{(k)}(\tilde{\mathbf{X}}) \quad (107)$$

In English: just apply the regular learning/training process for each layer/stage **sequentially and individually**<sup>33</sup>.

When this is complete, we can run **fine-tuning**: train all layers together (including any later layers that could not be pretrained) with a supervised learning algorithm. Note that we do indeed allow the pretrained encoding stages to be optimized here (i.e. not fixed).

---

<sup>33</sup>In other words, you proceed one layer at a time *in order*. You don’t touch layer  $i$  until the weights in layer  $i - 1$  have been learned.

## Structured Probabilistic Models for DL (Ch. 16)

Table of Contents Local

Written by Brandon McKinzie

**Motivation.** In addition to classification, we can ask probabilistic models to perform other tasks such as density estimation ( $\mathbf{x} \rightarrow p(\mathbf{x})$ ), denoising, missing value imputation, or sampling. What these [other] tasks have in common is they require a *complete understanding of the input*. Let's start with the most naive approach of modeling  $p(\mathbf{x})$ , where  $\mathbf{x}$  contains  $n$  elements, each of which can take on  $k$  distinct values: we store a lookup table of all possible  $\mathbf{x}$  and the corresponding probability value  $p(\mathbf{x})$ . This requires  $k^n$  parameters<sup>34</sup>. Instead, we use graphs to describe model structure (direct/indirect interactions) to drastically reduce the number of parameters.

**Directed Models.** Also called **belief networks** or **Bayesian networks**. Formally, a directed graphical model defined on a set of variables  $\{\mathbf{x}\}$  is defined by a DAG,  $\mathcal{G}$ , whose vertices are the random variables in the model, and a set of **local conditional probability distributions**,  $p(x_i | Pa_{\mathcal{G}}(x_i))$ , where  $Pa_{\mathcal{G}}(x_i)$  gives the parents of  $x_i$  in  $\mathcal{G}$ . The probability distribution over  $\mathbf{x}$  is given by

$$p(\mathbf{x}) = \prod_i p(x_i | Pa_{\mathcal{G}}(x_i)) \quad (108)$$

**Undirected Graphical Models.** Also called **Markov Random Fields (MRFs)** or **Markov Networks**. Appropriate for situations where interactions do not have a well-defined direction. Each **clique**  $\mathcal{C}$  (any set of nodes that are all [maximally] connected) in  $\mathcal{G}$  is associated with a factor  $\phi(\mathcal{C})$ . The factor  $\phi(\mathcal{C})$ , also called a **clique potential**, is just a function (not necessarily a probability) that outputs a number when given a possible set of values over the nodes in  $\mathcal{C}$ . The output number measures the affinity of the variables in that clique for being in the states specified by the inputs. The set of all factors in  $\mathcal{G}$  defines an **unnormalized probability distribution**:

$$\tilde{p}(\mathbf{x}) = \prod_{\mathcal{C} \in \mathcal{G}} \phi(\mathcal{C}) \quad (109)$$

Clique potentials are constrained to be nonnegative.

---

<sup>34</sup>Consider the common NLP case where our vector  $\mathbf{x}$  contains  $n$  word tokens, each of which can take on any symbol in our vocabulary of size  $v$ . If we assign  $n = 100$  and  $v = 100,000$ , which are relatively common values for this case, this amounts to  $(1e5)^{1e2} = 10^{500}$  parameters.

**The Partition Function.** To obtain a valid probability distribution, we must normalize the probability distribution:

$$p(\mathbf{x}) = \frac{1}{Z} \tilde{p}(\mathbf{x}) \quad (110)$$

$$Z = \int \tilde{p}(\mathbf{x}) d\mathbf{x} \quad (111)$$

where the normalizing function  $Z = Z(\{\phi\})$  is known as the **partition function** (physicz sh0ut0uT). It is typically intractable to compute, so we resort to approximations. Note that  $Z$  isn't even guaranteed to exist – it's only for those definitions of the clique potentials that cause the integral over  $\tilde{p}(\mathbf{x})$  to converge/be defined.

**Energy-Based Models** (EBMs). A convenient way to enforce  $\forall \mathbf{x}, \tilde{p}(\mathbf{x}) > 0$  is to use EBMs, where

$$\tilde{p}(\mathbf{x}) \triangleq \exp(-E(\mathbf{x})) \quad (112)$$

and  $E(\mathbf{x})$  is known as the **energy function**<sup>35</sup>. Many algorithms need to compute not  $p_{model}(\mathbf{x})$  but only  $\log \tilde{p}_{model}(\mathbf{x})$  (unnormalized log probabilities - logits!). For EBMs with latent variables  $\mathbf{h}$ , such algorithms are phrased in terms of the **free energy**:

$$\mathcal{F}(\mathbf{x} = x) = -\log \sum_h \exp(-E(\mathbf{x} = x, \mathbf{h} = h)) \quad (113)$$

where we sum over all possible assignments of the latent variables.

**Separation and D-Separation.** We want to know which subsets of variables are conditionally independent from each other, given the values of other subsets of variables. A set of variables  $\mathbb{A}$  is **separated** (if undirected model)/**d-separated** (if directed model) from another set of variables  $\mathbb{B}$  given a third set of variables  $\mathbb{S}$  if the graph structure implies that  $\mathbb{A}$  is independent from  $\mathbb{B}$  given  $\mathbb{S}$ .

- **Separation.** For *undirected* models. If variables  $a$  and  $b$  are connected by a path involving only unobserved variables (an **active** path), then  $a$  and  $b$  are *not* separated. Otherwise, they are separated. Any paths containing at least one observed variable are called **inactive**.
- **D-Separation**<sup>36</sup>. For *directed* models. Although there are rules that help determine whether a path between  $a$  and  $b$  is d-separated, it is simplest to just determine whether  $a$  is independent from  $b$  given any observed variables along the path.

---

<sup>35</sup>Physics throwback: this mirrors the Boltzmann factor,  $\exp(-\varepsilon/\tau)$ , which is proportional to the probability of the system being in quantum energy state  $\varepsilon$ .

<sup>36</sup>The D stands for dependence.

---

### 3.4.1 SAMPLING FROM GRAPHICAL MODELS

---

For directed graphical models, we can do **ancestral sampling** to produce a sample  $\mathbf{x}$  from the joint distribution represented by the model. Just sort the variables  $x_i$  into a topological ordering such that  $\forall i, j : j > i \iff x_i \in \text{Pa}_{\mathcal{G}}(x_j)$ . To produce the sample, just sequentially sample from the beginning,  $x_1 \sim P(x_1)$ ,  $x_2 \sim P(x_2 | \text{Pa}_{\mathcal{G}}(x_1))$ , etc.

For undirected graphical models, one simple approach is **Gibbs sampling**. Essentially, this involves drawing a conditioned sample from  $x_i \sim p(x_i | \text{neighbors}(x_i))$  for each  $x_i$ . This process is repeated many times, where each subsequent pass uses the previously sampled values in  $\text{neighbors}(x_i)$  to obtain an asymptotically converging [to the correct distribution] estimate for a sample from  $p(\mathbf{x})$ .

---

### 3.4.2 INFERENCE AND APPROXIMATE INFERENCE

---

One of the main tasks with graphical models is predicting the values of some subset of variables given another subset: inference. Although the graph structures we've discussed allow us to represent complicated, high-dimensional distributions with a reasonable number of parameters, the graphs used for deep learning are usually not restrictive enough to allow efficient inference. **Approximate inference** for deep learning usually refers to variational inference, in which we approximate the distribution  $p(\mathbf{h} | \mathbf{v})$  by seeking an approximate distribution  $q(\mathbf{h} | \mathbf{v})$  that is as close to the true one as possible.

**Example: Restricted Boltzmann Machine.** The quintessential example of how graphical models are used for deep learning. The canonical RBM is an energy-based model with **binary** visible and hidden units. Its energy function is

$$E(\mathbf{v}, \mathbf{h}) = -\mathbf{b}^T \mathbf{v} - \mathbf{c}^T \mathbf{h} - \mathbf{v}^T \mathbf{W} \mathbf{h} \quad (114)$$

where  $\mathbf{b}$ ,  $\mathbf{c}$ , and  $\mathbf{W}$  are unconstrained, real-valued, learnable parameters. One could interpret the values of the bias parameters as the affinities for the associated variable being its given value, and the value  $\mathbf{W}_{i,j}$  as the affinity of  $v_i$  being its value and  $h_j$  being its value at the same time<sup>37</sup>.

The restrictions on the RBM structure, namely the fact that there are no intra-layer connections, yields nice properties. Since  $\tilde{p}(\mathbf{h}, \mathbf{v})$  can be factored into clique potentials, we can say

---

<sup>37</sup>More concretely, remember that  $\mathbf{v}$  is a one-hot vector representing some state that can assume  $\text{len}(\mathbf{v})$  unique values, and similarly for  $\mathbf{h}$ . Then  $\mathbf{W}_{i,j}$  gives the affinity for the state associated with  $v$  being its  $i$ th value *and* the state associated with  $h$  being its  $j$ th value.

that:

$$p(\mathbf{h} \mid \mathbf{v}) = \prod_i p(h_i \mid \mathbf{v}) \quad (115)$$

$$p(\mathbf{v} \mid \mathbf{h}) = \prod_i p(v_i \mid \mathbf{h}) \quad (116)$$

Also, due to the restriction of binary variables, each of the conditionals is easy to compute, and can be quickly derived as

$$p(h_i = 1 \mid \mathbf{v}) = \sigma(c_i + \mathbf{v}^T \mathbf{W}_{:,i}) \quad (117)$$

allowing for efficient block Gibbs sampling.

## Monte Carlo Methods (Ch. 17)

Table of Contents Local

Written by Brandon McKinzie

**Monte Carlo Sampling (Basics).** We can approximate the value of a (usually prohibitively large) sum/integral by viewing it as an *expectation* under some distribution. We can then approximate its value by taking samples from the corresponding probability distribution and taking an empirical average. Mathematically, the basic idea is show below:

$$s = \int p(\mathbf{x})f(\mathbf{x})d\mathbf{x} = \mathbb{E}_p[f(\mathbf{x})] \rightarrow \hat{s}_n = \frac{1}{n} \sum_{i=1, \mathbf{x}^{(i)} \sim p}^n f(\mathbf{x}^{(i)}) \quad (118)$$

As we've seen before, the empirical average is an unbiased<sup>38</sup> estimator. Furthermore, the central limit theorem tells us that the distribution of  $\hat{s}_n$  converges to a normal distribution with mean  $s$  and variance  $\text{Var}[f(\mathbf{x})]/n$ .

**Importance Sampling.** What if it's not feasible for us to sample from  $p$ ? We can approach this a couple ways, both of which will exploit the following identity:

$$p(\mathbf{x})f(\mathbf{x}) = q(\mathbf{x}) \frac{p(\mathbf{x})f(\mathbf{x})}{q(\mathbf{x})} \quad (122)$$

- **Optimal importance sampling.** We can use the aforementioned identity/decomposition to find the **optimal  $q$**  – optimal in terms of number of samples required to achieve a given level of accuracy. First, we rewrite our estimator  $\hat{s}_p$  (they now use subscript to denote the sampling distribution) as  $\hat{s}_q$ :

$$\hat{s}_q = \frac{1}{n} \sum_{i=1, \mathbf{x}^{(i)} \sim q}^n \frac{p(\mathbf{x}^{(i)})f(\mathbf{x}^{(i)})}{q(\mathbf{x}^{(i)})} \quad (123)$$

---

<sup>38</sup> Recall that expectations on such an average are still taken over the underlying (assumed) probability distribution:

$$\mathbb{E}_p[\hat{s}_n] = \frac{1}{n} \sum_{i=1}^n \mathbb{E}_p[f(\mathbf{x}^{(i)})] \quad (119)$$

$$= \frac{1}{n} \sum_{i=1}^n s \quad (120)$$

$$= s \quad (121)$$

You should think of the expectation  $\mathbb{E}_p[f(\mathbf{x}^{(i)})]$  as the expected value of the *random sample* from the underlying distribution, which of course is  $s$ , because we defined it that way.

At first glance, it feels a little wonky, but recognize that we are *sampling from q instead of p* (i.e. if this were an integral, it would be over  $q(\mathbf{x})d\mathbf{x}$ ). The catch is that, now, the variance can be greatly sensitive to the choice of  $q$ :

$$\text{Var} [\hat{s}_q] = \text{Var} \left[ \frac{p(\mathbf{x})f(\mathbf{x})}{q(\mathbf{x})} \right] / n \quad (124)$$

with the optimal (minimum) value of  $q$  at:

$$q^* = \frac{p(\mathbf{x}) | f(\mathbf{x}) |}{Z} \quad (125)$$

- **Biased importance sampling.** Computing the optimal value of  $q$  can be as challenging/infeasible as sampling from  $p$ . Biased sampling does not require us to find a normalization constant for  $p$  or  $q$ . Instead, we compute:

$$\hat{s}_{BIS} = \frac{\sum_{i=1}^n \frac{\tilde{p}(\mathbf{x}^{(i)})}{\tilde{q}(\mathbf{x}^{(i)})} f(\mathbf{x}^{(i)})}{\sum_{i=1}^n \frac{\tilde{p}(\mathbf{x}^{(i)})}{\tilde{q}(\mathbf{x}^{(i)})}} \quad (126)$$

where  $\tilde{p}$  and  $\tilde{q}$  are the unnormalized forms of  $p$  and  $q$ , and the  $\mathbf{x}^{(i)}$  samples are still drawn from [the original/unknown]  $q$ .  $\mathbb{E} [\hat{s}_{BIS}] \neq s$  except asymptotically when  $n \rightarrow \infty$ .

## Confronting the Partition Function (Ch. 18)

Table of Contents Local

Written by Brandon McKinzie

**Noise Contrastive Estimation** (NCE) (18.6). We now estimate

$$\log p_{model}(\mathbf{x}) = \log \tilde{p}_{model}(\mathbf{x}; \boldsymbol{\theta}) + c \quad (127)$$

and explicitly learn an approximation,  $c$ , for  $-\log Z(\boldsymbol{\theta})$ . Obviously MLE would just try jacking up  $c$  to maximize this, so we adopt a surrogate supervised training problem: binary classification that a given sample  $\mathbf{x}$  belongs to the (true) data distribution  $p_{data}$  or to the noise distribution  $p_{noise}$ . We introduce binary variable  $y$  to indicate whether the sample is in the true data distribution ( $y=1$ ) or the noise distribution ( $y=0$ ). Our surrogate model is thus defined by

$$p_{joint}(y=1) = \frac{1}{2} \quad (128)$$

$$p_{joint}(\mathbf{x} | y=1) = p_{model}(\mathbf{x}) \quad (129)$$

$$p_{joint}(\mathbf{x} | y=0) = p_{noise}(\mathbf{x}) \quad (130)$$

We can now use MLE on the optimization problem,

$$\boldsymbol{\theta}, c = \arg \max_{\boldsymbol{\theta}, c} \mathbb{E}_{\mathbf{x}, y \sim p_{train}} [\log p_{joint}(y | \mathbf{x})] \quad (131)$$

$$p_{joint}(y=1 | \mathbf{x}) = \frac{p_{model}(\mathbf{x})}{p_{model}(\mathbf{x}) + p_{noise}(\mathbf{x})} \quad (132)$$

$$= \frac{1}{1 + p_{noise}(\mathbf{x})/p_{model}(\mathbf{x})} \quad (133)$$

$$= \sigma(\log p_{model}(\mathbf{x}) - \log p_{noise}(\mathbf{x})) \quad (134)$$

## Approximate Inference (Ch. 19)

Table of Contents Local

Written by Brandon McKinzie

**Overview.** Most graphical models with multiple layers of hidden variables have intractable posterior distributions. This is typically because the partition function scales exponentially with the number of units and/or due to marginalizing out latent variables. Many approximate inference approaches make use of the observation that exact inference can be described as an optimization problem.

Assume we have a probabilistic model consisting of observed variables  $\mathbf{v}$  and latent variables  $\mathbf{h}$ . We want to compute  $\log p(\mathbf{v}; \boldsymbol{\theta})$ , but it's too costly to marginalize out  $\mathbf{h}$ . Instead, we compute a lower bound  $\mathcal{L}(\mathbf{v}, \boldsymbol{\theta}, q)$  – often called the **evidence lower bound** (ELBO) or negative **variational free energy** – on  $\log p(\mathbf{v}; \boldsymbol{\theta})$ <sup>39</sup>:

$$\mathcal{L}(\mathbf{v}, \boldsymbol{\theta}, q) = \log p(\mathbf{v}; \boldsymbol{\theta}) - D_{KL}(q(\mathbf{h} | \mathbf{v}) || p(\mathbf{h} | \mathbf{v}; \boldsymbol{\theta})) \quad (135)$$

$$= -\mathbb{E}_{\mathbf{h} \sim q(\mathbf{h} | \mathbf{v})} [\log p(\mathbf{h}, \mathbf{v})] + H(q(\mathbf{h} | \mathbf{v})) \quad (136)$$

where the second form is the more canonical definition<sup>40</sup>. Note that  $\mathcal{L}(\mathbf{v}, \boldsymbol{\theta}, q)$  is a lower-bound on  $\log p(\mathbf{v}; \boldsymbol{\theta})$  by definition, since

$$\log p(\mathbf{v}; \boldsymbol{\theta}) - \mathcal{L}(\mathbf{v}, \boldsymbol{\theta}, q) = D_{KL}(q(\mathbf{h} | \mathbf{v}) || p(\mathbf{h} | \mathbf{v}; \boldsymbol{\theta})) \geq 0$$

With equality (to zero) iff  $q$  is the same distribution as  $p(\mathbf{h} | \mathbf{v})$ . In other words,  $\mathcal{L}$  can be viewed as a function parameterized by  $q$  that's maximized when  $q$  is  $p(\mathbf{h} | \mathbf{v})$ , and with maximal value  $\log p(\mathbf{v})$ . Therefore, we can cast the *inference* problem of computing the (log) probability of the observed data  $\log p(\mathbf{v})$  into an *optimization* problem of maximizing  $\mathcal{L}$ . Exact inference can be done by searching over a family of functions that contains  $p(\mathbf{h} | \mathbf{v})$ .

$q$  is an arbitrary probability distribution over  $\mathbf{h}$ . Note that the book will write  $q$  when they really mean  $q(\mathbf{h} | \mathbf{v})$ .

<sup>39</sup>Recall that  $D_{KL}(P || Q) = \mathbb{E}_{x \sim P(x)} \left[ \log \frac{P(x)}{Q(x)} \right]$

<sup>40</sup>This can be derived easily from the first form. Hint:

$$\log \frac{q(\mathbf{h} | \mathbf{v})}{p(\mathbf{h} | \mathbf{v})} = \log \frac{q(\mathbf{h} | \mathbf{v})}{p(\mathbf{h}, \mathbf{v}; \boldsymbol{\theta})/p(\mathbf{v}; \boldsymbol{\theta})}$$

**Expectation Maximization** (19.2). Technically not an approach to approximate inference, but rather an approach to learning with an approximate posterior. Popular for training models with latent variables. The EM algorithm consists of alternating between the following 2 steps until convergence:

1. **E-step.** For each training example  $\mathbf{v}^{(i)}$  (in current batch or full set), set

$$q(\mathbf{h} \mid \mathbf{v}^{(i)}) = p(\mathbf{h} \mid \mathbf{v}^{(i)}; \boldsymbol{\theta}^{(0)}) \quad (137)$$

where  $\boldsymbol{\theta}^{(0)}$  denotes the current parameter values of the model at the beginning of the E-step. This can also be interpreted as maximizing  $\mathcal{L}$  w.r.t.  $q$ .

2. **M-step.** Update the parameters  $\boldsymbol{\theta}$  by completely or partially finding

$$\arg \max_{\boldsymbol{\theta}} \sum_i \mathcal{L} \left( \mathbf{v}^{(i)}, \boldsymbol{\theta}, q(\mathbf{h} \mid \mathbf{v}^{(i)}; \boldsymbol{\theta}^{(0)}) \right) \quad (138)$$

## Deep Generative Models (Ch. 20)

Table of Contents Local

Written by Brandon McKinzie

**Boltzmann Machines** (20.1). An energy-based model over a  $d$ -dimensional binary random vector  $\mathbf{x} \in \{0, 1\}^d$ . The energy function is simply  $E(\mathbf{x}) = -\mathbf{x}^T \mathbf{U} \mathbf{x} - \mathbf{b}^T \mathbf{x}$ , i.e. parameters between all pairs of  $x_i, x_j$ , and bias parameters for each  $x_i$ <sup>41</sup>. In settings where our data consists of samples of fully observed  $\mathbf{x}$ , this is clearly limited to very simple cases, since e.g. the probability of some  $x_i$  being on is given by logistic regression from the values of the other units.

### Proof: prob of $x_i$ being on is logistic regression on other units

It's important to be as specific as possible here, since the task stated as-is is ambiguous. We want to prove that the probability of some fully observed state  $\mathbf{x}$  that has its  $i$ th element clamped to 1, which I'll denote as  $p_{i=on}(\mathbf{x})$ , is logistic regression over the other units.

To prove this, it's easier to use the conventional definition where  $\mathbf{U}$  is symmetric with zero diagonal, and we write  $E(\mathbf{x})$  as

$$E(\mathbf{x}) = - \sum_{i=1}^d \sum_{j=i+1}^d x_i U_{i,j} x_j - \sum_{i=1}^d b_i x_i \quad (139)$$

where the difference is that we explicitly only sum over the upper triangle of  $\mathbf{U}$ .

Intuitively, since  $p(\{\mathbf{x}\}_{j \neq i}) = p_{i=on}(\mathbf{x}) + p_{i=off}(\mathbf{x})$ , our final formula for  $p_{i=on}$  should only contain terms involving the parameters that interact with  $x_i$ , and only for those cases where  $x_i = 1$ . This motivates exploring the formula for  $\Delta E_i(\mathbf{x}) \triangleq E_{i=off} - E_{i=on}$  where I've dropped the explicit notation on  $\mathbf{x}$  for simplicity/space. Before jumping in to deriving this, step back and realize that  $\Delta E_i$  will only contain summation terms where either the row or column index of  $\mathbf{U}$  is  $i$ , and only for terms with bias element  $b_i$ . Since our summation is over the upper triangle of  $\mathbf{U}$ , this means terms along the slices  $U_{i,i+1:d}$  and  $U_{1:i-1,i}$ . Now there is no derivation needed and we can simply write

$$\Delta E_i = \sum_{k=i+1}^d U_{i,k} x_k + \sum_{k=1}^{i-1} x_k U_{k,i} + b_i \quad (140)$$

The goal is to use this to get a logistic-regression-like formula for  $p_{i=on}$ , so we should now think about the relationship between any given  $p(\mathbf{x})$  and the associated  $E(\mathbf{x})$ . The critical observation is that  $E(\mathbf{x}) = -\ln p(\mathbf{x}) - \ln Z$ , which therefore means

$$\Delta E_i = \ln p_{i=on}(\mathbf{x}) - \ln p_{i=off}(\mathbf{x}) = -\ln \left( \frac{1 - p_{i=on}(\mathbf{x})}{p_{i=on}(\mathbf{x})} \right) \quad (141)$$

$$\exp(-\Delta E_i) = \frac{1 - p_{i=on}(\mathbf{x})}{p_{i=on}(\mathbf{x})} = \frac{1}{p_{i=on}(\mathbf{x})} - 1 \quad (142)$$

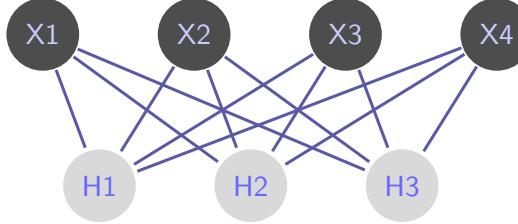
$$\therefore p_{i=on}(\mathbf{x}) = \frac{1}{1 + \exp(-\Delta E_i)} \quad (143)$$

Since  $\Delta E_i$  is a linear function of all other units, we have proven that  $p_{i=on}(\mathbf{x})$  for some state  $\mathbf{x}$  reduces to logistic regression over the other units.

---

<sup>41</sup> Authors are being lazy because it's assumed the reader is familiar (which is fair, I guess). i.e. they aren't mentioning that this formula implies that  $\mathbf{U}$  is either lower or upper triangular, and the diagonal is zero.

**Restricted Boltzmann Machines** (20.2). A BM with variables partitioned into two sets: hidden and observed. The graphical model is bipartite over the hidden and observed nodes, as I've drawn in the example below.



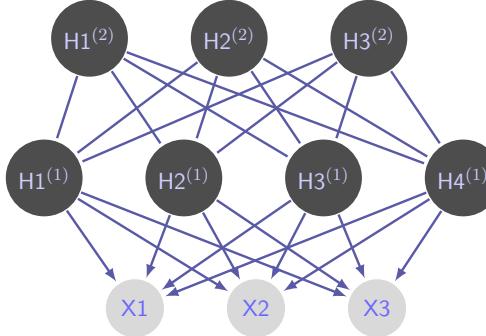
Although the joint distribution  $p(\mathbf{x}, \mathbf{h})$  has a potentially intractable partition function, the conditional distributions can be computed efficiently by exploiting independencies:

$$p(\mathbf{h} | \mathbf{x}) = \prod_{j=1}^{n_h} \sigma \left( [2\mathbf{h} - 1] \odot [\mathbf{c} + \mathbf{W}^T \mathbf{x}] \right)_j \quad (144)$$

$$p(\mathbf{x} | \mathbf{h}) = \prod_{i=1}^{n_x} \sigma \left( [2\mathbf{x} - 1] \odot [\mathbf{b} + \mathbf{W}\mathbf{h}] \right)_i \quad (145)$$

where  $\mathbf{b}$  and  $\mathbf{c}$  are the observed and hidden bias parameters, respectively.

**Deep Belief Networks** (20.3). Several layers of (usually binary) latent variables and a single observed layer. The "deepest" (away from the observed) layer connections are undirected, and all other layers are directed and pointing toward the data. I've drawn an example below.



We can sample from a DBN via first Gibbs sampling on the undirected layer, then ancestral sampling through the rest of the (directed) model to eventually obtain a sample from the visible units.

**Deep Boltzmann Machines** (20.4). Same as DBNs, but now all layers are undirected. Note that this is very close to the standard RBM, since we have a set of hidden and observed variables, except now we interpret certain subgroups of hidden units as being in a “layer”, thus allowing for connections between hidden units in adjacent layers. What’s interesting is that this still defines a bipartite graph, with odd-numbered layers on one side and even on the other<sup>42</sup>.

**Differentiable Generator Networks** (20.10.2). Use a differentiable function  $g(\mathbf{z}; \boldsymbol{\theta}^{(g)})$  to transform samples of latent variables  $\mathbf{z}$  to either (a) samples  $\mathbf{x}$ , or (b) distributions over samples  $\mathbf{x}$ . For an example of case (a), the standard procedure for sampling from  $\mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$  is to first sample from  $\mathcal{N}(\mathbf{0}, \mathbf{I})$  into a generator network consisting of a single affine layer:

**Case a:** interpret  $g(\mathbf{z})$  as emitting  $\mathbf{x}$  directly.

$$\mathbf{x} \leftarrow g(\mathbf{z}) = \boldsymbol{\mu} + \mathbf{L}\mathbf{z}$$

where  $\mathbf{L}$  is the *Cholesky decomposition*<sup>43</sup> of  $\boldsymbol{\Sigma}$ . In general, we think of the generator function  $g$  as providing a change of variables that transforms the distribution over  $\mathbf{z}$  into the desired distribution  $\mathbf{x}$ . Of course, there *is* an exact formula for doing this,

$$p_x(\mathbf{x}) = \frac{p_z(g^{-1}(\mathbf{x}))}{\left| \det \frac{\partial g}{\partial \mathbf{z}} \right|} \quad (146)$$

but it’s usually far easier to use indirect means of learning  $g$ , rather than trying to maximize/evaluate  $p_x(\mathbf{x})$  directly.

For case (b), the common approach is to train the generator net to emit conditional probabilities

**Case b:** interpret  $g(\mathbf{z})$  as emitting  $p(\mathbf{x} | \mathbf{z})$ .

$$p(x_i | \mathbf{z}) = g(\mathbf{z})_i \quad p(\mathbf{x}) = \mathbb{E}_{\mathbf{z}} [p(\mathbf{x} | \mathbf{z})] \quad (147)$$

which can also support generating discrete data (case a cannot). The challenge in training generator networks is that we often have a set of examples  $\mathbf{x}$ , but the value of  $\mathbf{z}$  for each  $\mathbf{x}$  is not fixed and known ahead of time. We’ll now look at some ways of training generator nets given only training samples for  $\mathbf{x}$ . Note that such a setting is very unlike unsupervised learning, where we typically interpret  $\mathbf{x}$  as inputs that we don’t have labels for, while here we interpret  $\mathbf{x}$  as *outputs* that we don’t know the associated inputs for.

---

<sup>42</sup>Recall that this immediately implies that units in all odd layers are conditionally independent given the even layers (and vice-versa for even to odd).

<sup>43</sup>The [unique] Cholesky decomposition of a (real-symmetric) p.d. matrix  $\mathbf{A}$  is a decomposition of the form  $\mathbf{A} = \mathbf{L}\mathbf{L}^T$ , where  $\mathbf{L}$  is lower triangular.

**Variational Autoencoders** (20.10.3). VAEs are directed models that use learned approximate inference and can be trained purely with gradient-based methods. To generate a sample  $\mathbf{x}$ , the VAE first samples  $\mathbf{z}$  from the *code distribution*  $p_{model}(\mathbf{z})$ . This sample is then fed through the a differentiable generator network  $g(\mathbf{z})$ . Finally,  $\mathbf{x}$  is sampled from  $p_{model}(\mathbf{x}; g(\mathbf{z})) = p_{model}(\mathbf{x} \mid \mathbf{z})$ .



# PAPERS AND

# TUTORIALS

## CONTENTS

4.1	WaveNet . . . . .	69
4.2	Neural Style . . . . .	72
4.3	Neural Conversation Model . . . . .	74
4.4	NMT By Jointly Learning to Align & Translate . . . . .	76
4.4.1	Detailed Model Architecture . . . . .	77
4.5	Effective Approaches to Attention-Based NMT . . . . .	79
4.6	Using Large Vocabularies for NMT . . . . .	81
4.7	Candidate Sampling – TensorFlow . . . . .	84
4.8	Attention Terminology . . . . .	86
4.9	TextRank . . . . .	88
4.9.1	Keyword Extraction . . . . .	90
4.9.2	Sentence Extraction . . . . .	91
4.10	Simple Baseline for Sentence Embeddings . . . . .	92
4.11	Survey of Text Clustering Algorithms . . . . .	94
4.11.1	Distance-based Clustering Algorithms . . . . .	97
4.11.2	Probabilistic Document Clustering and Topic Models . . . . .	98
4.11.3	Online Clustering with Text Streams . . . . .	100
4.12	Deep Sentence Embedding Using LSTMs . . . . .	102
4.13	Clustering Massive Text Streams . . . . .	105
4.14	Supervised Universal Sentence Representations (InferSent) . . . . .	107
4.15	Dist. Rep. of Sentences from Unlabeled Data (FastSent) . . . . .	108
4.16	Latent Dirichlet Allocation . . . . .	110
4.17	Conditional Random Fields . . . . .	113
4.18	Attention Is All You Need . . . . .	116
4.19	Hierarchical Attention Networks . . . . .	120
4.20	Joint Event Extraction via RNNs . . . . .	123
4.21	Event Extraction via Bidi-LSTM Tensor NNs . . . . .	125
4.22	Reasoning with Neural Tensor Networks . . . . .	127
4.23	Language to Logical Form with Neural Attention . . . . .	128
4.24	Seq2SQL: Generating Structured Queries from NL using RL . . . . .	130
4.25	SLING: A Framework for Frame Semantic Parsing . . . . .	133
4.26	Poincaré Embeddings for Learning Hierarchical Representations . . . . .	135

4.27	Enriching Word Vectors with Subword Information (FastText) . . . . .	137
4.28	DeepWalk: Online Learning of Social Representations . . . . .	139
4.29	Review of Relational Machine Learning for Knowledge Graphs . . . . .	141
4.30	Fast Top-K Search in Knowledge Graphs . . . . .	144
4.31	Dynamic Recurrent Acyclic Graphical Neural Networks (DRAGNN) . . . . .	146
4.31.1	More Detail: Arc-Standard Transition System . . . . .	149
4.32	Neural Architecture Search with Reinforcement Learning . . . . .	150
4.33	Joint Extraction of Events and Entities within a Document Context . . . . .	152
4.34	Globally Normalized Transition-Based Neural Networks . . . . .	155
4.35	An Introduction to Conditional Random Fields . . . . .	158
4.35.1	Inference (Sec. 4) . . . . .	162
4.35.2	Parameter Estimation (Sec. 5) . . . . .	165
4.35.3	Related Work and Future Directions (Sec. 6) . . . . .	168
4.36	Co-sampling: Training Robust Networks for Extremely Noisy Supervision . . . . .	169
4.37	Hidden-Unit Conditional Random Fields . . . . .	170
4.37.1	Detailed Derivations . . . . .	172
4.38	Pre-training of Hidden-Unit CRFs . . . . .	177
4.39	Structured Attention Networks . . . . .	179
4.40	Neural Conditional Random Fields . . . . .	181
4.41	Bidirectional LSTM-CRF Models for Sequence Tagging . . . . .	183
4.42	Relation Extraction: A Survey . . . . .	184
4.43	Neural Relation Extraction with Selective Attention over Instances . . . . .	187
4.44	On Herding and the Perceptron Cycling Theorem . . . . .	189
4.45	Non-Convex Optimization for Machine Learning . . . . .	191
4.45.1	Non-Convex Projected Gradient Descent (3) . . . . .	194
4.46	Improving Language Understanding by Generative Pre-Training . . . . .	195
4.47	Deep Contextualized Word Representations . . . . .	196
4.48	Exploring the Limits of Language Modeling . . . . .	198
4.49	Connectionist Temporal Classification . . . . .	200
4.50	BERT . . . . .	202
4.51	Wasserstein is all you need . . . . .	204
4.52	Noise Contrastive Estimation . . . . .	206
4.52.1	Self-Normalized NCE . . . . .	208
4.53	Neural Ordinary Differential Equations . . . . .	210
4.54	On the Dimensionality of Word Embedding . . . . .	212
4.55	Generative Adversarial Nets . . . . .	213
4.56	A Framework for Intelligence and Cortical Function . . . . .	216
4.57	Large-Scale Study of Curiosity Driven Learning . . . . .	217
4.58	Universal Language Model Fine-Tuning for Text Classification . . . . .	218
4.59	The Marginal Value of Adaptive Gradient Methods in Machine Learning . . . . .	220
4.60	A Theoretically Grounded Application of Dropout in Recurrent Neural Networks . . . . .	221
4.61	Improving Neural Language Models with a Continuous Cache . . . . .	222
4.62	Protection Against Reconstruction and Its Applications in Private Federated Learning . . . . .	223
4.63	Context Dependent RNN Language Model . . . . .	225

4.64	Strategies for Training Large Vocabulary Neural Language Models . . . . .	226
4.65	Product quantization for nearest neighbor search . . . . .	228
4.66	Large Memory Layers with Product Keys . . . . .	229
4.67	Show, Ask, Attend, and Answer . . . . .	231
4.68	Did the Model Understand the Question? . . . . .	233
4.69	XLNet . . . . .	234
4.70	Transformer-XL . . . . .	236
4.71	Efficient Softmax Approximation for GPUs . . . . .	237
4.72	Adaptive Input Representations for Neural Language Modeling . . . . .	238
4.73	Neural Module Networks . . . . .	239
4.74	Learning to Compose Neural Networks for QA . . . . .	241
4.75	End-to-End Module Networks for VQA . . . . .	243
4.76	Fast Multi-language LSTM-based Online Handwriting Recognition . . . . .	245
4.77	Multi-Language Online Handwriting Recognition . . . . .	246
4.78	Modular Generative Adversarial Networks . . . . .	248
4.79	Transfer Learning from Speaker Verification to TTS . . . . .	250

## WaveNet

Table of Contents Local

Written by Brandon McKinzie

**Introduction.**

- Inspired by recent advances in **neural autoregressive generative models**, and based on the PixelCNN architecture.
- Long-range dependencies dealt with via “dilated causal convolutions, which exhibit very large receptive fields.”

**WaveNet.** The joint probability of a waveform  $x = \{x_1, \dots, x_T\}$  is factorised as a product of conditional probabilities,

$$p(x) = \prod_{t=1}^T p(x_t | x_1, \dots, x_{t-1}) \quad (148)$$

which are modeled by a stack of convolutional layers (no pooling).

The model outputs a categorical distribution over the next value  $x_t$  with a softmax layer and it is optimized to maximize the log-likelihood of the data w.r.t. the parameters.

Main ingredient of WaveNet is *dilated* causal convolutions, illustrated below. Note the absence of recurrent connections, which makes them faster to train than RNNs, but at the cost of requiring many layers, or large filters to increase the receptive field<sup>44</sup>.

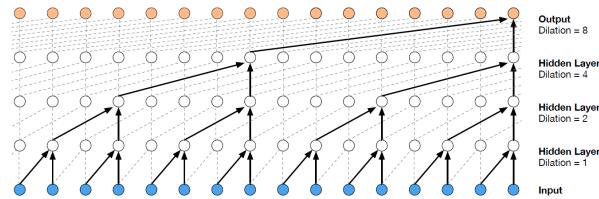


Figure 3: Visualization of a stack of *dilated* causal convolutional layers.

---

<sup>44</sup>Loose interpretation of receptive fields here is that large fields can take into account more info (back in time) as opposed to smaller fields, which can be said to be “short-sighted”

Excellent concise definition from paper:

A dilated convolution (a convolution with holes) is a convolution where the filter is applied over an area larger than its length by skipping input values with a certain step. It is equivalent to a convolution with a larger filter derived from the original filter by dilating it with zeros, but is significantly more efficient. A dilated convolution effectively allows the network to operate on a coarser scale than with a normal convolution. This is similar to pooling or strided convolutions, but here the output has the same size as the input. As a special case, dilated convolution with dilation 1 yields the standard convolution.

**Softmax distributions.** Chose to model the conditional distributions  $p(x_t \mid x_1, \dots, x_{t-1})$  with a softmax layer. To deal with the fact that there are  $2^{16}$  possible values, first apply a “ $\mu$ -law companding transformation” to data, and then quantize it to 256 possible values:

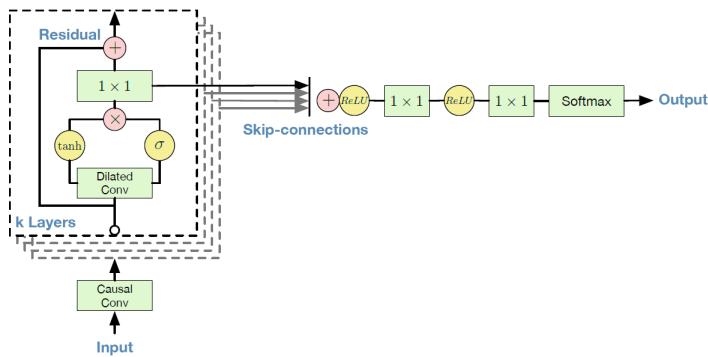
$$f(x_t) = \text{sign}(x_t) \frac{\ln(1 + \mu|x_t|)}{\ln(1 + \mu)} \quad (149)$$

which (after plotting in Wolfram) looks identical to the sigmoid function.

**Gated activation and res/skip connections.** Use the same gated activation unit as PixelCNN:

$$z = \tanh(W_{f,k} * x) \odot \sigma(W_{g,k} * x) \quad (150)$$

where  $*$  denotes conv operator,  $\odot$  denotes elem-wise mult.,  $k$  is layer index,  $f, g$  denote filter/gate, and  $W$  is learnable conv filter. This is illustrated below, along with the residual/skip connections used to speed up convergence/enable training deeper models.



**Conditional Wavenets.** Can also model conditional distribution of  $x$  given some additional  $h$  (e.g. speaker identity).

$$p(x | h) = \prod_{t=1}^T p(x_t | x_1, \dots, x_{t-1}, h) \quad (151)$$

→ **Global conditioning.** Single  $h$  that influences output dist. accross all times. Activation becomes:

$$z = \tanh(W_{f,k} * x + V_{f,k}^T h) \odot \sigma(W_{g,k} * x + V_{g,k}^T h) \quad (152)$$

→ **Local conditioning** (confusing). Have a second time-series  $h_t$ . They first transform this  $h_t$  using a “transposed conv net (learned unsampling) that maps it to a new time-series  $y = f(h)$  w/same resolution as  $x$ .

## Experiments.

- **Multi-Speaker Speech Generation.** Dataset: multi-speaker corpus of 44 hours of data from 109 different speakers<sup>45</sup>. Receptive field of 300 milliseconds.
- **Text-to-Speech.** Single-speaker datasets of 24.6 hours (English) and 34.8 hours (Chinese) speech. Locally conditioned on *linguistic features*. Receptive field of 240 milliseconds. Outperformed both LSTM-RNN and HMM.
- **Music.** Trained the WaveNets to model two music datasets: (1) 200 hours of annotated music audio, and (2) 60 hours of solo piano music from youtube. Larger receptive fields sounded more musical.
- **Speech Recognition.** “With WaveNets we have shown that layers of dilated convolutions allow the receptive field to grow longer in a much cheaper way than using LSTM units.”

**Conclusion** (verbatim): “This paper has presented WaveNet, a deep generative model of audio data that operates directly at the waveform level. WaveNets are autoregressive and combine causal filters with dilated convolutions to allow their receptive fields to grow exponentially with depth, which is important to model the long-range temporal dependencies in audio signals. We have shown how WaveNets can be conditioned on other inputs in a global (e.g. speaker identity) or local way (e.g. linguistic features). When applied to TTS, WaveNets produced samples that outperform the current best TTS systems in subjective naturalness. Finally, WaveNets showed very promising results when applied to music audio modeling and speech recognition.”

---

<sup>45</sup>Speakers encoded as ID in form of a one-hot vector

## Neural Style

Table of Contents Local

Written by Brandon McKinzie

**Notation.**

- **Content image:**  $\mathbf{p}$
- **Filter responses:** the matrix  $P^l \in \mathcal{R}^{N_l \times M_l}$  contains the activations of the filters in layer  $l$ , where  $P_{ij}^l$  would give the activation of the  $i$ th filter at position  $j$  in layer  $l$ .  $N_l$  is number of feature maps, each of size  $M_l$  (height  $\times$  width of the feature map)<sup>46</sup>.
- **Reconstructed image:**  $\mathbf{x}$  (initially random noise). Denote its corresponding filter response matrix at layer  $l$  as  $F^l$ .

**Content Reconstruction.**

1. Feed in **content image**  $\mathbf{p}$  into pre-trained network, saving any desired filter responses during the forward pass. These are interpreted as the various “encodings” of the image done by the network. Think of them analogously to “ground-truth” labels.
2. Define  $\mathbf{x}$  as the **generated image**, which we first initialize to random noise. We will be changing the pixels of  $\mathbf{x}$  via gradient descent updates.
3. Define the **loss function**. After each forward pass, evaluate with squared-error loss between the two representations at the layer of interest:

$$\mathcal{L}_{content}(\mathbf{p}, \mathbf{x}, l) = \frac{1}{2} \sum_{i,j} (F_{ij}^l - P_{ij}^l)^2 \quad (1)$$

$$\frac{\partial \mathcal{L}_{content}}{\partial F_{ij}^l} = \begin{cases} (F^l - P^l)_{ij} & F_{ij}^l > 0 \\ 0 & F_{ij}^l < 0 \end{cases} \quad (2)$$

where it appears we are assuming ReLU activations (?).

4. Compute iterative updates to  $\mathbf{x}$  via **gradient descent** until it generates the same response in a certain layer of the CNN as the original image  $\mathbf{p}$ .

---

<sup>46</sup>If not clear,  $M_l$  is a scalar, for any given value of  $l$ .

**Style Representation.** On top of the CNN responses in each layer, the authors built a style representation that computes the correlations between the different [aforementioned] filter responses. The correlation matrix for layer  $l$  is denoted in the standard way with a Gram matrix  $G^l \in \mathcal{R}^{N_l \times N_l}$ , with entries

$$G_{ij}^l = \langle F_i^l, F_j^l \rangle = \sum_k F_{ik}^l F_{jk}^l \quad (3)$$

To generate a texture that matches the style of a given image, do the following.

1. Let  $\mathbf{a}$  denote the original [style] image, with corresponding  $A^l$ . Let  $\mathbf{x}$ , initialized to random noise, denote the generated [style] image, with corresponding  $G^l$ .
2. The contribution to the loss of layer  $l$ , denoted  $E_l$ , to the total loss, denoted  $\mathcal{L}_{style}$ , is given by

$$E_l = \frac{1}{4N_l^2 M_l^2} \sum_{ij} (G_{ij}^l - A_{ij}^l)^2 \quad (4)$$

$$\mathcal{L}_{style}(\mathbf{a}, \mathbf{x}) = \sum_{l=0}^L w_l E_l \quad (5)$$

$$\frac{\partial E_l}{\partial F_{ij}^l} = \begin{cases} \frac{1}{N_l^2 M_l^2} \left( (F^l)^T (G^l - A^l) \right)_{ji} & F_{ij}^l > 0 \\ 0 & F_{ij}^l < 0 \end{cases} \quad (6)$$

where  $w_l$  are [as of yet unspecified] weighting factors of the contribution of layer  $l$  to the total loss.

**Mixing content with style.** Essentially a joint minimization that combines the previous two main ideas.

1. Begin with the following images: white noise  $\mathbf{x}$ , content image  $\mathbf{p}$ , and style image  $\mathbf{a}$ .
2. The loss function to minimize is a linear combination of 1 and 5:

$$\mathcal{L}_{total}(\mathbf{p}, \mathbf{a}, \mathbf{x}, l) = \alpha \mathcal{L}_{content}(\mathbf{p}, \mathbf{x}, l) + \beta \mathcal{L}_{style}(\mathbf{a}, \mathbf{x}) \quad (7)$$

Note that we can choose which layers  $\mathcal{L}_{style}$  uses by tweaking the layer weights  $w_l$ . For example, the authors chose to set  $w_l = 1/5$  for 'conv[1, 2, 4, 5]\_1' and 0 otherwise. For the ratio  $\alpha/\beta$ , they explored  $1 \times 10^{-3}$  and  $1 \times 10^{-4}$ .

## Neural Conversation Model

[Table of Contents](#)   [Local](#)

*Written by Brandon McKinzie*

[Reminder: **red text** means I need to come back and explain what is meant, once I understand it.]

**Abstract.** This paper presents a simple approach for conversational modeling which uses the sequence to sequence framework. It can be trained end-to-end, meaning fewer hand-crafted rules. The **lack of consistency** is a common failure of our model.

**Introduction.** Major advantage of the seq2seq model is it requires little feature engineering and domain specificity. Here, the model is tested on chat sessions from an IT helpdesk dataset of conversations, as well as movie subtitles.

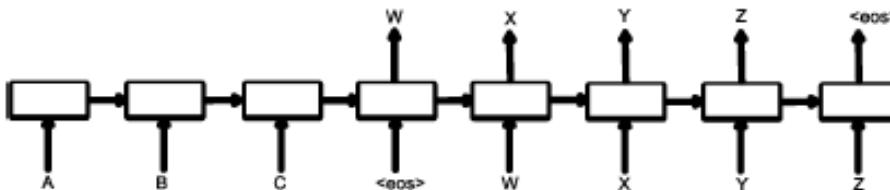
**Related Work.** The authors' approach is based on the following (linked and saved) papers on seq2seq:

- Kalchbrenner & Blunsom, 2013.
- Sutskever et al., 2014. (Describes Seq2Seq model)
- Bahdanau et al., 2014.

**Model.** Succinctly described by the authors:

The model reads the input sequence one token at a time, and predicts the output sequence, also one token at a time. During training, the true output sequence is given to the model, so learning can be done by backpropagation. The model is trained to maximize the cross entropy of the correct sequence given its context. During inference, in which the true output sequence is not observed, we simply feed the predicted output token as input to predict the next output. This is a “greedy” inference approach.

Example of less greedy approach: **beam search**.



The **thought vector** is the hidden state of the model when it receives [as input] the end of sequence symbol  $\langle eos \rangle$ , because it stores the info of the sentence, or *thought*, “ABC”. The authors acknowledge that this model will *not* be able to “solve” the problem of modeling dialogue due to the objective function not capturing the actual objective achieved through human communication, which is typically longer term and based on exchange of information [rather than next step prediction]<sup>47</sup>.

Ponder: what *would* be a reasonable objective function & model for conversation?

## IT Data & Experiment.

Reminder: Check out this git repo

- **Data Description:** Customers talking to IT support, where typical interactions are 400 words long and turn-taking is clearly signaled.
- **Training Data:** 30M tokens, 3M of which are used as validation. They built a vocabulary of the most common 20K words, and introduced special tokens indicating turn-taking and actor.
- **Model:** A single-layer LSTM with 1024 memory cells.
- **Optimization:** SGD with gradient clipping.
- **Perplexity:** At convergence, achieved **perplexity** of 8, whereas an n-gram model achieved 18.

---

<sup>47</sup>I'd imagine that, in order to model human conversation, one obvious element needed would be a *memory*. Reminds me of DeepMind's DNC. There would need to be some online filtering & output process to capture the crucial aspects/info to store in memory for later, and also some method of retrieving them when needed later. The method for retrieval would likely be some inference process where, given a sequence of inputs, the probability of them being related to some portion of memory could be trained. This would allow for conversations that stretch arbitrarily back in the past. Also, when storing the memories, I'd imagine a reasonable architecture would be some encoder-decoder for a sparse distributed representation of memory.

## NMT By Jointly Learning to Align & Translate

Table of Contents Local

Written by Brandon McKinzie

[Bahdanau et. al, 2014]. The primary motivation for me writing this is to better understand the **attention mechanism** in my sequence to sequence chatbot implementation.

**Abstract.** The authors claim that using a fixed-length vector [in the vanilla encoder-decoder for NMT] is a bottleneck. They propose allowing a model to (soft-)search for parts of a source sentence that are relevant to predicting a target word, without having to form these parts as a hard segment explicitly.

### Learning to Align<sup>48</sup> and translate.

- **Decoder.** Their encoder defines the conditional output distribution as

$$p(y_i \mid y_1, \dots, y_{i-1}, \mathbf{x}) = g(y_{i-1}, s_i, c_i) \quad (153)$$

$$s_i = f(s_{i-1}, y_{i-1}, c_i) \quad (154)$$

where  $s_i$  is the RNN [decoder] hidden state at time  $i$ .

- NOTE:  $c_i$  is *not* the  $i$ th element of the standard context vector; rather, it is *itself* a distinct context vector that depends on a sequence of **annotations**  $(h_1, \dots, h_{T_x})$ . It seems that each annotation  $h_i$  is a hidden (encoder) state “that contains information about the whole input sequence with a strong focus on the parts surrounding the  $i$ -th word of the input sequence.”
- The context vector  $c_i$  is computed as follows:

$$c_i = \sum_{j=1}^{T_x} \alpha_{ij} h_j \quad (155)$$

$$\alpha_{ij} = \frac{\exp(e_{ij})}{\sum_{k=1}^{T_x} \exp(e_{ik})} \quad (156)$$

$$e_{ij} = a(s_{i-1}, h_j) \quad (157)$$

where the function  $e_{ij}$  is given by an **alignment model** which scores how well the inputs around position  $j$  and the output at position  $i$  match.

- **Encoder.** It’s just a bidirectional RNN. What they call “annotation  $h_j$ ” is literally just a concatenated vector of  $h_j^{forward}$  and  $h_j^{backward}$

---

<sup>48</sup>By “align” the authors are referring to aligning the source-search to the relevant parts for prediction.

---

#### 4.4.1 DETAILED MODEL ARCHITECTURE

---

(Appendix A). Explained with the TensorFlow user in mind.

**Decoder Internals.** It's just a GRU. However, it will be helpful to detail how we format the inputs (given we now have attention). Wherever we'd usually pass the previous decoder state  $s_{t-1}$ , we now pass a *concatenated* state,  $[s_{t-1}, c_t]$ , that also contains the  $i$ th context vector. Below I go over the flow of information from GRU input to output:

1. **Notation:**  $y_t$  is the loop-embedded output of the decoder (prediction) at time  $t$ ,  $s_t$  is the internal hidden state of the decoder at time  $t$ , and  $c_t$  is the context vector at time  $t$ .  $\tilde{s}_t$  is the proposed/proposal state at time  $t$ .
2. **Gates:**

$$z_t = \sigma(W_z y_{t-1} + U_z [s_{t-1}, c_t]) \quad [\text{update gate}] \quad (158)$$

$$r_t = \sigma(W_r y_{t-1} + U_r [s_{t-1}, c_t]) \quad [\text{reset gate}] \quad (159)$$

$$(160)$$

3. **Proposal state:**

$$\tilde{s}_t = \tanh(W y_{t-1} + U[r_t \circ s_{t-1}, c_t]) \quad (161)$$

4. **Hidden state:**

$$s_t = (1 - z_t) \circ s_{t-1} + z_t \circ \tilde{s}_t \quad (162)$$

**Alignment Model.** All equations enumerated below are for some timestep  $t$  during the decoding process.

1. **Score:** For all  $j \in [0, L_{enc}-1]$  where  $L_{enc}$  is the number of words in the encoder sequence, compute:

$$a_j = a(s_{t-1}, h_j) = v_a^T \tanh(W_a s_{t-1} + U_a h_j) \quad (163)$$

2. **Alignments:** Feed the unnormalized alignments (scores) through a softmax so they represent a valid probability distribution.

$$a_j \leftarrow \frac{e^{a_j}}{\sum_{k=0}^{L_{enc}-1} e^{a_k}} \quad (164)$$

3. **Context:** The context vector input for our decoder at this timestep:

$$c = \sum_{j=1}^{L_{enc}} a_j h_j \quad (165)$$

**Decoder Outputs.** All below is for some timestep  $t$  during the decoding process. To find the probability of some (one-hot) word  $y$  [at timestep  $t$ ]:

$$\Pr(y | s, c) \propto e^{y^T W_o u} \quad (166)$$

$$u = [\max\{\tilde{u}_{2j-1}, \tilde{u}_{2j}\}]_{j=1,\dots,\ell}^T \quad (167)$$

$$\tilde{u} = U_o[s_{t-1}, c] + V_o y_{t-1} \quad (168)$$

**N.B.:** From reading other (and more recent) papers, these last few equations do not appear to be the way it is usually done (thank the lord). See Luong's work for a much better approach.

## Effective Approaches to Attention-Based NMT

Table of Contents Local

Written by Brandon McKinzie

[Luong et. al, 2015]

**Attention-Based Models.** For attention especially, the devil is in the details, so I'm going to go into somewhat excruciating detail here to ensure no ambiguities remain. For both global and local attention, the following information holds true:

- “At each time step  $t$  in the decoding phase, both approaches first take as input the hidden state  $\mathbf{h}_t$  at the top layer of a stacking LSTM.”
- Then, they derive [with different methods] a context vector  $\mathbf{c}_t$  to capture source-side info.
- Given  $\mathbf{h}_t$  and  $\mathbf{c}_t$ , they both compute the **attentional hidden state** as:

$$\tilde{\mathbf{h}}_t = \tanh(\mathbf{W}_c[\mathbf{c}_t; \mathbf{h}_t]) \quad (169)$$

- Finally, the predictive distribution (decoder outputs) is given by feeding this through a softmax:

$$p(y_t | y_{<t}, x) = \text{softmax}(\mathbf{W}_s \tilde{\mathbf{h}}_t) \quad (170)$$

**Global Attention.** Now I'll describe in detail the processes involved in  $\mathbf{h}_t \rightarrow \mathbf{a}_t \rightarrow \mathbf{c}_t \rightarrow \tilde{\mathbf{h}}_t$ .

1.  $\mathbf{h}_t$ : Compute the hidden state  $\mathbf{h}_t$  in the normal way (not obvious if you've read e.g. Bahdanau's work...)
2.  $\mathbf{a}_t$ :
  - (a) Compute the **scores** between  $\mathbf{h}_t$  and each source  $\bar{\mathbf{h}}_s$ , where our options are:

$$\text{score}(\mathbf{h}_t, \bar{\mathbf{h}}_s) = \begin{cases} \mathbf{h}_t^T \bar{\mathbf{h}}_s & \text{dot} \\ \mathbf{h}_t^T \mathbf{W}_a \bar{\mathbf{h}}_s & \text{general} \\ \mathbf{v}_a^T \tanh(\mathbf{W}_a[\mathbf{h}_t; \bar{\mathbf{h}}_s]) & \text{concat} \end{cases} \quad (171)$$

- (b) Compute the **alignment vector**  $\mathbf{a}_t$  of length  $L_{enc}$  (number of words in the encoder sequence):

$$\mathbf{a}_t(s) = \text{align}(\mathbf{h}_t, \bar{\mathbf{h}}_s) \quad (172)$$

$$= \frac{\exp(\text{score}(\mathbf{h}_t, \bar{\mathbf{h}}_s))}{\sum_{s'} \exp(\text{score}(\mathbf{h}_t, \bar{\mathbf{h}}_{s'}))} \quad (173)$$

3.  $\mathbf{c}_t$ : The weighted average over all source (encoder) hidden states<sup>49</sup>:

$$\mathbf{c}_t = \sum_{i=1}^{L_{enc}} \mathbf{a}_t(i) \bar{\mathbf{h}}_i \quad (174)$$

4.  $\tilde{\mathbf{h}}_t$ : For convenience, I'll copy the equation for  $\tilde{\mathbf{h}}_t$  again here:

$$\tilde{\mathbf{h}}_t = \tanh(\mathbf{W}_c[\mathbf{c}_t; \mathbf{h}_t]) \quad (175)$$

**Input-Feeding Approach.** A copy of each output  $\tilde{\mathbf{h}}_t$  is sent forward and concatenated with the inputs for the next timestep, i.e. the inputs go from  $\mathbf{x}_{t+1}$  to  $[\tilde{\mathbf{h}}_t; \mathbf{x}_{t+1}]$ .

---

<sup>49</sup>NOTE: Right after mentioning the context vector, the authors have the following cryptic footnote that may be useful to ponder: *For short sentences, we only use the top part of  $a_t$  and for long sentences, we ignore words near the end.*

## Using Large Vocabularies for NMT

Table of Contents Local

Written by Brandon McKinzie

Paper information:

- Full title: On Using Very Large Target Vocabulary for Neural Machine Translation.
- Authors: Jean, Cho, Memisevic, Bengio.
- Date: 18 Mar 2015.
- [arXiv link]

**NMT Overview.** Typical implementation is encoder-decoder network. Notation for inputs & encoder:

$$x = (x_1, \dots, x_T) \quad [\text{source sentence}] \quad (176)$$

$$h = (h_1, \dots, h_T) \quad [\text{encoder state seq}] \quad (177)$$

$$h_t = f(x_t, h_{t-1}) \quad (178)$$

where  $f$  is the function defined by the *cell state* (e.g. GRU/LSTM/etc.). Then the decoder generates the output sequence  $y$ , and with probability given below:

$$y = (y_1, \dots, y'_T) \quad [y_i \in \mathbb{Z}] \quad (179)$$

$$\Pr[y_t | y_{<t}, x] \propto e^{q(y_{t-1}, z_t, c_t)} \quad (180)$$

$$z_t = g(y_{t-1}, z_{t-1}, c_t) \quad [\text{decoder hidden?}] \quad (181)$$

$$c_t = r(z_{t-1}, h_1, \dots, h_T) \quad [\text{decoder inp?}] \quad (182)$$

The functions  $q$ ,  $g$ , and  $r$  are just placeholders – “some function of [inputs].”

As usual, model is jointly trained to maximize the conditional log-likelihood of correct translation. For  $N$  training sample pairs  $(x^n, y^n)$ , and denoting the length of the  $n$ -th target sentence as  $T_n$ , this can be written as,

$$\theta^* = \arg \max_{\theta} \sum_{n=1}^N \sum_{t=1}^{T_n} \log (\Pr[y_t^n | y_{<t}^n, x^n]) \quad (183)$$

**Model Details.** Above is the general structure. Here I'll summarize the specific model chosen by the authors.

- **Encoder.** Bi-directional, which just means  $h_t = [h_t^{backward}; h_t^{forward}]$ . The chosen cell state (the function  $f$ ) is GRU.
- **Decoder.** At each timestep, computes the following:  
→ **Context vector  $c_t$ .**

$$c_t = \sum_{i=1}^T \alpha_i h_i \quad (184)$$

$$\alpha_t = \frac{e^{a(h_t, z_{t-1})}}{\sum_k e^{a(h_k, z_{t-1})}} \quad (185)$$

$a$  is a standard single-hidden-layer NN.

→ **Decoder hidden state  $z_t$ .** Also a GRU cell. Computed based on the previous hidden state  $z_{t-1}$ , the previously generated symbol  $y_{t-1}$ , and also the computed context vector  $c_t$ .

- **Next-word probability.** They model equation 180 as<sup>50</sup>,

$$\Pr[y_t | y_{<t}, x] = \frac{1}{Z} e^{\mathbf{w}_t^T \phi(y_{t-1}, z_t, c_t) + b_t} \quad (186)$$

$$Z = \sum_{k: y_k \in V} e^{\mathbf{w}_k^T \phi(y_{t-1}, z_t, c_t) + b_k} \quad (187)$$

Reminder:  $y_i$  is an integer token, while  $\mathbf{w}_i$  is the target word vector of length vocab size

where  $\phi$  is affine transformation followed by a nonlinear activation,  $\mathbf{w}_t$  and  $b_t$  are the **target word vector** and bias.  $V$  is the set of all target *vocabulary*.

### Approximate Learning Approach.

Main idea:

“In order to avoid the growing complexity of computing the normalization constant, we propose here to use only a small subset  $V'$  of the target vocabulary at each update.”

Consider the gradient of the log-likelihood<sup>51</sup>, written in terms of the energy  $\mathcal{E}$ .

$$\nabla \log (\Pr[y_t | y_{<t}, x]) = \nabla \mathcal{E}(y_t) - \sum_{k: y_k \in V} \Pr[y_k | y_{<t}, x] \nabla \mathcal{E}(y_k) \quad (188)$$

$$\mathcal{E}(y_j) = \mathbf{w}_j^T \phi(y_{t-1}, z_t, c_t) + b_j \quad (189)$$

---

<sup>50</sup>Note: The formula for  $Z$  is correct. Notice that the only part of the RHS of  $\Pr(y_t)$  with a  $t$  is as the subscript of  $w$ . To be clear,  $w_k$  is a full word vector and the sum is over all words in the output *vocabulary*, the index  $k$  has absolutely nothing to do with timestep. They use the word target but make sure not to misinterpret that as somehow meaning target words in the sentence or something.

<sup>51</sup>**NOTE TO SELF:** After long and careful consideration, I'm concluding that the authors made a typo when defining  $\mathcal{E}(y_j)$ , which they choose to subscript all parts of the RHS with  $j$ , but that is in direct contradiction with a step-by-step derivation, which is why I have written it the way it is. I'm pretty sure my version is right, but I know you'll have to re-derive it yourself next time you see this. And you'll somehow prove me wrong. Actually, after reading on further, I doubt you'll prove me wrong. Challenge accepted, me. Have fun!

The crux of the approach is interpreting the second term as  $\mathbb{E}_P [\nabla \mathcal{E}(y)]$ , where  $P$  denotes  $Pr(y | y_{<t}, x)$ . They approximate this expectation by taking it over a subset  $V'$  of the predefined proposal distribution  $Q$ . So  $Q$  is a p.d.f. over the possible  $y_i$ , and we sample from  $Q$  to generate the elements of the subset  $V'$ .

$$\mathbb{E}_P [\nabla \mathcal{E}(y)] \approx \sum_{k: y_k \in V'} \frac{\omega_k}{\sum_{k': y_{k'} \in V'} \omega_{k'}} \nabla \mathcal{E}(y_k) \quad (190)$$

$$\omega_k = e^{\mathcal{E}(y_k) - \log Q(y_k)} \quad (191)$$

Here is some math I did that was illuminating to me; I'm not sure why the authors didn't point out these relationships.

$$\omega_k = \frac{e^{\mathcal{E}(y_k)}}{Q(y_k)} \quad \text{thus} \quad p(y_k | y_{<t}, x) = \omega_k \frac{Q(y_k)}{Z} \quad (192)$$

$$\rightarrow e^{\mathcal{E}(y_k)} = Z \cdot p(y_k | y_{<t}, x) = Q(y_k) \cdot \omega_k \quad (193)$$

### Now check this out

Below are the exact and approximate formulas for  $\mathbb{E}_P [\nabla \mathcal{E}(y)]$  written in a seductive suggestive manner. Pay careful attention to subscripts and primes.

$$\mathbb{E}_P [\nabla \mathcal{E}(y)] = \sum_{k: y_k \in V} \frac{\omega_k \cdot Q(y_k)}{\sum_{k': y_{k'} \in V} \omega_{k'} \cdot Q(y_{k'})} \nabla \mathcal{E}(y_k) \quad (194)$$

$$\mathbb{E}_P [\nabla \mathcal{E}(y)] = \sum_{k: y_k \in V'} \frac{\omega_k}{\sum_{k': y_{k'} \in V'} \omega_{k'}} \nabla \mathcal{E}(y_k) \quad (195)$$

They're almost the same! It's much easier to see why when written this way. I interpret the difference as follows: in the exact case, we explicitly attach the probabilities  $Q(y_k)$  and sum over all values in  $V$ . In the second case, by sampling a subset  $V'$  from  $Q$ , we have encoded these probabilities implicitly as the relative frequency of elements  $y_k$  in  $V'$

### How to do in practice (very important).

"In practice, we partition the training corpus and define a subset  $V'$  of the target vocabulary for each partition prior to training. Before training begins, we sequentially examine each target sentence in the training corpus and accumulate unique target words until the number of unique target words reaches the predefined threshold  $\tau$ . The accumulated vocabulary will be used for this partition of the corpus during training. We repeat this until the end of the training set is reached. Let us refer to the subset of target words used for the  $i$ -th partition by  $V'_i$ .

## Candidate Sampling – TensorFlow

Table of Contents Local

Written by Brandon McKinzie

[\[Link to article\]](#)

**What is Candidate Sampling** The goal is to learn a compatibility function  $F(x, y)$  which says something about the compatibility of a class  $y$  with a context  $x$ . Candidate sampling: for each training example  $(x_i, y_i)$ , only need to evaluate  $F(x, y)$  for a small set of classes  $\{C_i\} \subset \{L\}$ , where  $\{L\}$  is the set of all possible classes (vocab size number of elements). We represent  $F(x, y)$  as a *layer that is trained by back-prop from/within the loss function*.

**C.S. for Sampled Softmax.** I'll further narrow this down to my use case of having exactly 1 target class (word) at a given time. Any other classes are referred to as **negative** classes (for that example).

**Sampling algorithm.** For each training example  $(x_i, y_i)$ , do:

- Sample the subset  $S_i \subset L$ . How? By sampling from  $Q(y|x)$  which gives the probability of any particular  $y$  being included in  $S_i$ .
- Create the set of **candidates**, which is just  $C_i := S_i \cup y_i$ .

**Training task.** We are given this set  $C_i$  and want to find out which element of  $C_i$  is the target class  $y_i$ . In other words, we want the posterior probability that any of the  $y$  in  $C_i$  are the target class, given what we know about  $C_i$  and  $x_i$ . We can evaluate and rearrange as usual with Bayes' rule to get:

$$\Pr(y_i^{true} = y | C_i, x_i) = \frac{\Pr(y_i^{true} = y | x_i) \cdot \Pr(C_i | y_i^{true} = y, x_i)}{\Pr(C_i | x_i)} \quad (196)$$

$$= \frac{\Pr(y | x_i)}{Q(y | x_i)} \cdot \frac{1}{K(x_i, C_i)} \quad (197)$$

where they've just defined

$$K(x_i, C_i) \triangleq \frac{\Pr(C_i | x_i)}{\prod_{y' \in C_i} Q(y' | x_i) \prod_{y' \in (L - C_i)} (1 - Q(y' | x_i))} \quad (198)$$

## Clarifications.

- The learning function  $F(x, y)$  is the *input* to our softmax. It is our neural network, excluding the softmax function.
- After training our network, it should have learned the general form

$$F(x, y) = \log(\Pr(y | x)) + K(x) \quad (199)$$

which is the general form because

$$\text{Softmax}(\log(\Pr(y | x)) + K(x)) = \frac{e^{\log(\Pr(y | x)) + K(x)}}{\sum_{y'} e^{\log(\Pr(y' | x)) + K(x)}} \quad (200)$$

$$= \Pr(y | x) \quad (201)$$

Note that I've been a little sloppy here, since  $\Pr(y | x)$  up until the last line actually represented the (possibly) unnormalized/*relative* probabilities.

- **[MAIN TAKEAWAY].** Time to bring it all together. Notice that we've only trained  $F(x, y)$  to include *part* of what's needed to compute the probability of any  $y$  being the target given  $x_i$  and  $C_i$  ... equation 199 doesn't take into account  $C_i$  at all! Luckily we know the form of the full equation because it just the log of equation 197. We can easily satisfy that by subtracting  $\log(Q(y | x))$  from  $F(x, y)$  right before feeding into the softmax.

**TL;DR.** Train network to learn  $F(x, y)$  before softmax, but instead of feeding  $F(x, y)$  to softmax directly, feed

$$\text{Softmax Input: } F(x, y) - \log(Q(y | x)) \quad (202)$$

instead. That's it.

## Attention Terminology

Table of Contents Local

Written by Brandon McKinzie

Generally useful info. Seems like there are a few notations floating around, and here I'll attempt to set the record straight. The order of notes here will loosely correspond with the order that they're encountered going from encoder output to decoder output.

**Jargon.** The people in the attention business *love* obscure names for things that don't need names at all. Terminology:

- **Attentions keys/values:** Encoder output sequence.
- **Query:** Decoder [cell] state. Typically the most recent one.
- **Scores:** Values of  $e_{ij}$ . For the Bahdanau version, in code this would be computed via

$$e_i = v^T \tanh(\text{FC}(s_{i-1}) + \text{FC}(h)) \quad (203)$$

where we'd have FC be `tf.layers.fully_connected` with `num_outputs` equal to our attention size (up to us). Note that  $v$  is a vector.

- **Alignments:** output of the softmax layer on the attention scores.
- **Memory:** The  $\alpha$  matrix in the equation  $c_i = \sum_{j=1}^{T_x} \alpha_{ij} h_j$ .

When someone lazily calls some layer output the “attention”, they are usually referring to the layer *just after* the linear combination/map of encoder hidden states. You’ll often see this as some vague function of the previous decoder state, context vector, and possibly even decoder output (after project), like  $f(s_{i-1}, y_{i-1}, c_i)$ . In 99.9% of cases, this function is just a fully connected layer (if even needed) to map back to the state size for decoder input. That is it.

**From encoder to decoder.** The path of information flow from encoder outputs to decoder inputs, a non-trivial process that isn’t given the *attention* (heh) it deserves<sup>52</sup>

1. **Encoder outputs.** Tensor of shape **[batch size, sequence length, state size]**. The state is typically some RNNCell state.
  - Note: TensorFlow’s `AttentionMechanism` classes will actually convert this to **[batch size,  $L_{enc}$ , attention size]**, and refer to it as the “memory”. It is also what is returned when calling `myAttentionMech.values`.

---

<sup>52</sup>For some reason, the literature favors explaining the path “backwards”, starting with the highly abstracted “decoder inputs as a weighted sum of encoder states” and then breaking down what the weights are. Unfortunately, the weights are computed via a multi-stage process so that becomes very confusing very quick.

2. **Compute the scores.** The attention scores are the computation described by Luong/Bahdanau techniques. They both take an inner product of sorts on *copies* of the encoder outputs and decoder previous state (query). The main choices are:

$$\text{score}(\mathbf{h}_t, \bar{\mathbf{h}}_s) = \begin{cases} \mathbf{h}_t^T \bar{\mathbf{h}}_s & \text{dot} \\ \mathbf{h}_t^T \mathbf{W}_a \bar{\mathbf{h}}_s & \text{general} \\ \mathbf{v}_a^T \tanh(\mathbf{W}_a[\mathbf{h}_t; \bar{\mathbf{h}}_s]) & \text{concat} \end{cases} \quad (204)$$

**Synonyms:**  
 - scores  
 - unnormalized alignments

where the shapes are as follows (for single timestep during decoding process):

- $\bar{\mathbf{h}}_s$ : [batch size, 1, state size]
- $\mathbf{h}_t$ : [batch size, 1, state size]
- $\mathbf{W}_a$ : [batch size, state size, state size]
- $\text{score}(\mathbf{h}_t, \bar{\mathbf{h}}_s)$ : [batch size]

3. **Softmax the scores.** In the vast majority of cases, the attention scores are next fed through a softmax to convert them into a valid probability distribution. Most papers will call this some vague probability function, when in reality they are using softmaxonly.

**Synonyms:**  
 - softmax outputs  
 - attention dist.  
 - alignments

$$\mathbf{a}_t(s) = \text{align}(\mathbf{h}_t, \bar{\mathbf{h}}_s) \quad (205)$$

$$= \frac{\exp(\text{score}(\mathbf{h}_t, \bar{\mathbf{h}}_s))}{\sum_{s'} \exp(\text{score}(\mathbf{h}_t, \bar{\mathbf{h}}_{s'}))} \quad (206)$$

where the alignment vector  $\mathbf{a}_t$  has shape [batch size,  $L_{enc}$ ]

4. **Compute the context vector.** The inner product of the softmax outputs and the raw encoder outputs. This will have shape [batch size, attention size] in TensorFlow, where attention size is from the constructor for your AttentionMechanism.

**Synonyms:**  
 - context vector  
 - attention

5. **Combine context vector and decoder output:** Typically with a concat. *The result is what people mean when they say “attention”*. Luong et al. denotes this as  $\tilde{\mathbf{h}}_t$ , the decoder output at timestep  $t$ . This is what TensorFlow means by “Luong-style mechanisms output the attention.” And yes, these are used (at least for Luong) to compute the prediction:

$$\tilde{\mathbf{h}}_t = \tanh(\mathbf{W}_c [\mathbf{c}_t, \mathbf{h}_t]) \quad (207)$$

$$p(y_t | y_{<t}, x) = \text{softmax}(\mathbf{W}_s \tilde{\mathbf{h}}_t) \quad (208)$$

## TextRank

Table of Contents Local

Written by Brandon McKinzie

**Introduction.** A graph-based ranking algorithm is a way of deciding on the importance of a vertex within a graph, by taking into account global information recursively computed from the entire graph, rather than relying only on local vertex-specific information. **TextRank** is a graph-based ranking model for graphs extracted from natural language texts. The authors investigate/evaluate TextRank on unsupervised keyword and sentence extraction.

**The TextRank PageRank Model.** In general [graph-based ranking], a vertex can be ranked based on certain properties such as the number of vertices pointing to it (in-degree), how highly-ranked *those* vertices are, etc. Formally, the authors [of *PageRank*] define the score of a vertex  $V_i$  as follows:

$$S(V_i) = (1 - d) + d * \sum_{V_j \in \text{In}(V_i)} \frac{1}{|\text{Out}(V_i)|} S(V_j) \quad \text{where } d \in \mathbb{R}[0, 1] \quad (209)$$

and the damping factor  $d$  is interpreted as the probability of jumping from a given vertex<sup>53</sup> to another random vertex in the graph. In practice, the algorithm is implemented through the following steps:

- (1) Initialize all vertices with arbitrary values.<sup>54</sup>
- (2) Iterate over vertices, computing equation 4.9 until convergence [of the error rate] below a predefined threshold. The error-rate, defined as the difference between the "true score" and the score computed at iteration  $k$ ,  $S^k(V_i)$ , is *approximated* as:

$$\text{Error}^k(V_i) \approx S^k(V_i) - S^{k-1}(V_i) \quad (210)$$

<sup>53</sup>Note that  $d$  is a single parameter for the graph, i.e. it is the same for all vertices.

<sup>54</sup>The authors do not specify what they mean by arbitrary. What range? What sampling distribution? Arbitrary as in uniformly random? **EDIT:** The authors claim that the vertex values upon completion are not affected by the choice of initial value. Investigate!

**Semantic graph:** one whose structure encodes meaning between the nodes (semantic elements).

The factor  $d$  is usually set to 0.85.

**Weighted Graphs.** In contrast with the PageRank model, here we are concerned with natural language texts, which may include multiple or partial links between the units (vertices). The authors hypothesize that modifying equation 4.9 to incorporate *weighted* connections may be useful for NLP applications.

$$WS(V_i) = (1 - d) + d * \sum_{j \in \text{In}(V_i)} \frac{w_{ji}}{\sum_{V_k \in \text{Out}(V_j)} w_{jk}} WS(V_j) \quad (211)$$

$w_{ij}$  denotes the connection between vertices  $V_i$  and  $V_j$ .

where I've shown the modified part in green. The authors mention they set all weights to random values in the interval 0-10 (no explanation).

**Text as a Graph.** In general, the application of graph-based ranking algorithms to natural language texts consists of the following main steps:

- (1) Identify text units that best define the task at hand, and add them as vertices in the graph.
- (2) Identify relations that connect such text unit in order to draw edges between vertices in the graph. Edges can be directed or undirected, weighted or unweighted.
- (3) Iterate the algorithm until convergence.
- (4) Sort [in reversed order] vertices based on final score. Use the values attached to each vertex for ranking/selection decisions.

---

#### 4.9.1 KEYWORD EXTRACTION

---

**Graph.** The authors apply TextRank to extract words/phrases that are representative for a given document. The individual graph components are defined as follows:

- **Vertex:** sequence of one or more lexical units from the text.
  - In addition, we can restrict which vertices are added to the graph with syntactic filters.
  - Best filter [for the authors]: *nouns and adjectives only*.
- **Edge:** two vertices are connected if their corresponding lexical units co-occur within a window of  $N$  words<sup>55</sup>. Typically  $N \in \mathbb{Z}[2, 10]$

#### Procedure:

- (1) **Pre-Processing:** Tokenize and annotate [with POS] the text.
- (2) **Build the Graph:** Add all [single] words to the graph that pass the syntactic filter, and connect [undirected/unweighted] edges as defined earlier (co-occurrence).
- (3) **Run algorithm:** Initialize all scores to 1. For a convergence threshold of 0.0001, usually takes about 20-30 iterations.
- (4) **Post-Processing:**
  - (i) Keep the top  $T$  vertices (by score), where the authors chose  $T = |V|/3$ .<sup>56</sup> Remember that vertices are still individual words.
  - (ii) From the new subset of  $T$  keywords, collapse any that were adjacent in the original text in a single lexical unit.

**Evaluation.** The data set used is a collection of 500 abstracts, each with a set of keywords. Results are evaluated using **precision**, **recall**, and **F-measure**<sup>57</sup>. The best results were obtained with a co-occurrence window of 2 [on an undirected graph], which yielded:

Precision: 31.2%   Recall: 43.1%   F-measure: 36.2

The authors found that larger window size corresponded with lower precision, and that directed graphs performed worse than undirected graphs.

---

<sup>55</sup>That is . . . simpler than expected. *Can we do better?*

<sup>56</sup>Another approach is to have  $T$  be a fixed value, where typically  $5 < T < 20$ .

<sup>57</sup>Brief terminology review:

- **Precision:** fraction of keywords extracted that are in the "true" set of keywords.
- **Recall:** fraction of "true" keywords that are in the extracted set of keywords.
- **F-score:** combining precision and recall to get a single number for evaluation:

$$F = \frac{2pr}{p+r}$$

A PR Curve plots precision as a function of recall.

---

#### 4.9.2 SENTENCE EXTRACTION

---

**Graph.** Now we move to “sentence extraction for automatic summarization.”

- **Vertex:** a vertex is added to the graph for each sentence in the text.
- **Edge:** each weighted edge represents the similarity between two sentences. The authors use the following similarity measure between two sentences  $S_i$  and  $S_j$ :

$$\text{Similarity}(S_i, S_j) = \frac{|S_i \cap S_j|}{\log(|S_i|) + \log(|S_j|)} \quad (212)$$

where the numerator is the number of words that occur in both  $S_i$  and  $S_j$ .

The **procedure** is identical to the algorithm described for keyword extraction, except we run it on full sentences.

**Evaluation.** The data set used is 567 news articles. For each article, TextRank generates a 100-word summary (i.e. they set  $T = 100$ ). They evaluate with the ROUGE evaluation toolkit (Ngram statistics).

## Simple Baseline for Sentence Embeddings

Table of Contents Local

Written by Brandon McKinzie

**Overview.** It turns out that simply taking a weighted average of word vectors and doing some PCA/SVD is a competitive way of getting unsupervised word embeddings. Apparently it *beats* supervised learning with LSTMs (?!). The authors claim the theoretical explanation for this method lies in a latent variable generative model for sentences (of course).

Discussion based on paper by Arora et al., (2017).

**Algorithm.**

1. Compute the weighted average of the word vectors in the sentence:

$$\frac{1}{N} \sum_i^N \frac{a}{a + p(\mathbf{w}_i)} \mathbf{w}_i \quad (213)$$

The authors call their weighted average the **Smooth Inverse Frequency (SIF)**.

where  $\mathbf{w}_i$  is the word vector for the  $i$ th word in the sentence,  $a$  is a parameter, and  $p(\mathbf{w}_i)$  is the (estimated) word frequency [over the entire corpus].

2. Remove the projections of the average vectors on their first principal component (“common component removal”) (y tho?).

**Algorithm 1** Sentence Embedding

**Input:** Word embeddings  $\{v_w : w \in \mathcal{V}\}$ , a set of sentences  $\mathcal{S}$ , parameter  $a$  and estimated probabilities  $\{p(w) : w \in \mathcal{V}\}$  of the words.

**Output:** Sentence embeddings  $\{v_s : s \in \mathcal{S}\}$

- 1: **for all** sentence  $s$  in  $\mathcal{S}$  **do**
- 2:    $v_s \leftarrow \frac{1}{|s|} \sum_{w \in s} \frac{a}{a + p(w)} v_w$
- 3: **end for**
- 4: Compute the first principal component  $u$  of  $\{v_s : s \in \mathcal{S}\}$
- 5: **for all** sentence  $s$  in  $\mathcal{S}$  **do**
- 6:    $v_s \leftarrow v_s - uu^\top v_s$
- 7: **end for**

**Theory.** Latent variable generative model. The model treats corpus generation as a dynamic process, where the  $t$ -th word is produced at time step  $t$ , driven by the random walk of a **discourse vector**  $c_t \in \Re^d$  ( $d$  is size of the embedding dimension). The discourse vector is *not* pointing to a specific word; rather, it describes what is being talked about. We can tell how related (correlation) the discourse is to any word  $w$  and corresponding vector  $v_w$  by taking the inner product  $c_t \cdot v_w$ . Similarly, we model the probability of observing word  $w$  at time  $t$ ,  $w_t$ , as:

$$\Pr [w_t | c_t] \propto e^{c_t \cdot v_w} \quad (214)$$

- **The Random Walk.** If we assume that  $c_t$  doesn't change much over the words in a single sentence, we can assume it stays at some  $c_s$ . The authors claim that in their previous paper they showed that the MAP<sup>58</sup> estimate of  $c_s$  is – up to multiplication by a scalar – the average of the embeddings of the words in the sentence.
- **Improvements/Modifications to 214.**
  1. Additive term  $\alpha p(w)$  where  $\alpha$  is a scalar. Allows words to occur even if  $c_t \cdot v_w$  is very small.
  2. Common discourse vector  $c_0 \in \Re^d$ . Correction term for the most frequent discourse that is often related to syntax.
- **Model.** Given the discourse vector  $c_s$  for a sentence  $s$ , the probability that  $w$  is in the sentence (at all (?)):

$$\Pr [w | c_s] = \alpha p(w) + (1 - \alpha) \frac{e^{\tilde{c}_s \cdot v_w}}{Z_{\tilde{c}_s}} \quad (215)$$

$$\tilde{c}_s = \beta c_0 + (1 - \beta) c_s \quad (216)$$

with  $c_0 \perp c_s$  and  $Z_{\tilde{c}_s}$  is a normalization constant, taken over all  $w \in V$ .

---

<sup>58</sup>Review of MAP:

$$\theta_{MAP} = \arg \max_{\theta} \sum_i \log (p_X(x | \theta)p(\theta))$$

## Survey of Text Clustering Algorithms

[Table of Contents](#)   [Local](#)

*Written by Brandon McKinzie*

Aggarwal et al., “A Survey of Text Clustering Algorithms,” (2012).

**Introduction.** The unique characteristics for clustering *text*, as opposed to more traditional (numeric) clustering, are (1) large dimensionality but highly sparse data, (2) words are typically highly correlated, meaning the number of principal components is much smaller than the feature space, and (3) the number of words per document can vary, so we must normalize appropriately.

Common types of clustering algorithms include agglomerative clustering algorithms, partitioning algorithms, and standard parametric modeling based methods such as the EM-algorithm.

### Feature Selection.

- **Document Frequency-Based.** Using document frequency to filter *out* irrelevant features. Dealing with certain words, like “the”, should probably be taken a step further and simply removed (stop words).
- **Term Strength.** A more aggressive technique for stop-word removal. It’s used to measure how informative a word/term  $t$  is for identifying two related documents,  $x$  and  $y$ . Denoted  $s(t)$ , it is defined as:

$$s(t) = \Pr [t \in y \mid t \in x] \quad (217)$$

See ref 94 of the paper for more.

So, how do we know  $x$  and  $y$  are related to begin with? One way is a user-defined cosine similarity threshold. Say we gather a set of such *pairs* and randomly identify one of the pair as the “first” document of the pair, then we can approximate  $s(t)$  as

$$s(t) = \frac{\text{Num pairs in which } t \text{ occurs in both}}{\text{Num pairs in which } t \text{ occurs in the first of the pair}} \quad (218)$$

*In order to prune features, the term strength may be compared to the expected strength of a term which is randomly distributed in the training documents with the same frequency. If the term strength of  $t$  is not at least two standard deviations greater than that of the random word, then it is removed from the collection.*

- **Entropy-Based Ranking.** The quality of a term is measured by the entropy reduction when it is removed [from the collection]. The entropy  $E(t)$  of term  $t$  in a collection of  $n$  documents is:

$$E(t) = - \sum_{i=1}^n \sum_{j=1}^n (S_{ij} \cdot \log(S_{ij}) + (1 - S_{ij}) \cdot \log(1 - S_{ij})) \quad (219)$$

$$S_{ij} = 2^{-d_{ij}/\bar{d}} \quad (220)$$

where

- $S_{ij} \in (0, 1)$  is the similarity between doc  $i$  and  $j$ .
- $d_{ij}$  is the distance between  $i$  and  $j$  after the term  $t$  is removed
- $\bar{d}$  is the average distance between the documents after the term  $t$  is removed.

**LSI-based Methods.** Latent Semantic Indexing is based on dimensionality reduction where the new (transformed) features are a linear combination of the originals. This helps magnify the semantic effects in the underlying data. LSI is quite similar to PCA<sup>59</sup>, except that we use an approximation of the covariance matrix  $C$  which is appropriate for the sparse and high-dimensional nature of text data.

Let  $\mathbf{A} \in \mathbb{R}^{n \times d}$  be term-document matrix, where  $\mathbf{A}_{i,j}$  is the (normalized) frequency for term  $j$  in document  $i$ . Then  $\mathbf{A}^T \mathbf{A} = n \cdot \Sigma$  is the (scaled) approximation to covariance matrix<sup>60</sup>, assuming the data is mean-centered. Quick check/reminder:

$$(\mathbf{A}^T \mathbf{A})_{ij} = \mathbf{A}_{:,i}^T \mathbf{A}_{:,j} \triangleq \mathbf{a}_i^T \mathbf{a}_j \quad (221)$$

$$\approx n \cdot \mathbb{E} [\mathbf{a}_i \mathbf{a}_j] \quad (222)$$

where the expectation is technically over the underlying data distribution, which gives e.g.  $P(a_i = x)$ , the probability the  $i$ th word in our vocabulary having frequency  $x$ . Apparently, since the data is sparse, we don't have to worry much about it actually being mean-centered (**why?**).

As usual, we use the eigenvectors of  $\mathbf{A}^T \mathbf{A}$  with the largest variance in order to represent the text<sup>61</sup>. In addition:

*One excellent characteristic of LSI is that the truncation of the dimensions removes the noise effects of synonymy and polysemy, and the similarity computations are more closely affected by the semantic concepts in the data.*

---

<sup>59</sup>The difference between LSI and PCA is that PCA subtracts out the means, which destroys the sparseness of the design matrix.

<sup>60</sup>Approximation because it is based on our training data, not on true expectations over the underlying data-distribution.

<sup>61</sup>In typical collections, only about 300 to 400 eigenvectors are required for the representation.

See ref 28. of paper for  
more on LSI.

**Non-negative Matrix Factorization.** Another latent-space method (like LSI), but particularly suitable for clustering. The main characteristics of the NMF scheme:

- In LSI, the new basis system consists of a set of orthonormal vectors. This is *not* the case for NMF.
- In NMF, the vectors in the basis system **directly correspond to cluster topics**. Therefore, the cluster membership for a document may be determined by examining the largest component of the document along any of the [basis] vectors.

Assume we want to create  $k$  clusters, using our  $n$  documents and vocabulary size  $d$ . The goal of NMF is to find matrices  $\mathbf{U} \in \mathbb{R}^{n \times k}$  and  $\mathbf{V} \in \mathbb{R}^{d \times k}$  that minimize:

$$J = \frac{1}{2} \|\mathbf{A} - \mathbf{UV}^T\|_F^2 \quad (223)$$

$$= \frac{1}{2} \left( \text{tr}(\mathbf{AA}^T) - 2\text{tr}(\mathbf{AVU}^T) + \text{tr}(\mathbf{UV}^T \mathbf{VU}^T) \right) \quad (224)$$

$$\begin{aligned} u_{ij} &\geq 0 \\ v_{ij} &\geq 0 \end{aligned}$$

Note that the columns of  $\mathbf{V}$  provide the  $k$  basis vectors which correspond to the  $k$  different clusters. We can interpret this as trying to factorize  $\mathbf{A} \approx \mathbf{UV}^T$ . For each row,  $\mathbf{a}_i$ , of  $\mathbf{A}$  (document vector), this is

$$\mathbf{a}_i \approx \mathbf{u}_i \cdot \mathbf{V}^T \quad (225)$$

$$= \sum_{i=1}^k \mathbf{u}_i \mathbf{V}_i^T \quad (226)$$

Therefore, the document vector  $\mathbf{a}_i$  can be rewritten as an approximate linear (non-negative) combination of the basis vector which corresponds to the  $k$  columns of  $\mathbf{V}^T$ .

Lagrange-multiplier stuff: Our optimization problem can be solved using the Lagrange method.

- Variables to optimize: All elements of both  $\mathbf{U} = [u_{ij}]$  and  $\mathbf{V} = [v_{ij}]$
- Constraint: non-negativity, i.e.  $\forall i, j, u_{ij} \geq 0$  and  $v_{ij} \geq 0$ .
- Multipliers: Denote as matrices  $\alpha$  and  $\beta$ , with same dimensions as  $\mathbf{U}$  and  $\mathbf{V}$ , respectively.
- Lagrangian: I'll just show it here first, and then explain in this footnote<sup>62</sup>:

$$\mathcal{L} = J + \text{tr}(\alpha \cdot \mathbf{U}^T) + \text{tr}(\beta \cdot \mathbf{V}^T) \quad (227)$$

$$\text{where } \text{tr}(\alpha \cdot \mathbf{U}^T) = \sum_{i=1}^n \alpha_i \cdot \mathbf{u}_i = \sum_{i=1}^n \sum_{j=1}^n \alpha_{ij} u_{ij} \quad (228)$$

Any matrix multiplication with  $\mathbf{a} \cdot$  is just a reminder to think of the matrices as column vectors.

You should think of  $\alpha$  as a column vector of length  $n$ , and  $\mathbf{U}^T$  as a row vector of length  $n$ . The reason we prefer  $\mathcal{L}$  over just  $J$  is because now we have an *unconstrained* optimization problem.

---

<sup>62</sup> Recall that in Lagrangian minimization,  $\mathcal{L}$  takes the form of [the-function-to-be-minimized] +  $\lambda$  ([constraint-function] - [expected-value-of-constraint-at-optimum]). So the second term is expected to tend toward zero (i.e. critical point) at the optimal values. In our case, since our optimal value is sort-of (?) at 0 for any value of  $u_{ij}$  and/or  $v_{ij}$ , we just have a sum over [lagrange-mult]  $\times$  [variable].

- Optimization: Set the partials of  $\mathcal{L}$  w.r.t both  $U$  and  $V$  (separately) to zero<sup>63</sup>:

$$\frac{\partial \mathcal{L}}{\partial U} = -\mathbf{A} \cdot \mathbf{V} + \mathbf{U} \cdot \mathbf{V}^T \cdot \mathbf{V} + \boldsymbol{\alpha} = 0 \quad (229)$$

$$\frac{\partial \mathcal{L}}{\partial V} = -\mathbf{A}^T \cdot \mathbf{U} + \mathbf{V} \cdot \mathbf{U}^T \cdot \mathbf{U} + \boldsymbol{\beta} = 0 \quad (230)$$

Since, ultimately, these just say [some matrix] = 0, we can multiply both sides (element-wise) by a constant ( $x \times 0 = 0$ ). Using<sup>64</sup> the **Kuhn-Tucker conditions**  $\alpha_{ij} \cdot u_{ij} = 0$  and  $\beta_{ij} \cdot v_{ij} = 0$ , we get:

$$(\mathbf{A} \cdot \mathbf{V})_{ij} \cdot u_{ij} - (\mathbf{U} \cdot \mathbf{V}^T \cdot \mathbf{V})_{ij} \cdot u_{ij} = 0 \quad (231)$$

$$(\mathbf{A}^T \cdot \mathbf{U})_{ij} \cdot v_{ij} - (\mathbf{V} \cdot \mathbf{U}^T \cdot \mathbf{U})_{ij} \cdot v_{ij} = 0 \quad (232)$$

- Update rules:

$$u_{ij} = \frac{(\mathbf{A} \cdot \mathbf{V})_{ij} \cdot u_{ij}}{(\mathbf{U} \cdot \mathbf{V}^T \cdot \mathbf{V})_{ij}} \quad (233)$$

$$v_{ij} = \frac{(\mathbf{A}^T \cdot \mathbf{U})_{ij} \cdot v_{ij}}{(\mathbf{V} \cdot \mathbf{U}^T \cdot \mathbf{U})_{ij}} \quad (234)$$

#### 4.11.1 DISTANCE-BASED CLUSTERING ALGORITHMS

One challenge in clustering short segments of text (e.g., tweets) is that exact keyword matching may not work well. One general strategy for solving this problem is to expand text representation by exploiting related text documents, which is related to smoothing of a document language model in information retrieval.

See ref. 66 in the paper for computing similarities of short text segments.

**Agglomerative and Hierarchical Clustering.** “Agglomerative” refers to the process of bottom-up clustering to build a tree – at the bottom are leaves (documents) and internal nodes correspond to the merged groups of clusters. The different methods for merging groups of documents for the different agglomerative methods are as follows:

- **Single Linkage Clustering.** Defines similarity between two groups (clusters) of documents as the largest similarity between any pair of documents from these two groups. First, (1) compute all similarity pairs [between documents; ignore cluster labels], then (2) sort in decreasing order, and (3) walk through the list in that order, merging clusters if the pair belong to different clusters. One drawback is *chaining*: the resulting clusters assume transitivity of similarity<sup>65</sup>.

<sup>63</sup>Recall that the Lagrangian consists entirely of traces (re: scalars). Derivatives of traces with respect to matrices output the same dimension as that matrix, and derivatives are taken element-wise as always.

<sup>64</sup>i.e. the equations that follow are *not* the KT conditions, they just *use/exploit* them...

<sup>65</sup>Here, transitivity of similarity means if  $A$  is similar to  $B$ , and  $B$  is similar to  $C$ , then  $A$  is similar to  $C$ . This is not guaranteed by any means for textual similarity, and so we can end up with  $A$  and  $Z$  in the same cluster, even though they aren't similar at all.

- **Group-Average Linkage Clustering.** Similarity between two clusters is the *average* similarity over all unique pairwise combinations of documents from one cluster to the other. One way to speed up this computation with an approximation is to just compute the similarity between the mean vector of either cluster.
- **Complete Linkage Clustering.** Similarity between two clusters is the *worst-case* similarity between any pair of documents.

## Distance-Based Partitioning Algorithms.

- **K-Medoid Clustering.** Use a set of points from training data as anchors (medoids) around which the clusters are built. Key idea is we are using an optimal set of representative documents *from the original corpus*. The set of  $k$  reps is successively improved via randomized inter-changes. In each iteration, we replace a randomly picked rep in the current set of medoids with a randomly picked rep from the collection, if it improves the clustering objective function. This approach is applied until convergence is achieved.
- **K-Means Clustering.** Successively (1) assigning points to the nearest cluster centroid and then (2) re-computing the centroid of each cluster. Repeat until convergence. Requires typically few iterations (about 5 for many large data sets). Disadvantage: sensitive to initial set of seeds (initial cluster centroids). One method for improving the initial set of seeds is to use some supervision - e.g. initialize with  $k$  pre-defined topic vectors (see ref. 4 in paper for more).

K-Medoid isn't great for clustering text,  
especially short texts.

**Hybrid Approach: Scatter-Gather Method.** Use a hierarchical clustering algorithm on a sample of the corpus in order to find a robust initial set of seeds. This robust set of seeds is used in conjunction with a standard k-means clustering algorithm in order to determine good clusters. **TODO:** resume note-taking; page 19/52 of PDF.

Scatter-Gather is  
discussed in detail in ref.  
25 of the paper

### 4.11.2 PROBABILISTIC DOCUMENT CLUSTERING AND TOPIC MODELS

**Overview.** Primary assumptions in any topic modeling approach:

From pg. 31/52 of paper.

- The  $n$  documents in the corpus are assumed to each have a probability of belonging to one of  $k$  topics. Denote the probability of document  $D_i$  belonging to topic  $T_j$  as  $\Pr [T_j | D_i]$ . This allows for *soft cluster membership* in terms of probabilities.
- Each topic is associated with a probability vector, which quantifies the probability of the different terms in the lexicon for that topic. For example, consider a document that belongs completely to topic  $T_j$ . We denote the probability of term  $t_l$  occurring in that document as  $\Pr [t_l | T_j]$ .

The two main methods for topic modeling are **Probabilistic Latent Semantic Indexing** (PLSA) and **Latent Dirichlet Allocation** (LDA).

**PLSA**. We note that the two aforementioned probabilities,  $\Pr [T_j | D_i]$  and  $\Pr [t_l | T_j]$  allow us to calculate  $\Pr [t_l | D_i]$ : the probability that term  $t_l$  occurs in some document  $D_i$ :

$$\Pr [t_l | D_i] = \sum_{j=1}^k \Pr [t_l | T_j] \cdot \Pr [T_j | D_i] \quad (235)$$

which should be interpreted as a weighted average<sup>66</sup>. From here, we can generate a  $n \times d$  matrix of probabilities.

Recall that we also have our  $n \times d$  term-document matrix  $\mathbf{X}$ , where  $\mathbf{X}_{i,l}$  gives the number of times term  $l$  occurred in document  $D_i$ . This allows us to do maximum likelihood! Our negative log-likelihood,  $J$  can be derived as follows:

$$J = -\log (\Pr [\mathbf{X}]) \quad (236)$$

$$= -\log \left( \prod_{i,l} \Pr [t_l | D_i]^{\mathbf{X}_{i,l}} \right) \quad (237)$$

$$= -\sum_{i,l} \mathbf{X}_{i,l} \cdot \log (\Pr [t_l | D_i]) \quad (238)$$

Interpret  $\Pr [\mathbf{X}]$  as the joint probability of observing the words in our data and with their assoc. frequencies.

and we can plug-in eqn 235 to for evaluating  $\Pr [t_l | D_i]$ . We want to optimize the value of  $J$ , subject to the constraints:

$$(\forall T_j) : \sum_l \Pr [t_l | T_j] = 1 \quad (\forall D_i) : \sum_j \Pr [T_j | D_i] = 1 \quad (239)$$

This can be solved with a Lagrangian method, similar to the process for NMF described earlier. See page 33/52 of the paper for details.

**Latent Dirichlet Allocation** (LDA). The term-topic probabilities and topic-document probabilities are modeled with a *Dirichlet distribution* as a prior<sup>67</sup>. Typically preferred over PLSI because PLSI more prone to overfitting.

---

<sup>66</sup>This is actually pretty bad notation, and borderline incorrect.  $\Pr [T_j | D_i]$  is NOT a conditional probability! It is our prior! It is literally  $\Pr [\text{ClusterOf}(D_i) = T_j]$ .

<sup>67</sup>LDA is the Bayesian version of PLSI

---

#### 4.11.3 ONLINE CLUSTERING WITH TEXT STREAMS

Reference List: [3]: A Framework for Clustering Massive Text and Categorical Data Streams; [112]: Efficient Streaming Text Clustering; [48]: Bursty feature representation for clustering text streams; [61]: Clustering Text Data Streams (Liu et al.)

See ref. 112 for more on OSKM

**Overview.** Maintaining text clusters in real time. One method is the **Online Spherical K-Means Algorithm** (OSKM)<sup>68</sup>.

**Condensed Droplets Algorithm.** I'm calling it that because they don't call it anything – it is the algorithm in [3].

- **Fading function:**  $f(t) = 2^{-\lambda \cdot t}$ . A time-dependent weight for each data point (text stream). Non-monotonic decreasing; decays uniformly with time.
- **Decay rate:**  $\lambda = 1/t_0$ . Inverse of the half-life of the data stream.

When a cluster is created by a new point, it is allowed to remain as a trend-setting outlier for at least one half-life. During that period, if at least one more data point arrives, then the cluster becomes an active and mature cluster. If not, the trend-setting outlier is recognized as a true anomaly and is removed from the list of current clusters (*cluster death*). Specifically, this happens when the (weighted) number of points in the [single-point] cluster is 0.5. The same criterion is used to define the death of mature clusters. The statistics of the data points are referred to as **condensed droplets**, which represent the word distributions within a cluster, and can be used in order to compute the similarity of an incoming data point to the cluster. Main idea of algorithm is as follows:

1. Initialize empty set of clusters  $\mathcal{C} = \{\}$ . As new data points arrive, unit clusters containing individual data points are created. Once a maximum number  $k$  of such clusters have been created, we can begin the process of online cluster maintenance.
2. For a new data point  $X$ , compute its similarity to each cluster  $C_j$ , denoted as  $S(X, C_j)$ .
  - If  $S(X, C_{best}) > \text{thresh}_{outlier}$ , or if there are no inactive clusters left<sup>69</sup>, insert  $X$  to the cluster with maximum similarity.
  - Otherwise, a new cluster is created<sup>70</sup> containing the solitary data point  $X$ .

---

<sup>68</sup>Authors only provide a very brief description, which I'll just copy here:  
*This technique divides up the incoming stream into small segments, each of which can be processed effectively in main memory.*

*A set of k-means iterations are applied to each such data segment in order to cluster them. The advantage of using a segment-wise approach for clustering is that since each segment can be held in main memory, we can process each data point multiple times as long as it is held in main memory. In addition, the centroids from the previous segment are used in the next iteration for clustering purposes. A decay factor is introduced in order to age-out the old documents, so that the new documents are considered more important from a clustering perspective.*

<sup>69</sup>We specify some max allowed number of clusters  $k$ .

<sup>70</sup>The new cluster replaces the least recently updated inactive cluster.

## Misc.

- Mandatory read: reference [61]. Details phrase extraction/**topic signatures**. The use of using phrases instead of individual words is referred to as **semantic smoothing**.
- For *dynamic* (and more recent) topic modeling, see reference [107] of the paper, titled “A probabilistic model for online document clustering with application to novelty detection.”

**Semi-Supervised Clustering.** Useful when we have any prior knowledge about the kinds of clusters available in the underlying data. Some approaches:

- Incorporate this knowledge when seeding the cluster centroids for  $k$ -means clustering.
- Iterative EM approach: unlabeled documents are assigned labels using a naive Bayes approach on the currently labeled documents. These newly labeled documents are then again used for re-training a Bayes classifier. Iterate to convergence.
- Graph-based approach: graph nodes are documents and the edges are similarities between the connected documents (nodes). We can incorporate prior knowledge by adding certain edges between nodes that we know are similar. A *normalized cut algorithm* is then applied to this graph in order to create the final clustering.

We can also use partially supervised methods in conjunction with pre-existing categorical *hierarchies*.

## Deep Sentence Embedding Using LSTMs

[Table of Contents](#)   [Local](#)

*Written by Brandon McKinzie*

Palangi et al., “Deep Sentence Embeddings Using Long Short-Term Memory Networks: Analysis and Application to Information Retrieval,” (2016).

**Abstract.** Sentence embeddings using LSTM cells, which automatically attenuate unimportant words and detect salient keywords. Main emphasis on applications for document retrieval (matching a query to a document<sup>71</sup>).

**Introduction.** Sentence embeddings are learned using a loss function defined on *sentence pairs*. For example, the well-known Paragraph Vector<sup>72</sup> is learned in an unsupervised manner as a distributed representation of sentences and documents, which are then used for sentiment analysis.

The authors appear to use a dataset of their own containing examples of (search-query, clicked-title) for a search engine. Their training objective is to maximize the similarity between the two vectors mapped by the LSTM-RNN from the query and the clicked document, respectively. One very interesting claim to pay close attention to:

*We further show that different cells in the learned model indeed correspond to **different topics**, and the keywords associated with a similar topic activate the same cell unit in the model.*

### Related Work. (Identified by reference number)

- [2] Good for sentiment, but doesn’t capture fine-grained sentence structure.
- [6] Unsupervised embedding method trained on the BookCorpus [7]. Not good for document retrieval task.
- [9] Semi-supervised Recursive Autoencoder (RAE) for sentiment prediction.
- [3] DSSM (uses bag-of-words) and [10] CLSM (uses bag of n-grams) models for IR and also sentence embeddings.
- [12] Dynamic CNN for sentence embeddings. Good for sentiment prediction and question type classification. In [13], a CNN is proposed for **sentence matching**<sup>73</sup>

---

<sup>71</sup>Note that this similar to topic extraction.

<sup>72</sup>Q. V. Le and T. Mikolov, “Distributed representations of sentences and documents.”

<sup>73</sup>Might want to look into this.

**Basic RNN.** The information flow (sequence of operations) is enumerated below.

1. Encode  $t$ th word [of the given sentence] in one-hot vector  $\mathbf{x}(t)$ .
2. Convert  $\mathbf{x}(t)$  to a letter tri-gram vector  $\mathbf{l}(t)$  using fixed hashing operator<sup>74</sup>  $\mathbf{H}$ :

$$\mathbf{l}(t) = \mathbf{H}\mathbf{x}(t) \quad (240)$$

3. Compute the hidden state  $\mathbf{h}(t)$ , which is the sentence embedding for  $t = T$ , the length of the sentence.

$$\mathbf{h}(t) = \tanh(\mathbf{U}\mathbf{l}(t) + \mathbf{W}\mathbf{h}(t-1) + \mathbf{b}) \quad (241)$$

where  $\mathbf{U}$  and  $\mathbf{W}$  are the usual parameter matrices for the input/recurrent paths, respectively.

**LSTM.** With peephole connections that expose the internal cell state  $s$  to the sigmoid computations. I'll rewrite the standard LSTM equations from my textbook notes, but with the modifications for peephole connections:

$$f_i^{(t)} = \sigma \left( b_i^f + \sum_j U_{i,j}^f x_j^{(t)} + \sum_j W_{i,j}^f h_j^{(t-1)} + \sum_j P_{i,j}^f s_j^{(t-1)} \right) \quad (242)$$

$$s_i^{(t)} = f_i^{(t)} \odot s_i^{(t-1)} + g_i^{(t)} \odot \sigma \left( b_i + \sum_j U_{i,j} x_j^{(t)} + \sum_j W_{i,j} h_j^{(t-1)} \right) \quad (243)$$

$$g_i^{(t)} = \sigma \left( b_i^g + \sum_j U_{i,j}^g x_j^{(t)} + \sum_j W_{i,j}^g h_j^{(t-1)} + \sum_j P_{i,j}^g s_j^{(t-1)} \right) \quad (244)$$

$$q_i^{(t)} = \sigma \left( b_i^o + \sum_j U_{i,j}^o x_j^{(t)} + \sum_j W_{i,j}^o h_j^{(t-1)} + \sum_j P_{i,j}^o s_j^{(t)} \right) \quad (245)$$

The final hidden state can then be computed via

$$h_i^{(t)} = \tanh(s_i^{(t)}) \odot q_i^{(t)} \quad (246)$$

---

<sup>74</sup>Details aside, the hashing operator serves to lower the dimensionality of the inputs a bit. In particular we use it to convert one-hot word vectors into their letter tri-grams. For example, the word “good” gets surrounded by hashes, ‘#good#’, and then hashed from the one-hot vector to vectorized tri-grams, “#go”, “goo”, “ood”, “od#”.

**Learning method.** We want to maximize the likelihood of the clicked document given query, which can be formulated as the following optimization problem:

$$L(\boldsymbol{\Lambda}) = \min_{\boldsymbol{\Lambda}} \left\{ -\log \prod_{r=1}^N \Pr \left[ D_r^+ \mid Q_r \right] \right\} = \min_{\boldsymbol{\Lambda}} \sum_{r=1}^N l_r(\boldsymbol{\Lambda}) \quad (247)$$

$$l_r(\boldsymbol{\Lambda}) = \log \left( 1 + \sum_{j=1}^n e^{-\gamma \cdot \Delta_{r,j}} \right) \quad (248)$$

where

- $N$  is the number of (query, clicked-doc) pairs in the corpus, while  $n$  is the number of negative samples used during training.
- $D_r^+$  is the clicked document for  $r$ th query.
- $\Delta_{r,j} = R(Q_r, D_r^+) - R(Q_r, D_{r,j}^-)$  ( $R$  is just cosine similarity)<sup>75</sup>.
- $\boldsymbol{\Lambda}$  is all the parameter matrices (and biases) in the LSTM.

The authors then describe standard BPTT updates with momentum, which need not be detailed here. See the “Algorithm 1” figure in the paper for extremely detailed pseudo-code of the training procedure.

---

<sup>75</sup> Note that  $\Delta_{r,j} \in [-2, 2]$ . We use  $\gamma$  as a scaling factor so as to expand this range.

## Clustering Massive Text Streams

Table of Contents Local

Written by Brandon McKinzie

Aggarwal et al., “A Framework for Clustering Massive Text and Categorical Data Streams,” (2006).

**Overview.** Authors present an online approach for clustering massive text and categorical data streams with the use of a statistical summarization methodology. First, we will go over the process of storing and maintaining the data structures necessary for the clustering algorithm. Then, we will discuss the differences which arise from using different kinds of data, and the empirical results.

**Maintaining Cluster Statistics.** The data stream consists of  $d$ -dimensional records, where each dimension corresponds to the numeric frequency of a given word in the vector space representation. Each data point is weighted by the **fading function**  $f(t)$ , a non-monotonic decreasing function which decays uniformly with time  $t$ . The authors define the **half-life** of a data point (e.g. a tweet) as:

$$t_0 \text{ s.t. } f(t_0) = \frac{1}{2}f(0) \quad (249)$$

and, similarly, the **decay-rate** as its inverse,  $\lambda = 1/t_0$ . Thus we have  $f(t) = 2^{-\lambda \cdot t}$ .

To achieve greater accuracy in the clustering process, we require a high level of granularity in the underlying data structures. To do this, we will use a process in which condensed clusters of data points are maintained, referred to as **cluster droplets**. We define them differently for the case of text and categorical data, beginning with categorical:

- **Categorical.** A cluster droplet  $\mathcal{D}(t, \mathcal{C})$  for a set of categorical data points  $\mathcal{C}$  at time  $t$  is defined as the tuple:

$$\mathcal{D}(t, \mathcal{C}) \triangleq (\bar{DF}_2, \bar{DF}_1, n, w(t), l) \quad (250)$$

where

- Entry  $k$  of the vector  $\bar{DF}_2$  is the (weighted) number of points in cluster  $\mathcal{C}$  where the  $i$ th dimension had value  $x$  and the  $j$ th dimension had value  $y$ . In other words, all pairwise combinations of values in the categorical vector<sup>76</sup>.  $\sum_{i=1}^d \sum_{j \neq i}^d v_i v_j$  entries total<sup>77</sup>.
- Similarly,  $\bar{DF}_1$  consists of the (weighted) counts that some dimension  $i$  took on the value  $x$ .  $\sum_{i=1}^d v_i$  entries total.

<sup>76</sup>This is intentionally written hand-wavy because I’m really concerned with *text* streams and don’t want to give this much space.

<sup>77</sup> $v_i$  is the number of values the  $i$ th categorical dimension can take on.

- $w(t)$  is the sum of the weights of the data points at time  $t$ .
- $l$  is the time stamp of the last time a data point was added to the cluster.
- **Text.** Can be viewed as an example of a sparse numeric data set. A cluster droplet  $\mathcal{D}(t, \mathcal{C})$  for a set of text data points  $\mathcal{C}$  at time  $t$  is defined as the tuple:

$$\mathcal{D}(t, \mathcal{C}) \triangleq (\bar{DF}2, \bar{DF}1, n, w(t), l) \quad (251)$$

where

- $\bar{DF}2$  contains  $3 \cdot wb \cdot (wb - 1)/2$  entries, where  $wb$  is the number of distinct words in the cluster  $\mathcal{C}$ .
- $\bar{DF}1$  contains  $2 \cdot wb$  entries.
- $n$  is the number of data points in the cluster  $\mathcal{C}$ .

### Cluster Droplet Maintenance.

1. We first start off with  $k$  trivial clusters (the first  $k$  data points that arrived).
2. When a new point  $\bar{X}$  arrives, the cosine similarity to each cluster's  $\bar{DF}1$  is computed.
3.  $\bar{X}$  is inserted into the cluster for which this is a maximum, so long as the associated  $S(\bar{X}, \bar{DF}1) > \text{thresh}$ , a predefined threshold. If not above the threshold *and* some **inactive cluster** exists, a new cluster is created containing the solitary point  $\bar{X}$ , which replaces the inactive cluster. If not above threshold but no inactive clusters, then we just insert it into the max similarity cluster anyway.
4. If  $\bar{X}$  was inserted (i.e. didn't replace an inactive cluster), then we need to:
  - (a) Update the statistics to reflect the decay of the data points at the current moment in time<sup>78</sup>. This is done by multiplying entries in the droplet vectors by  $2^{-\lambda \cdot (t-l)}$ .
  - (b) Add the statistics for each newly arriving data point to the cluster statistics.

---

<sup>78</sup>In other words, the statistics for a cluster do not decay, until a new point is added to it.

## Supervised Universal Sentence Representations (InferSent)

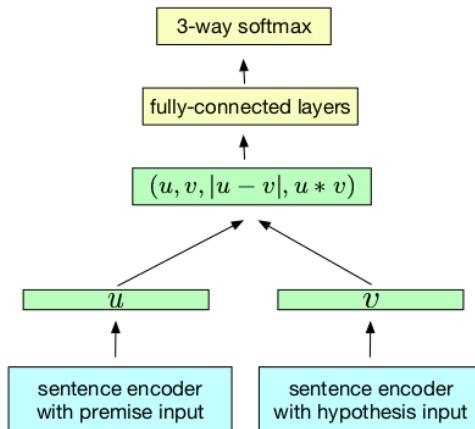
Table of Contents   Local

*Written by Brandon McKinzie*

Conneau et al., “Supervised Learning of Universal Sentence Representations from Natural Language Inference Data,” Facebook AI Research (2017).

**Overview.** Authors claim universal sentence representations trained using the supervised data of the Stanford Natural Language Inference (SNLI) dataset can consistently outperform unsupervised methods like SkipThought on a wide range of transfer tasks. They emphasize that training on NLI tasks in particular results in embeddings that perform well in transfer tasks. Their best encoder is a Bi-LSTM architecture with max pooling, which they claim is SOTA when trained on the SNLI data.

**The Natural Language Inference Task.** Also known as *Recognizing Textual Entailment* (RTE). The SNLI data consists of sentence pairs labeled as one of entailment, contradiction, or neutral. Below is a typical architecture for training on SNLI.



Note that the same sentence encoder is used for both  $u$  and  $v$ . To obtain a sentence vector from a BiLSTM encoder, they experiment with (1) the average  $h_t$  over all  $t$  (**mean pooling**), and (2) selecting the max value over each dimension of the hidden units [over all timesteps] (**max pooling**)<sup>79</sup>.

<sup>79</sup>Since the authors have already mentioned that BiLSTM did the best, I won’t go over the other architectures they tried: self-attentive networks, hierarchical convnet, vanilla LSTM/GRU.

## Dist. Rep. of Sentences from Unlabeled Data (FastSent)

Table of Contents Local

Written by Brandon McKinzie

Hill et al., “Learning Distributed Representations of Sentences from Unlabelled Data,” (2016).

**Overview.** A systematic comparison of models that learn distributed representations of phrases/sentences from unlabeled data. Deeper, more complex models are preferable for representations to be used in supervised systems, but shallow log-linear models work best for building representation spaces that can be decoded with simple spatial distance metrics.

Authors propose two new phrase/sentence representation learning objectives:

1. **Sequential Denoising Autoencoders** (SDAEs)
2. **FastSent**: a sentence-level log-linear BOW model.

**Distributed Sentence Representations.** Existing models trained on text:

- **SkipThought Vectors** (Kiros et al., 2015). Predict target sentences  $S_{i\pm 1}$  given source sentence  $S_i$ . Sequence-to-sequence model.
- **ParagraphVector** (Le and Mikolov, 2014). Defines 2 log-linear models:
  1. **DBOW**: learns a vector  $s$  for every sentence  $S$  in the training corpus which, together with word embeddings  $v_w$ , define a softmax distribution to predict words  $w \in S$  given  $S$ .
  2. **DM**: select k-grams of consecutive words  $\{w_i \dots w_{i+k} \in S\}$  and the sentence vector  $s$  to predict  $w_{i+k+1}$ .
- **Bottom-Up Methods.** Train **CBOW** and **Skip-Gram** word embeddings on the Books corpus.

The authors use `gensim` to implement ParagraphVector.

Models trained on *structured* (and freely-available) resources:

- **DictRep** (Hill et al., 2015a). Map dictionary definitions to pre-trained word embeddings, using either BOW or RNN-LSTM encoding.
- **NMT**. Consider sentence representations learned by sequence-to-sequence NMT models.

## Novel Text-Based Methods.

- **Sequential (Denoising) Autoencoders.** To avoid needing coherent inter-sentence narrative, try this representation-learning objective based on DAEs. For a given sentence  $S$  and **noise function**  $N(S | p_o, p_x)$  (where  $p_o, p_x \in [0, 1]$ ), the approach is as follows:

1. For each  $w \in S$ ,  $N$  deletes  $w$  with probability  $p_o$ .
2. For each non-overlapping bigram  $w_i w_{i+1} \in S$ ,  $N$  swaps  $w_i$  and  $w_{i+1}$  with probability  $p_x$ .

*We then train the same LSTM-based encoder-decoder architecture as NMT, but with the denoising objective to predict (as target) the original source sentence  $S$  given a corrupted version  $N(S | p_o, p_x)$  (as source).*

Authors recommend  
 $p_o = p_x = 0.1$

- **FastSent.** Designed to be a more efficient/quicker to train version of SkipThought.

## Latent Dirichlet Allocation

Table of Contents Local

Written by Brandon McKinzie

Blei et al., “Latent Dirichlet Allocation,” (2003).

**Introduction.** At minimum, one should be familiar with generative probabilistic models, mixture models, and the notion of latent variables before continuing. The “Dirichlet” in LDA of course refers to the **Dirichlet distribution**, which is a generalization of the beta distribution,  $B$ . It’s PDF is defined as<sup>8081</sup>:

$$\text{Dir}(\boldsymbol{x}; \boldsymbol{\alpha}) = \frac{1}{B(\boldsymbol{\alpha})} \prod_{i=1}^K \boldsymbol{x}_i^{\boldsymbol{\alpha}_i - 1} \quad \text{where} \quad B(\boldsymbol{\alpha}) = \frac{\prod_{i=1}^K \Gamma(\boldsymbol{\alpha}_i)}{\Gamma(\sum_{i=1}^K \boldsymbol{\alpha}_i)} \quad (252) \quad \sum_{i=1}^K \boldsymbol{x}_i = 1 \\ (\forall i \in [1, K]) : \boldsymbol{x}_i \geq 0$$

Main things to remember about LDA:

- Generative probabilistic model for collections of discrete data such as text corpora.
- Three-level **hierarchical Bayesian model**. Each document is a mixture of topics, each topic is an infinite mixture over a set of topic probabilities.

Condensed comparisons/history of related models leading up to LDA:

- **TF-IDF.** Design matrix  $\mathbf{X} \in \mathbb{R}^{V \times M}$ , where  $M$  is the number of docs, and  $\mathbf{X}_{i,j}$  gives the TF-IDF value for  $i$ th word in vocabulary and corresp. to document  $j$ .
- **LSI:**<sup>82</sup> Performs SVD on the TF-IDF design matrix  $\mathbf{X}$  to identify a linear subspace in the space of tf-idf features that captures most of the variance in the collection.
- **pLSI: TODO**

*pLSI is incomplete in that it provides no probabilistic model at the level of documents. In pLSI, each document is represented as a list of numbers (the mixing proportions for topics), and there is no generative probabilistic model for these numbers.*

---

<sup>80</sup> Recall that for positive integers  $n$ ,  $\Gamma(n) = (n - 1)!$ .

<sup>81</sup>The Dirichlet distribution is **conjugate** to the multinomial distribution. TODO: Review how to interpret this.

<sup>82</sup>Recall that LSI is basically PCA but without subtracting off the means

**Model.** LDA assumes the following generative process for each document (word sequence)  $\mathbf{w}$ :

1.  $N \sim \text{Poisson}(\lambda)$ : Sample  $N$ , the number of words (length of  $\mathbf{w}$ ), from  $\text{Poisson}(\lambda) = e^{-\lambda} \frac{\lambda^n}{n!}$ . The parameter  $\lambda$  *should* represent the average number of words per document.
2.  $\theta \sim \text{Dir}(\alpha)$ : Sample  $k$ -dimensional vector  $\theta$  from the Dirichlet distribution (eq. 252),  $\text{Dir}(\alpha)$ .  $k$  is the number of topics (pre-defined by us). Recall that this means  $\theta$  lies in the  $(k-1)$  simplex. The Dirichlet distribution thus tells us the probability density of  $\theta$  over this simplex – it defines the probability of  $\theta$  being at a given position on the simplex.
3. Do the following  $N$  times to generate the words for this document.
  - (a)  $z_n \sim \text{Multinomial}(\theta)$ . Sample a topic  $z_n$ .
  - (b)  $w_n \sim \Pr [w_n | z_n, \beta]$ : Sample a word  $w_n$  from  $\Pr [w_n | z_n, \beta]$ , a “multinomial probability conditioned on topic  $z_n$ .<sup>83</sup> The parameter  $\beta$  gives the distribution of words given a topic:

$$\beta_{ij} = \Pr [w_j | z_i] \quad (253)$$

In other words, we really sample  $w_n \sim \beta_{i,:}$

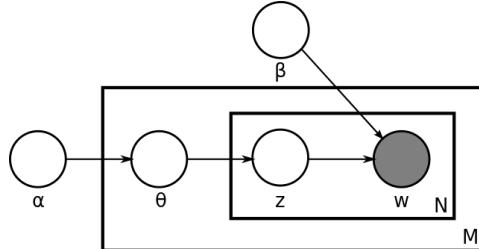
The defining equations for LDA are thus:

$$\Pr [\theta, \mathbf{z}, \mathbf{w} | \alpha, \beta] = \Pr [\theta | \alpha] \prod_{n=1}^N \Pr [z_n | \theta] \Pr [w_n | z_n, \beta] \quad (254)$$

$$\Pr [\mathbf{w} | \alpha, \beta] = \int \Pr [\theta' | \alpha] \left( \prod_{n=1}^N \sum_{z'_n} \Pr [z'_n | \theta'] \Pr [w_n | z'_n, \beta] \right) d\theta' \quad (255)$$

$$\Pr [\mathcal{D} = \{\mathbf{w}^{(1)}, \dots, \mathbf{w}^{(M)}\} | \alpha, \beta] = \prod_{d=1}^M \Pr [\mathbf{w}^{(d)} | \alpha, \beta] \quad (256)$$

Below is the plate notation for LDA, followed by an interpretation:



- **Outermost Variables:**  $\alpha$  and  $\beta$ . Both represent a (Dirichlet) prior distribution:  $\alpha$  parameterizes the probability of a given *topic*, while  $\beta$  a given *word*.
- **Document Plate.**  $M$  is the number of documents,  $\theta_m$  gives the true distribution of topics for document  $m$ <sup>84</sup>.

<sup>83</sup>**TODO:** interpret meaning of the multinomial distributions here. Seems a bit different than standard interp...

<sup>84</sup>In other words, the meaning of  $\theta_{m,i} = x$  is “x percent of document  $m$  is about topic  $i$ .”

- Topic/Word Place.  $z_{mn}$  is the topic for word  $n$  in doc  $m$ , and  $w_{mn}$  is the word. It is shaded gray to indicate it is the only **observed variable**, while all others are **latent variables**.

**Theory.** I'll quickly summarize and interpret the main theoretical points. Without having read all the details, this won't be of much use (i.e. it is for someone who has read the paper already).

- **LDA and Exchangeability.** We assume that each document is a bag of words (order doesn't matter; frequency still does) *and* a bag of topics. In other words, a document of  $N$  words *is* an unordered list of words and topics. De Finetti's theorem tells us that we can model the joint probability of the words and topics as if a random parameter  $\theta$  were drawn from some distribution and then the variables within  $w, z$  were **conditionally independent given  $\theta$** . LDA posits that a good distribution to sample  $\theta$  from is a Dirichlet distribution.
- **Geometric Interpretation:** **TODO**

**Inference and Parameter Estimation.** As usual, we need to find a way to compute the posterior distribution of the hidden variables given a document  $w$ :

$$\Pr [\theta, z | w, \alpha, \beta] = \frac{\Pr [\theta, z, w | \alpha, \beta]}{\Pr [w | \alpha, \beta]} \quad (257)$$

Computing the denominator exactly is intractable. Common approximate inference algorithms for LDA include Laplace approximation, variational approximation, and Markov Chain Monte Carlo.

## Conditional Random Fields

Table of Contents Local

Written by Brandon McKinzie

Lafferty et al., “Conditional Random Fields: Probabilistic Models for Segmenting and Labeling Sequence Data,” (2001).

**Introduction.** CRFs offer improvements to HMMs, MEMMs, and other discriminative Markov models. MEMMs and other non-generative models share a weakness called the **label bias problem**: the transitions leaving a given state compete only against each other, rather than against all other transitions in the model. The key difference between CRFs and MEMMs is that a CRF has a single exponential model for the joint probability of the entire sequence of labels given the observation sequence.

**The Label Bias Problem.** Recall that MEMMs are run left-to-right. One way of interpreting such a model is to consider how the probabilities (of state sequences) are distributed as we continue through the sequence of observations. The issue with MEMMs is that there’s nothing we can do if, somewhere along the way, we observe something that makes one of these state paths extremely likely/unlikely; we can’t redistribute the probability mass amongst the various allowed paths. The CRF solution:

*Account for whole state sequences at once by letting some transitions “vote” more strongly than others depending on the corresponding observations. This implies that score mass will not be conserved, but instead individual transitions can “amplify” or “dampen” the mass they receive.*

**Conditional Random Fields.** Here we formalize the model and notation. Let  $\mathbf{X}$  be a random variable over data sequences to be labeled (e.g. over all words/sentences), and let  $\mathbf{Y}$  the random variable over corresponding label sequences<sup>85</sup>. Formal definition:

*Let  $G = (V, E)$  be a graph such that  $Y = (Y_v)_{v \in V}$ , so that  $Y$  is indexed by the vertices of  $G$ . Then  $(X, Y)$  is a CRF if, when conditioned on  $X$ , the random variables  $Y_v$  obey the Markov property with respect to the graph:*

$$\Pr[Y_v | X, Y_w, w \neq v] = \Pr[Y_v | X, Y_w, w \sim v] \quad (258)$$

*where  $w \sim v$  means that  $w$  and  $v$  are neighbors in  $G$ .*

All this means is a CRF is a random field (discrete set of random-valued points in a space) where all points (i.e. globally) are conditioned on  $\mathbf{X}$ . If the graph  $G = (V, E)$  of  $\mathbf{Y}$  is a tree, its cliques<sup>86</sup> are the edges and vertices. Take note that  $\mathbf{X}$  is not a member of the vertices

<sup>85</sup>We assume all components  $\mathbf{Y}_i$  can only take on values in some finite label set  $\mathcal{Y}$ .

<sup>86</sup>A clique is a subset of vertices in an undirected graph such that every two distinct vertices in the clique are adjacent

in  $G$ .  $G$  only contains vertices corresponding to elements of  $\mathbf{Y}$ . Accordingly, when the authors refer to cases where  $G$  is a “chain”, remember that they just mean the  $\mathbf{Y}$  vertex sequence.

By the fundamental theorem of random fields:

$$p_\theta(\mathbf{y} \mid \mathbf{x}) \propto \exp \left( \sum_{e \in E, k} \lambda_k f_k(e, \mathbf{y}|_e, \mathbf{x}) + \sum_{v \in V, k} \mu_k g_k(v, \mathbf{y}|_v, \mathbf{x}) \right) \quad (259)$$

where  $\mathbf{y}|_S$  is the set of components of  $\mathbf{y}$  associated with the vertices in subgraph  $S$ . We assume the  $K$  feature [functions]  $f_k$  and  $g_k$  are given and fixed. Note that  $f_k$  are the feature functions over *transitions*  $y_{t-1}$  to  $y_t$ , and  $g_k$  are the feature functions over *states*  $y_t$  and  $x_t$ . Our estimation problem is thus to determine parameters  $\theta = (\lambda_1, \lambda_2, \dots; \mu_1, \mu_2, \dots)$  from the labeled training data.

**Linear-Chain CRF.** Let  $|\mathcal{Y}|$  denote the number of possible labels. At each position  $t$  in the observation sequence  $\mathbf{x}$ , we define the  $|\mathcal{Y}| \times |\mathcal{Y}|$  matrix random variable  $\mathbf{M}_t(\mathbf{x})$

$$\mathbf{M}_t(y', y, \mid \mathbf{x}) = \exp(\Lambda_t(y', y \mid x)) \quad (260)$$

$$\Lambda_t(y', y \mid x) = \sum_k \lambda_k f_k(y', y, \mathbf{x}) + \sum_k \mu_k g_k(y, \mathbf{x}) \quad (261)$$

where  $y_{t-1} := y'$  and  $y_t := y$ . We can see that the individual elements correspond to specific values of  $e$  and  $v$  in the double-summations of  $p_\theta(\mathbf{y} \mid \mathbf{x})$  above. Then the normalization (partition function)  $Z_\theta(\mathbf{x})$  is the  $(y_0, y_{T+1})$  entry (the fixed boundary states) of the product:

$$Z_\theta(\mathbf{x}) = \left[ \prod_{t=1}^{T+1} \mathbf{M}_t(\mathbf{x}) \right]_{y_0, y_{T+1}} \quad (262)$$

which includes all possible sequences  $\mathbf{y}$  that start with the fixed  $y_0$  and end with the fixed  $y_{T+1}$ . Now we can write the conditional probability as a function of just these matrices:

$$p_\theta(\mathbf{y} \mid \mathbf{x}) = \frac{\prod_{t=1}^{T+1} \mathbf{M}_t(y_{t-1}, y_t \mid \mathbf{x})}{\left[ \prod_{t=1}^{T+1} \mathbf{M}_t(\mathbf{x}) \right]_{y_0, y_{T+1}}} \quad (263)$$

**Parameter Estimation** (for linear-chain CRFs). For each  $t$  in  $[0, T + 1]$ , define the **forward vectors**  $\alpha_t(\mathbf{x})$  with base case  $\alpha_0(y \mid \mathbf{x}) = 1$  if  $y = y_0$ , else 0. Similarly, define the **backward vectors**  $\beta_t(\mathbf{x})$  with base case  $\beta_{T+1}(y \mid \mathbf{x}) = 1$  if  $y = y_{T+1}$  else 0<sup>87</sup>. Their recurrence relations are

$$\alpha_t(\mathbf{x}) = \alpha_{t-1}(\mathbf{x}) \mathbf{M}_t(\mathbf{x}) \quad (264)$$

$$\beta_t(\mathbf{x})^T = \mathbf{M}_{t+1}(\mathbf{x}) \beta_{t+1}(\mathbf{x}) \quad (265)$$

---

<sup>87</sup>Remember that  $y_0$  and  $y_{T+1}$  are their own fixed symbolic constants representing a fixed start/stop state.

## Attention Is All You Need

Table of Contents Local

Written by Brandon McKinzie

Vaswani et al., “Attention Is All You Need,” (2017)

**Overview.** Authors refer to sequence *transduction* models a lot – just a fancy way of referring to models that transform input sequences into output sequences. Authors propose new architecture, the **Transformer**, based solely on attention mechanisms (no recurrence!).

### Model Architecture.

- **Encoder.**  $N=6$  identical layers, each with 2 sublayers: (1) a **multi-head self-attention** mechanism and (2) a position-wise FC feed-forward network. They apply a residual connection and layer norm such that each sublayer, instead of outputting Sublayer(x), instead outputs LayerNorm(x + Sublayer(x)).
- **Decoder.**  $N=6$  with 3 sublayers each. In addition to the two sublayers described for the encoder, the decoder has a third sublayer, which performs **multi-head** attention over the output of the encoder stack. Same residual connections and layer norm.

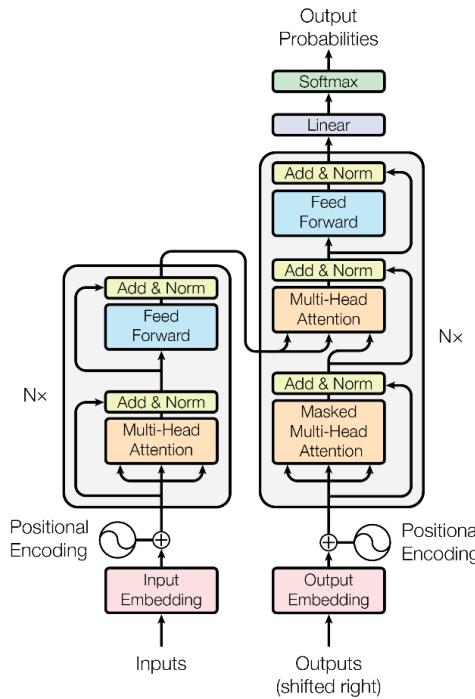


Figure shows encoder-decoder template layers. The actual model instantiates chain of 6 encoder layers and 6 decoders layers. The decoder's self-attention masks embeddings at future timesteps to zero.

**Attention.** An attention function can be described as a mapping:

$$\text{Attn}(\text{query}, \{(k_1, v_1), \dots\}) \Rightarrow \sum_i f_n(\text{query}, k_i) v_i \quad (266)$$

where the query, keys, values, and output are all vectors.

- **Scaled Dot-Product Attention.**

1. **Inputs:** queries  $q$ , keys  $k$  of dimension  $d_k$ , values  $v$  of dimension  $d_v$
2. **Dot Products:** Compute  $\forall k : (q \cdot k) / \sqrt{d_k}$ .
3. **Softmax:** on each dot product above. This gives the weights on the values shown earlier.

Appears that  $d_q \equiv d_k$ .

In practice, this is done simultaneously for all queries in a set via the following matrix equation:

$$\text{Attention}(Q, K, V) = \text{softmax} \left( \frac{QK^T}{\sqrt{d_k}} \right) V \quad (267)$$

Note that this is identical to the standard dot-product attention mechanism, except for the *scaling* factor (hence the name) of  $1/\sqrt{d_k}$ . The scaling factor is motivated by the fact that additive attention outperforms dot-product attention for large  $d_k$  and the authors stipulate this is due to the softmax having small gradients in this case, due to the large dot products<sup>88</sup>.

First, let's explicitly show which indices are being normalized over, since it can get confusing when presented with the highly vectorized version above. For a given input sequence of length  $T$ , and for the self-attention version where  $K=Q=V \in \mathbb{R}^{T \times d_k}$ , the output attention vector for timestep  $t$  is explicitly (ignoring the  $\sqrt{d_k}$  for simplicity)

$$\text{Attention}(Q, K, V)_t = \left[ \text{softmax} \left( QK^T \right) V \right]_t \quad (271)$$

$$= \sum_{t'}^T \frac{e^{Q_t \cdot K_{t'}}}{\sum_{t''}^T e^{Q_t \cdot K_{t''}}} V_{t'} \quad (272)$$

---

<sup>88</sup> Assume that  $\mathbf{q}$  and  $\mathbf{k}$  are vectors in  $\mathbb{R}^d$  whose components are independent RVs with  $\mathbb{E}[q_i] = \mathbb{E}[k_j] = 0$  ( $\forall i, j$ ), and  $\text{Var}[q_i] = \text{Var}[k_j] = 1$  ( $\forall i, j$ ). Then

$$\mathbb{E}[\mathbf{q} \cdot \mathbf{k}] = \mathbb{E} \left[ \sum_i^d q_i k_i \right] = \sum_i^d \mathbb{E}[q_i k_i] = \sum_i^d \mathbb{E}[q_i] \mathbb{E}[k_i] = 0 \quad (268)$$

$$\text{Var}[\mathbf{q} \cdot \mathbf{k}] = \text{Var} \left[ \sum_i^d q_i k_i \right] = \sum_i^d \text{Var}[q_i k_i] = \sum_i^d \mathbb{E}[q_i^2 k_i^2] - \mathbb{E}[q_i k_i]^2 \quad (269)$$

$$= \sum_i^d \mathbb{E}[q_i^2] \mathbb{E}[k_i^2] = \sum_i^d \text{Var}[q_i] \text{Var}[k_i] = d \quad (270)$$

See this S.O answer and/or these useful formulas for more details.

Next, the gradient of the  $d$ th softmax output w.r.t its inputs is

$$\frac{\partial \text{Softmax}_d(\mathbf{x})}{\partial x_j} = \text{Softmax}_d(\mathbf{x}) (\delta_{dj} - \text{Softmax}_d(\mathbf{x})) \quad (273)$$

- **Multi-Head Attention.** Basically just doing some number  $h$  of parallel attention computations. Before each of these, the queries, keys, and values are linearly projected with different, learned linear projections to  $d_k$ ,  $d_k$  and  $d_v$  dimensions respectively (and then fed to their respective attention function). The  $h$  outputs are then concatenated and once again projected, resulting in the final values.

$$\text{MultiHead}(Q, K, V) = \text{Concat}(\text{head}_1, \dots, \text{head}_h)W^O \quad (274)$$

$$\text{where } \text{head}_i = \text{Attention}(QW_i^Q, KW_i^K, VW_i^V) \quad (275)$$

The authors employ  $h = 8$ ,  $d_k = d_v = d_{model}/h = 64$ .

$$W_i^Q \in \mathbb{R}^{d_{model} \times d_k}$$

$$W_i^K \in \mathbb{R}^{d_{model} \times d_k}$$

$$W_i^V \in \mathbb{R}^{d_{model} \times d_v}$$

$$W^O \in \mathbb{R}^{hd_v \times d_{model}}$$

The Transformer uses multi-headed attention in 3 ways:

1. **Encoder-decoder attention:** the normal kind. Queries are previous decoder layer, and memory keys and values come from output of the [final layer of] encoder.
2. **Encoder self-attention:** all of the keys, values, and queries come from the previous *layer* in the encoder. Each position in the encoder can attend to all positions in the previous layer of the encoder.
3. **Decoder self-attention:** Similarly, self-attention layers in the decoder allow each position in the decoder to attend to all positions in the decoder up to and including that position (timestep). The masking is done on the inputs to the softmax, setting all inputs beyond the current timestep to  $-\infty$ .

## Other Components.

- **Position-wise Feed-Forward Networks (FFN):** each layer of the encoder and decoder contains a FC FFN, applied to each position separately and identically:

$$\text{FFN}(x) = \max(0, xW_1 + b_1)W_2 + b_2 \quad (276)$$

The FFN is linear → ReLU → linear.

- **Embeddings and Softmax:** use learned embeddings to convert input/output tokens to vectors of dimension  $d_{model}$ , and for the pre-softmax layer at the output of the decoder<sup>89</sup>.
- **Positional Encoding:** how the authors deal with the lack of recurrence (to make use of the sequence order). They add a sinusoid function of the position (timestep)  $pos$  and vector index  $i$  to the input embeddings for the encoder and decoder<sup>90</sup>:

$$\text{PE}(pos, 2i) = \sin\left(\frac{pos}{10000^{2i/d_{model}}}\right) \quad (277)$$

$$\text{PE}(pos, 2i + 1) = \cos\left(\frac{pos}{10000^{2i/d_{model}}}\right) \quad (278)$$

For inputs to encoder/decoder, the embedding weights are multiplied by  $\sqrt{d_{model}}$

The authors justify this choice:

<sup>89</sup>In other words, they use the same weight matrix for all three of (1) encoder input embedding, (2) decoder input embedding, and (3) (opposite direction) from decoder output to pre-softmax.

<sup>90</sup>Note that the positional encodings must necessarily be of dimension  $d_{model}$  to be summed with the input embeddings.

*We chose this function because we hypothesized it would allow the model to easily learn to attend by relative positions, since for any fixed offset  $k$ ,  $PE_{pos+k}$  can be represented as a linear function of  $PE_{pos}$ .*

**Summary of Add-ons.** Below is a list of all the little bells and whistles they add to the main components of the model that are easy to miss since they mention them throughout the paper in a rather unorganized fashion.

- Shared weights for encoder inputs, decoder inputs, and final softmax projection outputs.
- Multiply the encoder and decoder input embedding [shared] weights by  $\sqrt{d_{model}}$ . **TODO:** why? Also this must be highly correlated with their decision regarding weight initialization (mean/stddev/technique). Add whatever they use here if they mention it.
- Adam optimizer with  $\beta_1=0.9$ ,  $\beta_2=0.98$ ,  $\epsilon=10^{-9}$ .
- Learning rate schedule  $LR(s) = d_{model}^{-0.5} \cdot \min(s^{-0.5}, s \cdot w^{-1.5})$  for global step  $s$  and warmup steps  $w=4000$ .
- Dropout on sublayer outputs pre-layernorm-and-residual. Specifically, they *actually* return  $\text{LayerNorm}(x + \text{Dropout}(\text{Sublayer}(x)))$ . Use  $P_{drop} = 0.1$ .
- Dropout the summed embeddings+positional-encodings for both encoder and decoder stacks.
- Dropout on softmax outputs. So do  $\text{Dropout}(\text{Softmax}(QK))V$ .
- Label smoothing with  $\epsilon_{ls} = 0.1$ .

## Hierarchical Attention Networks

Table of Contents Local

Written by Brandon McKinzie

Yang et al., "Hierarchical Attention Networks for Document Classification."

**Overview.** Authors introduce the Hierarchical Attention Network (HAN) that is designed to capture insights regarding (1) the hierarchical structure of documents (words -> sentences -> documents), and (2) the context dependence between words and sentences. The latter is implemented by including two levels of attention mechanisms, one at the word level and one at the sentence level.

**Hierarchical Attention Networks.** Below is an illustration of the network. The first stage is familiar to sequence to sequence models - a bidirectional encoder for outputting sentence-level representations of a sequence of words. The HAN goes a step further by feeding this another bidirectional encoder for outputting document-level representations for sequences of sentences.

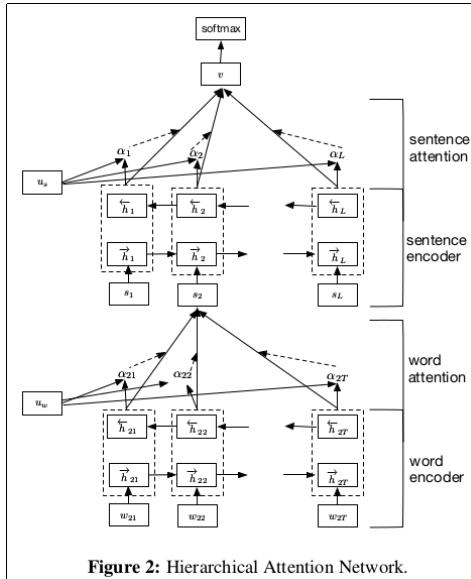


Figure 2: Hierarchical Attention Network.

The authors choose the GRU as their underlying RNN. For ease of reference, the defining equations of the GRU are shown below:

$$h_t = (1 - z_t) \odot h_{t-1} + z_t \odot \tilde{h}_t \quad (279)$$

$$z_t = \sigma(W_z x_t + U_z h_{t-1} + b_z) \quad (280)$$

$$\tilde{h}_t = \tanh(W_h x_t + r_t \odot (U_h h_{t-1}) + b_h) \quad (281)$$

$$r_t = \sigma(W_r x_t + U_r h_{t-1} + b_r) \quad (282)$$

**Hierarchical Attention.** Here I'll overview the main stages of information flow.

1. **Word Encoder.** Let the  $t$ th word in the  $i$ th sentence be denoted  $w_{it}$ . They embed the vectors with a word embedding matrix  $W_e$ ,  $x_{it} = W_e w_{it}$ , and then feed  $x_{it}$  through a bidirectional GRU to ultimately obtain  $h_{it} := [\overrightarrow{h}_{it}; \overleftarrow{h}_{it}]$ .
2. **Word Attention.** Extracts words that are important to the meaning of the sentence and aggregates the representation of these informative words to form a sentence vector.

$$u_{it} = \tanh(W_w h_{it} + b_w) \quad (283)$$

$$\alpha_{it} = \frac{\exp(u_{it}^T u_w)}{\sum_t \exp(u_{it}^T u_w)} \quad (284)$$

$$s_i = \sum_t \alpha_{it} h_{it} \quad (285)$$

Note the *context vector*  $u_w$ , which is shared for all words<sup>91</sup> and randomly initialized and jointly learned during the training process.

3. **Sentence Encoder.** Similar to the word encoder, but uses the sentence vectors  $s_i$  as the input for the  $i$ th sentence in the document. Note that the output of this encoder,  $h_i$  contains information from the neighboring sentences too (bidirectional) but focuses on sentence  $i$ .
4. **Sentence Attention.** For rewarding sentences that are clues to correctly classify a document. Similar to before, we now use a sentence level context vector  $u_s$  to measure the importance of the sentences.

$$u_i = \tanh(W_s h_i + b_s) \quad (286)$$

$$\alpha_i = \frac{\exp(u_i^T u_s)}{\sum_i \exp(u_i^T u_s)} \quad (287)$$

$$v = \sum_t \alpha_i h_i \quad (288)$$

where  $v$  is the document vector that summarizes all the information of sentences in a document.

As usual, we convert  $v$  to a normalized probability vector by feeding through a softmax:

$$p = \text{softmax}(W_c v + b_c) \quad (289)$$

---

<sup>91</sup>To emphasize, there is only a single context vector  $u_w$  in the network, period. The subscript just tells us that it is the word-level context vector, to distinguish it from the sentence-level context vector in the later stage.

**Configuration and Training.** Quick overview of some parameters chosen by the authors:

- **Tokenization:** Stanford CoreNLP. Vocabulary consists of words occurring more than 5 times, all others are replaced with UNK token.
- **Word Embeddings:** train word2vec on the training and validation splits. Dimension of 200.
- **GRU.** Dimension of 50 (so 100 because bidirectional).
- **Context vectors.** Both  $u_w$  and  $u_s$  have dimension of 100.
- **Training:** batch size of 64, grouping documents of similar length into a batch. SGD with momentum of 0.9.

## Joint Event Extraction via RNNs

Table of Contents Local

Written by Brandon McKinzie

Nguyen, Cho, and Grishman, “Joint Event Extraction via Recurrent Neural Networks,” (2016).

**Event Extraction Task.** Automatic Context Extraction (ACE) evaluation. Terminology:

- **Event:** something that happens or leads to some change of state.
- **Mention:** phrase or sentence in which an event occurs, including one trigger and an arbitrary number of arguments.
- **Trigger:** main word that most clearly expresses an event occurrence.
- **Argument:** an entity mention, temporal expression, or value that serves as a participant/attribute with a specific role in an event mention.

Example:

In Baghdad, a **cameraman** **died**{*Die*} when an American tank **-fired**{*Attack*} on the Palestine hotel.

TRIGGER  
ARGUMENT

Each event subtype has its own set of roles to be filled by the event arguments. For example, the roles for the *Die* event subtype include *Place*, *Victim*, and *Time*.

## Model.

- **Sentence Encoding.** Let  $w_i$  denote the  $i$ th token in a sentence. It is transformed into a real-valued vector  $x_i$ , defined as

$$x_i := [\text{GloVe}(w_i); \text{Embed}(\text{EntityType}(w_i)); \text{DepVec}(w_i)] \quad (290)$$

where “Embed” is an embedding we learn, and “DepVec” is the binary vector whose dimensions correspond to the possible relations between words in the dependency trees.

- **RNN.** Bidirectional LSTM on the inputs  $x_i$ .
- **Prediction.** Binary memory vector  $G_i^{trg}$  for triggers; binary memory matrices  $G_i^{arg}$  and  $G_i^{arg/trg}$  for arguments (at each timestep  $i$ ). At each time step  $i$ , do the following in order:
  1. Predict trigger  $t_i$  for  $w_i$ . First compute the feature representation vector  $R_i^{trig}$ , defined as:

$$R_i^{trig} := [h_i; L_i^{trg}; G_{i-1}^{trg}] \quad (291)$$

where  $h_i$  is the RNN output,  $L_i^{trg}$  is the local context vector for  $w_i$ , and  $G_{i-1}^{trg}$  is the memory vector from the previous step.  $L_i^{trg} := [\text{GloVe}(w_{i-d}); \dots; \text{GloVe}(w_{i+d})]$  for

some predefined window size  $d$ . This is then fed to a fully-connected layer with softmax activation,  $\mathbf{F}^{trg}$ , to compute the probability over possible trigger subtypes:

$$P_{i;t}^{trg} := F_t^{trg}(R_i^{trg}) \quad (292)$$

As usual, the predicted trigger type for  $w_i$  is computed as  $t_i = \arg \max_t (P_{i;t}^{trg})$ . If  $w_i$  is not a trigger,  $t_i$  should predict “Other.”

2. Argument role predictions,  $a_{i1}, \dots, a_{ik}$ , for all of the [already known] entity mentions in the sentence,  $e_1, \dots, e_k$  with respect to  $w_i$ .  $a_{ij}$  denotes the argument role of  $e_j$  with respect to [the predicted trigger of]  $w_i$ . If NOT( $w_i$  is trigger AND  $e_j$  is one of its arguments), then  $a_{ij}$  is set to Other. For example, if  $w_i$  was the word “died” from our example sentence, we’d hope that its predicted trigger would be  $t_i = Die$ , and that **the entity associated with “cameraman” would get a predicted argument role of Victim.**

---

```
def getArguments(triggerType=t, entities=e):
    k = len(e)
    if isOther(t):
        return [Other] * k
    else:
        for e_j in e:
```

---

$$R_{ij}^{arg} := [h_i; h_{i_j}; L_{ij}^{arg}; B_{ij}; G_{i-1}^{arg}[j]; G_{i-1}^{arg/trg}[j]] \quad (293)$$

3. Update memory. TO BE CONTINUED... (moving onto another paper because this model is getting a *bit* too contrived for my tastes. Also not a fan of the reliance on a dependency parse.)

## Event Extraction via Bidi-LSTM Tensor NNs

Table of Contents Local

Written by Brandon McKinzie

Y. Chen, S. Liu, S. He, K. Liu, and J. Zhao, “Event Extraction via Bidirectional Long Short-Term Memory Tensor Neural Networks.”

**Overview.** The task/goal is the event extraction task as defined in *Automatic Content Extraction* (ACE). Specifically, given a text document, our goal is to do the following in order *for each sentence*:

1. Identify any event triggers in the sentence.
2. If triggers found, predict their subtype. For example, given the trigger “fired,” we may classify it as having the *Attack* subtype.
3. If triggers found, identify their candidate argument(s). ACE defines an event argument as “an entity mention, temporal expression, or value that is involved in an event.”
4. For each candidate argument, predict its role: “the relationship between an argument to the event in which it participates.”

**Context-aware Word Representation.** Use pre-trained word embeddings for the input word tokens, the predicted trigger, and the candidate argument. Note: *we assume we already have predictions for the event trigger  $t$  and are doing a pass for one of (possibly many) candidate arguments  $a$ .*

1. Embed each word in the sentence with pre-trained embeddings. Denote the embedding for  $i$ th word as  $e(w_i)$ .
2. Feed each  $e(w_i)$  through a bidirectional LSTM. Denote the  $i$ th output of the forward LSTM as  $c_l(w_{i+1})$  and the output of the backward LSTM at the same time step as  $c_r(w_{i-1})$ . As usual, they take the general functional form:

$$c_l(w_i) = \overrightarrow{LSTM}(c_l(w_{i-1}), e(w_{i-1})) \quad (294)$$

$$c_r(w_i) = \overleftarrow{LSTM}(c_r(w_{i+1}), e(w_{i+1})) \quad (295)$$

$$(296)$$

3. Concatenate  $e(w_i)$ ,  $c_l(w_i)$ ,  $c_r(w_i)$  together along with the embedding of the candidate argument  $e(a)$  and predicted trigger  $e(t)$ . Also include the relative distance of  $w_i$  to  $t$  or (??)  $a$ , denoted as  $pi$  for position information, and the embedding of the predicted event type  $pe$  of the trigger. Denote this massive concatenation result as  $x_i$ :

$$x_i := c_l(w_i) \oplus e(w_i) \oplus c_r(w_i) \oplus pi \oplus pe \oplus e(a) \oplus e(t) \quad (297)$$

**Dynamic Multi-Pooling.** This is easiest shown by example. Continue with our example sentence:

In California, **Peterson** was arrested for the **murder** of his wife and unborn son.

where the colors are given for *this specific case where murder is our predicted trigger and we are considering the candidate argument Peterson*<sup>92</sup>. Given our  $n$  outputs from the previous stage,  $y^{(1)} \in \mathbb{R}^{n \times m}$ , where  $n$  is the length of the sentence and  $m$  is the size of that huge concatenation given in equation 297. We split our sentence by trigger and candidate argument, then (confusingly) redefine our notation as

$$y_{1j}^{(1)} \leftarrow \begin{bmatrix} y_{1j}^{(1)} & y_{2j}^{(1)} \end{bmatrix} \quad (298)$$

$$y_{2j}^{(1)} \leftarrow \begin{bmatrix} y_{3j}^{(1)} & \dots & y_{7j}^{(1)} \end{bmatrix} \quad (299)$$

$$y_{3j}^{(1)} \leftarrow \begin{bmatrix} y_{8j}^{(1)} & \dots & y_{nj}^{(1)} \end{bmatrix} \quad (300)$$

Peterson is the 3rd word, and murder is the 8th word.

where it's important to see that, for some  $1 \leq j \leq m$ , each new  $y_{ij}^{(1)}$  is a *vector* of length equal to the number of words in segment  $i$ . Finally, the dynamic multi-pooling layer,  $y^{(2)}$ , can be expressed as

$$y_{i,j}^{(2)} := \max(y_{i,j}^{(1)}) \quad 1 \leq i \leq 3, 1 \leq j \leq m \quad (301)$$

where the max is taken over each of the aforementioned vectors, leaving us with  $3m$  values total. These are concatenated to form  $y^{(2)} \in \mathbb{R}^{3m}$ .

**Output.** To predict of each argument role [for the given argument candidate],  $y^{(2)}$  is fed through a dense softmax layer,

$$O = W_2 y^{(2)} + b_2 \quad (302)$$

where  $W_2 \in \mathbb{R}^{n_1 \times 3m}$  and  $n_1$  is the number of possible argument roles (including "None"). The authors also use dropout on  $y^{(2)}$ .

---

<sup>92</sup>Yes, arrested could be another predicted trigger, but the network considers each possibility at separate times/locations in the architecture.

## Reasoning with Neural Tensor Networks

Table of Contents Local

Written by Brandon McKinzie

Socher et al., “Reasoning with Neural Tensor Networks for Knowledge Base Completion”

**Overview.** Reasoning over relationships between two entities. Goal: predict the likely truth of additional facts based on existing facts in the KB. This paper contributes (1) the new NTN and (2) a new way to represent entities in KBs. Each relation is associated with a distinct model. Inputs to a given relation’s model are pairs of database entities, and the outputs score how likely the pair has the relationship.

**Neural Tensor Networks for Relation Classification.** Let  $e_1, e_2 \in \mathbb{R}^d$  be the vector representations of the two entities, and let  $R$  denote the relation (and thus model) of interest. The NTN computes a score of how likely it is that  $e_1$  and  $e_2$  are related by  $R$  via:

$$g(e_1, R, e_2) = u_R^T \tanh \left( e_1^T W_R^{[1:k]} e_2 + V_R \begin{bmatrix} e_1 \\ e_2 \end{bmatrix} + b_R \right) \quad (303)$$

$$W_R^{[1:k]} \in \mathbb{R}^{d \times d \times k}$$

$$V_R \in \mathbb{R}^{k \times 2d}$$

where the bilinear tensor product  $e_1^T W_R^{[1:k]} e_2$  results in a vector  $h \in \mathbb{R}^k$  with each entry computed by one slice  $i = 1, \dots, k$  of the tensor.

*Intuitively, we can see each slice of the tensor as being responsible for one type of entity pair or instantiation of a relation... Another way to interpret each tensor slice is that it mediates the relationship between the two entity vectors differently.*

**Training Objective and Derivatives.** All models are trained with **contrastive max-margin objective functions** and minimize the following objective:

$$J(\Omega) = \sum_{i=1}^N \sum_{c=1}^C \max \left( 0, 1 - g(T^{(i)}) + g(T_c^{(i)}) \right) + \lambda \|\Omega\|_2^2 \quad (304)$$

where  $c$  is for “corrupted” samples,  $T_c^{(i)} := (e_1^{(i)}, R^{(i)}, e_c^{(i)})$ . Notice that this function is minimized when the difference,  $g(T^{(i)}) - g(T_c^{(i)})$ , is maximized. The authors used minibatched **L-BFGS** for optimization.

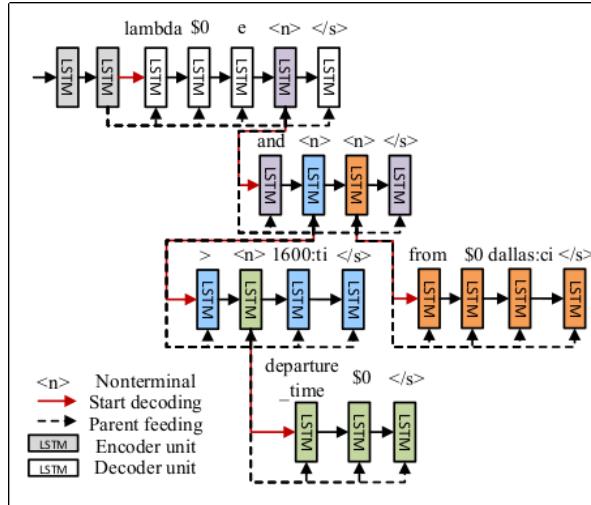
## Language to Logical Form with Neural Attention

Table of Contents Local

Written by Brandon McKinzie

Dong and Lapata, “Language to Logical Form with Neural Attention,” (2016)

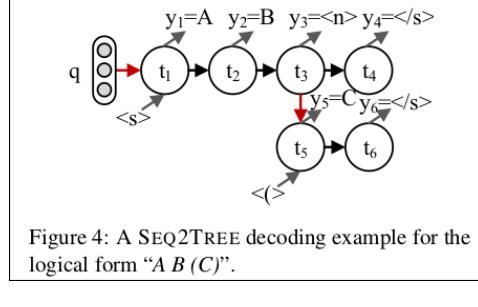
**Sequence-to-Tree Model.** Variant of Seq2Seq that is more faithful to the compositional nature of meaning representations. Its schematic is shown below. The authors define a “nonterminal”  $< n >$  token which indicates [the root of] a subtree.



where the author's have employed “parent-feeding”: for a given subtree (logical form), at each timestep, the hidden vector of the parent nonterminal is concatenated with the inputs and fed into the LSTM (best understood via above illustration).

*After encoding input  $q$ , the hierarchical tree decoder generates tokens at depth 1 of the subtree corresponding to parts of logical form  $a$ . If the predicted token is  $< n >$ , decode the sequence by conditioning on the nonterminal's hidden vector. This process terminates when no more nonterminals are emitted.*

Also note that the output posterior probability over the encoded input  $q$  is the product of subtree posteriors. For example, consider the decoding example in the figure below:



We would compute the output posterior as:

$$p(a | q) = p(y_1 y_2 y_3 y_4 | q) p(y_5 y_6 | y_{\leq 3}, q) \quad (305)$$

The model is trained by minimizing log-likelihood over the training data, using RMSProp for optimization. At inference time, greedy search or beam search is used to predict the most probable output sequence.

## Seq2SQL: Generating Structured Queries from NL using RL

Table of Contents Local

Written by Brandon McKinzie

Zhong, Xiong, and Socher, “Seq2SQL: Generating Structured Queries From Natural Language Using Reinforcement Learning”

**Overview.** Deep neural network for translating natural language questions to corresponding SQL queries. Outperforms state-of-the-art semantic parser.

**Seq2Tree and Pointer Baseline.** Baseline model is the Seq2Tree model from the previous note on Dong & Lapata’s (2016) paper. Authors here argue their output space is unnecessarily large, and employ the idea of pointer networks with augmented inputs. The input sequence is the concatenation of (1) the column names, (2) the limited vocabulary of the SQL language such as `SELECT`, `COUNT`, etc., and (3) the question.

$$x := [<\text{col}>; x_1^c; x_2^c; \dots; x_N^c; <\text{sql}>; x^s; <\text{question}>; x^q] \quad (306)$$

$$x_j^c \in \mathbb{R}^{T_j}$$

where we also insert special (“sentinel”) tokens to demarcate the boundaries of each section. The pointer network can then produce the SQL query by selecting exclusively from the input. Let  $g_s$  denote the  $s$ th decoder [hidden] state, and  $y_s$  denote the output (index/pointer to input query token).

$$[\text{ptr net}] y_s = \arg \max_t (\alpha_{s,t}^{ptr}) \quad \text{where} \quad \alpha_{s,t}^{ptr} = w^{ptr} \cdot \tanh(U^{ptr} g_s + V^{ptr} h_t) \quad (307)$$

**Seq2SQL.**

- Aggregation Classifier.** Our goal here is to predict which aggregation operation to use out of `COUNT`, `MIN`, `MAX`, `NULL`, etc. This is done by projecting the attention-weighted average of encoder states,  $\kappa^{agg}$ , to  $\mathbb{R}^C$  where  $C$  denotes the number of unique aforementioned aggregation operations. The sequence of computations is summarized as follows:

$$\alpha_t^{inp} = w^{inp} \cdot h_t^{enc} \quad (308)$$

$$\beta^{inp} = \text{softmax}(\alpha^{inp}) \quad (309)$$

$$\kappa^{agg} = \sum_t \beta_t^{inp} h_t^{enc} \quad (310)$$

$$\alpha^{agg} = W^{agg} \tanh(V^{agg} \kappa^{agg} + b^{agg}) + c^{agg} \quad (311)$$

$$\beta^{agg} = \text{softmax}(\alpha^{agg}) \quad (312)$$

$$W^{agg} \in \mathbb{R}^{C \times T}$$

$$\beta^{agg} \in \mathbb{R}^C$$

where  $\beta_i^{agg}$  gives the probability for the  $i$ th aggregation operation. We use cross entropy loss  $L^{agg}$  for determining the aggregation operation. Note that this part isn't really a sequence-to-sequence architecture. It's nothing more than an MLP applied to an attention-weighted average of the encoder states.

2. **Get Pointer to Column.** A pointer network is used for identifying which column in the input representation should be used in the query. Recall that  $x_{j,t}^c$  denotes the  $t$ th word in column  $j$ . We use the last encoder state for a given column's LSTM<sup>93</sup> as its representation;  $T_j$  denotes the number of words in the  $j$ th column.

$$e_j^c = h_{j,T_j}^c \quad \text{where} \quad h_{j,t}^c = \text{LSTM}\left(\text{emb}(x_{j,t}^c), h_{j,t-1}^c\right) \quad (313)$$

To construct a representation for the question, compute another input representation  $\kappa^{sel}$  using the same architecture (but distinct weights) as for  $\kappa^{agg}$ . As usual, we compute the scores for each column  $j$  via:

$$\alpha_j^{sel} = W^{sel} \tanh\left(V^{sel} \kappa^{sel} + V^c e_j^c\right) \quad (314)$$

$$\beta^{sel} = \text{softmax}(\alpha^{sel}) \quad (315)$$

Similar to the aggregation, we train the SELECT network using cross entropy loss  $L^{sel}$ .

3. **WHERE Clause Pointer Decoder.** Recall from equation 307 that this is a model with recurrent connections from its *outputs leading back into its inputs*, and thus a common approach is to train it with **teacher forcing**<sup>94</sup>. However, since the boolean expressions within a WHERE clause can be swapped around while still yielding the same SQL query, reinforcement learning (instead of cross entropy) is used to *learn a policy to directly optimize the expected correctness of the execution result*. Note that this also implies that we will be sampling from the output distribution at decoding step  $s$  to obtain the next input for  $s + 1$  [instead of teacher forcing].

---

<sup>93</sup>Yes, we encode each column with an LSTM separately.

<sup>94</sup>Teacher forcing is just a name for how we train the decoder portion of a sequence-to-sequence model, wherein we feed the ground-truth output  $y^{(t)}$  as input at time  $t + 1$  during training.

$$R(q(y), q_g) = \begin{cases} -2 & \text{if } q(y) \text{ is not a valid SQL query} \\ -1 & \text{if } q(y) \text{ is a valid SQL query and executes to an incorrect result} \\ +1 & \text{if } q(y) \text{ is a valid SQL query and executes to the correct result} \end{cases} \quad (316)$$

$$L^{whe} = -\mathbb{E}_y [R(q(y), q_g)] \quad (317)$$

$$\nabla L_\Theta^{whe} = -\nabla_\Theta (\mathbb{E}_{y \sim p_y} [R(q(y), q_g)]) \quad (318)$$

$$= -\mathbb{E}_{y \sim p_y} \left[ R(q(y), q_g) \nabla_\Theta \sum_t \log p_y(y_t; \Theta) \right] \quad (319)$$

$$\approx -R(q(y), q_g) \nabla_\Theta \sum_t \log p_y(y_t; \Theta) \quad (320)$$

where

- $y = [y^1, y^2, \dots, y^T]$  denotes the sequences of generated tokens in the WHERE clause.
- $q(y)$  denotes the query generated by the model.
- $q_g$  denotes the ground truth query corresponding to the question.

and the gradient has been approximated in the last line using a single Monte-Carlo sample  $y$ .

Finally, the model is trained using gradient descent to minimize  $L = L^{agg} + L^{sel} + L^{whe}$ .

**Speculations for Event Extraction.** I want to explore using this paper's model for the task of event extraction. Below, I've replaced some words (shown in green) from a sentence in the paper in order to formalize this as event extraction.

*Seq2Event takes as input a sentence and the possible event types of an ontology. It generates the corresponding event annotation, which, during training, is compared against an event template. The result of the comparison is utilized to train the reinforcement learning algorithm<sup>95</sup>.*

---

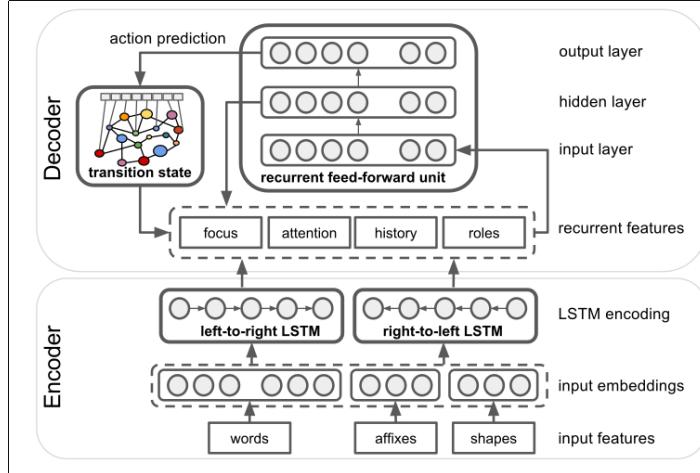
<sup>95</sup>Original: Seq2SQL takes as input a question and the columns of a table. It generates the corresponding SQL query, which, during training, is executed against a database. The result of the execution is utilized as the reward to train the reinforcement learning algorithm.

## SLING: A Framework for Frame Semantic Parsing

Table of Contents Local

Written by Brandon McKinzie

M. Ringgaard, R. Gupta, F. Pereira, “SLING: A framework for frame semantic parsing” (2017)



### Model.

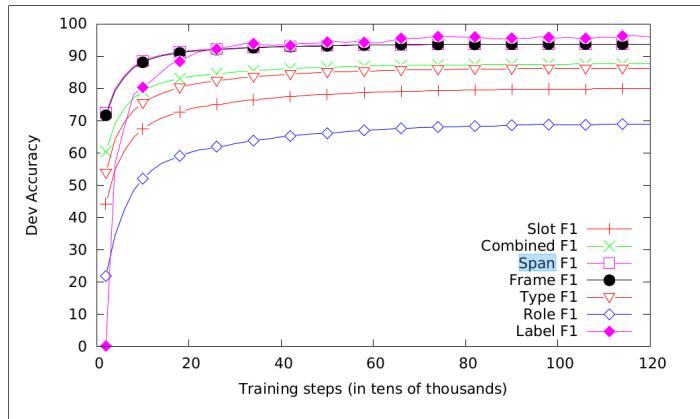
- **Inputs.** [words; affixes; shapes]
- **Encoder.**
  1. Embed.
  2. Bidirectional LSTM.
- **Inputs to TBRU.**
  - BLSTM [forward and backward] hidden state for the current token in the parser state.
  - **Focus.** Hidden layer activations corresponding to the transition steps that evoked/brought into focus the top- $k$  frames in the attention buffer.
  - **Attention.** Recall that we maintain an **attention buffer**: an ordered list of frames, where the order represents closeness to center of attention. The attention portion of inputs for the TBRU looks at the top- $k$  frames in the attention buffer, finds the phrases in the text (if any) that evoked them. The activations from the BLSTM for the last token of each of those phrases are included as TBRU inputs<sup>96</sup>
  - **History.** Hidden layer activations from the previous  $k$  steps.
  - **Roles.** Embeddings of  $(s_i, r_i, t_i)$ , where the frame at position  $s_i$  in the attention buffer has a role (key)  $r_i$  with frame at position  $t_i$  as its value. Back-off features are added for the source roles  $(s_i, r_i)$ , target role  $(r_i, t_i)$ , and unlabeled roles  $(s_i, t_i)$ .
- **Decoder (TBRU).** Outputs a softmax over possible transitions (actions).

<sup>96</sup>Okay, how is this attention at all? Seems misleading to call it attention.

**Transition System.** Below is the list of possible actions. Note that, since the system is trained to predict the correct *frame graph* result, it isn’t directly told what order it should take a given set of actions<sup>97</sup>.

- **SHIFT**. Move to next input token.
- **STOP**. Signal that we’ve reached end of parse.
- **EVOKE(type, num)**. New frame of **type** from next **num** tokens in the input; placed at front of attention buffer.
- **REFER(frame, num)**. New mention from next **num** tokens, evoking existing **frame** from attention buffer. Places at front.
- **CONNECT(source-frame, role, target-frame)**. Inserts **(role, target-frame)** slot into **source-frame**, move **source-frame** to front.
- **ASSIGN(source-frame, role, value)**. Same as CONNECT, but with primitive/constant **value**.
- **EMBED(target-frame, role, type)**. New frame of **type**, and inserts **(role, target-frame)** slot. New frame placed to front.
- **ELABORATE(source-frame, role, type)**. New frame of **type**. Inserts **(role, new-frame)** slot to **source-frame**. New frame placed at front.

**Evaluation.** Need some way of comparing an annotated document with its gold-standard annotation. This is done by constructing a virtual graph where the document is the start node. It is then connected to the spans (which are presumably nodes themselves), and the spans are connected to the frames they evoke. Frames that refer to other frames are given corresponding edges between them. Quality is computed by aligning the golden and predicted graphs and computing precision, recall, and F1. Specifically, these scores are computed separately for spans, frames, frame types, roles linking to other frames (referred to here as just “roles”), and roles that link to global constants (referred to here as just “labels”). Results are shown below.



<sup>97</sup>This is important to keep in mind, since more than one sequence of actions can result in a given predicted frame graph.

## Poincaré Embeddings for Learning Hierarchical Representations

Table of Contents Local

Written by Brandon McKinzie

M. Nickel and D. Kiela, “Poincaré Embeddings for Learning Hierarchical Representations” (2017)

**Introduction.** Dimensionality of embeddings can become prohibitively large when needed for complex data. Authors focus on mitigating this problem for large datasets whose objects can be organized according to a latent hierarchy<sup>98</sup>. They propose to compute embeddings in a particular model of hyperbolic space, the **Poincaré ball model**, claiming it is well-suited for gradient-based optimization (they make use of **Riemannian optimization**).

**Prerequisite Math.** Recall that a hyperbola is a set of points, such that for any point  $P$  of the set, the absolute difference of the distances  $|PF_2|$ ,  $|PF_1|$  to two fixed points  $F_1$ ,  $F_2$  (the foci), is constant, usually denoted by  $2a$ ,  $a > 0$ . We can define a hyperbola by this set of points or by its canonical form, which are both given, respectively, as:

$$H = \{P \mid ||PF_2| - |PF_1|| = 2a\} \quad (321)$$

$$\frac{x^2}{a^2} - \frac{y^2}{b^2} = 1 \quad (322)$$

where  $b^2 := c^2 - a^2$ ,  $(\pm a, 0)$  are the two vertices, and  $(\pm c, 0)$  are the two foci. Cannon et al. define  $n$ -dimensional hyperbolic space by the formula

$$H^n = \{x \in \mathbb{R}^{n+1} : x * x = -1\} \quad (323)$$

where  $*$  denotes the non-euclidean inner product (subtracts last term; same as Minkowski sapce-time). Notice that this is the defining equation for a hyperboloid of two sheets, and Cannon et al. says “usually we deal only with one of the two sheets.” Hyperbolic spaces are well-suited to model hierarchical data, since both circle length and disc area grow *exponentially* with  $r$ .

---

<sup>98</sup>This begs the question: how useful would a Poincaré embedding be for situations where this assumption isn’t valid?

**Poincaré Embeddings.** Let  $\mathcal{B}^d = \{\mathbf{x} \in \mathbb{R}^d \mid \|\mathbf{x}\| < 1\}$  be the open d-dimensional unit ball. The **Poincaré ball** model of hyperbolic space corresponds then to the Riemannian manifold<sup>99</sup>  $(\mathcal{B}^d, g_{\mathbf{x}})$ , where

$$g_{\mathbf{x}} = \left( \frac{2}{1 - \|\mathbf{x}\|^2} \right)^2 g^E \quad (324)$$

is the **Riemannian metric tensor**, and  $g^E$  denotes the Euclidean metric tensor. The distance between two points  $\mathbf{u}, \mathbf{v} \in \mathcal{B}^d$  is given as

$$d(\mathbf{u}, \mathbf{v}) = \operatorname{arccosh} \left( 1 + 2 \frac{\|\mathbf{u} - \mathbf{v}\|^2}{(1 - \|\mathbf{u}\|^2)(1 - \|\mathbf{v}\|^2)} \right) \quad (325)$$

The boundary of the ball corresponds to the sphere  $S^{d-1}$  and is denoted by  $\partial\mathcal{B}$ . Geodesics in  $\mathcal{B}^d$  are then circles that are orthogonal to  $\partial\mathcal{B}$ . To compute Poincaré embeddings for a set of symbols  $\mathcal{S} = \{\mathbf{x}_i\}_{i=1}^n$ , we want to find embeddings  $\Theta = \{\boldsymbol{\theta}_i\}_{i=1}^n$ , where  $\boldsymbol{\theta}_i \in \mathcal{B}^d$ . Given some loss function  $\mathcal{L}$  that encourages semantically similar objects to be close as defined by the Poincaré distance, our goal is to solve the optimization problem

$$\Theta' \leftarrow \arg \min_{\Theta} \mathcal{L}(\Theta) \quad s.t. \quad \forall \boldsymbol{\theta}_i \in \Theta : \|\boldsymbol{\theta}_i\| < 1 \quad (326)$$

**Optimization.** Let  $\mathcal{T}_{\boldsymbol{\theta}}\mathcal{B}$  denote the **tangent space** of a point  $\boldsymbol{\theta} \in \mathcal{B}^d$ . Let  $\nabla_R \in \mathcal{T}_{\boldsymbol{\theta}}\mathcal{B}$  denote the **Riemannian gradient** of  $\mathcal{L}(\boldsymbol{\theta})$ , and  $\nabla_E$  the Euclidean gradient of  $\mathcal{L}(\boldsymbol{\theta})$ . Using **RSGD**, parameter updates take the form

$$\boldsymbol{\theta}_{t+1} = \mathfrak{R}_{\boldsymbol{\theta}_t}(-\eta_t \nabla_R \mathcal{L}(\boldsymbol{\theta}_t)) \quad (327)$$

where  $\mathfrak{R}_{\boldsymbol{\theta}_t}$  denotes the **retraction** onto  $\mathcal{B}$  at  $\boldsymbol{\theta}$  and  $\eta_t$  denotes the learning rate at time  $t$ .

---

<sup>99</sup>All five analytic models of hyperbolic geometry in Cannon et al. are differentiable manifolds with a Riemannian metric. A **Riemannian metric**  $ds^2$  on Euclidean space  $\mathbb{R}^n$  is a function that assigns at each point  $p \in \mathbb{R}^n$  a positive definite symmetric inner product on the tangent space at  $p$ , this inner product varying differentiably with the point  $p$ . If  $x_1, \dots, x_n$  are the standard coordinates in  $\mathbb{R}^n$ , then  $ds^2$  has the form  $\sum_{i,j} g_{ij} dx_i dx_j$ , and the matrix  $(g_{ij})$  depends differentiably on  $x$  and is positive definite and symmetric.

## Enriching Word Vectors with Subword Information (FastText)

Table of Contents Local

Written by Brandon McKinzie

P. Bojanowski, E. Grave, A. Joulin, and T. Mikolov, “Enriching Word Vectors with Subword Information” (2017)

**Overview.** Based on the skipgram model, but where each word is represented as a bag of character  $n$ -grams. A vector representation is associated *each* character  $n$ -gram; words being represented as the *sum*<sup>100</sup> of these representations.

**Skipgram with Negative Sampling.** Since this is based on skipgram, recall the objective of skipgram which is to maximize:

$$\sum_{t=1}^T \sum_{c \in \mathcal{C}_t} \log \Pr [w_c | w_t] \quad (328)$$

for a sequence of words  $w_1, \dots, w_T$ . One way of parameterizing  $\Pr [w_c | w_t]$  is by computing a softmax over a scoring function  $s : (w_t, w_c) \mapsto \mathbb{R}$ ,

$$\Pr [w_c | w_t] = \frac{e^{s(w_t, w_c)}}{\sum_{j=1}^W e^{s(w_t, j)}} \quad (329)$$

However, this implies that, given  $w_t$ , we only predict one context word  $w_c$ . Instead, we can frame the problem as a set of independent binary classification tasks, and independently predict the presence/absence of context words. Let  $\ell : x \mapsto \log(1 + e^{-x})$  denote the standard logistic negative log-likelihood. Our objective is:

$$\sum_{t=1}^T \left[ \sum_{c \in \mathcal{C}_t} \ell(s(w_t, w_c)) + \sum_{n \in \mathcal{N}_{t,c}} \ell(-s(w_t, n)) \right] \quad (330)$$

where  $\mathcal{N}_{t,c}$  is a set of negative examples sampled from the vocabulary. A common scoring function involves associating a distinct input vector  $u_w$  and output vector  $v_w$  for each word  $w$ . Then the score is computed as  $s(w_t, w_c) = \mathbf{u}_{w_t}^T \mathbf{v}_{w_c}$ .

---

<sup>100</sup>It would be interesting to explore other aggregation operations than just summation.

**FastText.** Main contribution is a different scoring function  $s$  that utilizes subword information. Each word  $w$  is represented as a bag of character  $n$ -grams. Special symbols  $<$  and  $>$  delimit word boundaries, and the authors also insert the special sequence containing the full word (with the delimiters) in its bag of  $n$ -grams. The word *where* is thus represented by first building its bag of  $n$ -grams, for the choice of  $n = 3$ :

$$\text{where} \longrightarrow \{<wh, whe, her, ere, re>, <where>\} \quad (331)$$

Such a set of  $n$ -grams for a word  $w$  is denoted  $\mathcal{G}_w$ . Each  $n$ -gram  $g$  for a word  $w$  has its own vector  $\mathbf{z}_g$ , and the final vector representation of  $w$  is the sum of these. The scoring function becomes

$$s(w, c) = \sum_{g \in \mathcal{G}_w} \mathbf{z}_g^T \mathbf{v}_c \quad (332)$$

## DeepWalk: Online Learning of Social Representations

[Table of Contents](#)   [Local](#)

*Written by Brandon McKinzie*

B. Perozzi, R. Al-Rfou, and S. Skiena, “DeepWalk: Online Learning of Social Representations,” (2014).

**Problem Definition.** Classifying members of a social network into one or more categories.

*Let  $G = (V, E)$ , where  $V$  are the members of the network, and  $E$  be its edges,  $E \subseteq (V \times V)$ . Given a partially labeled social network  $G_L = (V, E, X, Y)$ , with attributes  $X \in \mathbb{R}^{|V| \times S}$  where  $S$  is the size of the feature space for each attribute vector, and  $Y \in \mathbb{R}^{|V| \times |\mathcal{Y}|}$ ,  $\mathcal{Y}$  is the set of labels.*

In other words, the elements of our training dataset,  $(X, Y)$ , are the members of the social network, and we want to label each member, represented by a vector in  $\mathbb{R}^S$ , with one or more of the  $|\mathcal{Y}|$  labels. We aim to learn features that capture the graph structure *independent* of the labels’ distribution, and to do so in an unsupervised fashion.

**Learning Social Representations.** We want the representations to be adaptable, community-aware, low-dimensional, and continuous. The authors’ method learns representations for vertices from a stream of short random walks, optimized with techniques from language modeling.

- **Random Walks.** Denote a random walk rooted at vertex  $v_i$  as  $\mathcal{W}_{v_i}$ , where the  $k$ th visited vertex is chosen at random from the neighbors of the  $(k - 1)^{th}$  visited vertex, and so on. Motivation for their use here is that they’re “the foundation of a class of *output sensitive* algorithms which use them to compute local community structure information in time sublinear to the size of the input graph.”
- **Language Modeling.** Authors present a generalization of language modeling, which traditionally aims to maximize  $\Pr [w_n | w_0, \dots, w_{n-1}]$  over all words in a training corpus. The motivation of the generalization is to explore the graph through a stream of short random walks. The walks are thought of as short sentences/phrases in a special language, and we want to estimate the probability of observing vertex  $v_i$  given all previous vertices so far in the random walk. Since we want to learn a latent social representation of each vertex, and not simply a probability distribution over node co-occurrences, we condition on the *embeddings* of visited nodes in this latent space (rather than the nodes themselves directly)

$$\Pr [v_i | \Phi(v_1), \Phi(v_2), \dots, \Phi(v_{i-1})] \quad (333)$$

where, in practice, the mapping function  $\Phi$  is represented by a  $|V| \times d$  matrix of free parameters (an embedding matrix). Since this becomes infeasible to compute as the walk length grows, the authors opt for an approach resembling CBOW: minimizing the NLL of vertices in the context of a given vertex.

$$\min_{\Phi} -\log \Pr [v_{i-w}, \dots, v_{i-1}, v_{i+1}, \dots, v_{i+w} | \Phi(v_i)] \quad (334)$$

Remember that, here,  $v_j$  is the  $j$ th vertex visited in some given random walk.

**DeepWalk Algorithm.** Below is a conceptual summary of the procedure, followed by a figure/illustration of the formal algorithm definition.

1. **Inputs.** Graph  $G(V, E)$ , window size  $w$ , embedding size  $d$ , walks per vertex  $\gamma$ , walk length  $t$ .
2. **Random Walk.** For each vertex  $v_i$ , compute  $\mathcal{W}_{v_i} := \text{RandomWalk}(G, v_i, t)$ .
3. **Updates.** Upon finishing a walk,  $\mathcal{W}_{v_i}$ , run skipgram on the sequence of walked vertices to update the embedding matrix  $\Phi$ .
4. **Outputs.** The embedding matrix  $\Phi \in \mathbb{R}^{|V| \times d}$ .

---

**Algorithm 1** DEEPWALK( $G, w, d, \gamma, t$ )

---

**Input:** graph  $G(V, E)$   
 window size  $w$   
 embedding size  $d$   
 walks per vertex  $\gamma$   
 walk length  $t$

**Output:** matrix of vertex representations  $\Phi \in \mathbb{R}^{|V| \times d}$

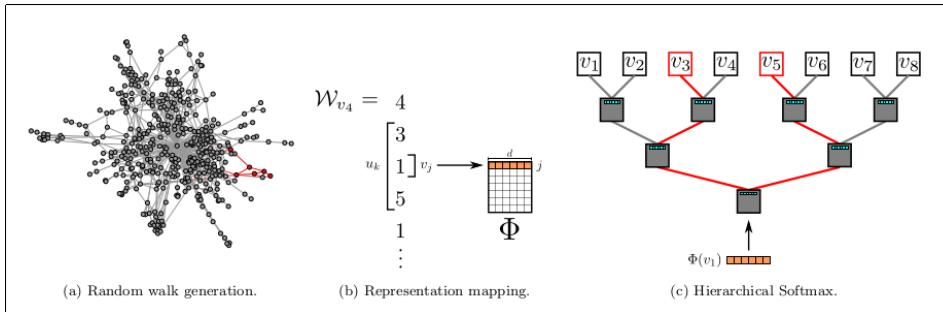
```

1: Initialization: Sample  $\Phi$  from  $\mathcal{U}^{|V| \times d}$ 
2: Build a binary Tree  $T$  from  $V$ 
3: for  $i = 0$  to  $\gamma$  do
4:    $\mathcal{O} = \text{Shuffle}(V)$ 
5:   for each  $v_i \in \mathcal{O}$  do
6:      $\mathcal{W}_{v_i} = \text{RandomWalk}(G, v_i, t)$ 
7:     SkipGram( $\Phi, \mathcal{W}_{v_i}, w$ )
8:   end for
9: end for

```

---

where  $\text{SkipGram}(\Phi, W_{v_i}, w)$  performs SGD updates on  $\Phi$  to minimize  $-\log \Pr [u_k | \Phi(v_j)]$  for each visited  $v_j$ , for each  $u_k$  in the “context” of  $v_j$ . Notice that a binary tree  $T$  is build from the set of vertices  $V$  (line 2) – this is done as preparation for computing each  $\Pr [u_k | \Phi(v_j)]$  via a *hierarchical softmax*, to reduce computational burden of its partition function. Finally, a visual overview of the DeepWalk algorithm is shown below. The authors use this algorithm, com-



bined with a one-vs-rest logistic regression implementation by LibLinear, for various multiclass multilabel classification tasks.

## Review of Relational Machine Learning for Knowledge Graphs

[Table of Contents](#)   [Local](#)

*Written by Brandon McKinzie*

Nickel, Murphy, Tresp, and Gabrilovich, “Review of Relational Machine Learning for Knowledge Graphs,” (2015).

**Introduction.** Paper discusses latent feature models such as tensor factorization and multiway neural networks, and mining observable patterns in the graph. In **Statistical Relational Learning (SRL)**, the representation of an object can contain its relationships to other objects. The main goals of SRL include:

- Prediction of missing edges (relationships between entities).
- Prediction of properties of nodes.
- Clustering nodes based on their connectivity patterns.

We’ll be reviewing how SRL techniques can be applied to large-scale **knowledge graphs** (KGs), i.e. graph structured knowledge bases (KBs) that store factual information in the form of relationships between entities.

**Probabilistic Knowledge Graphs.** Let  $\mathcal{E} = \{e_1, \dots, e_{N_e}\}$  be the set of all entities and  $\mathcal{R} = \{r_1, \dots, r_{N_r}\}$  be the set of all relation types in a KG. We model each *possible* triple  $x_{ijk} = (e_i, r_k, e_j)$  as a binary random variable  $y_{ijk} \in \{0, 1\}$  that indicates its existence. The full tensor  $\mathbf{Y} \in \{0, 1\}^{N_e \times N_e \times N_r}$  is called the *adjacency tensor*, where each possible realization of  $\mathbf{Y}$  can be interpreted as a possible world.

Clearly,  $\mathbf{Y}$  will be large and sparse in most applications. Ideally, a relational model for large-scale KGs should scale at most linearly with the data size, i.e., linearly in the number of entities  $N_e$ , linearly in the number of relations  $N_r$ , and linearly in the number of *observed* triples  $|\mathcal{D}| = N_d$ .

**Types of SRL Models.** The presence or absence of certain triples in relational data is correlated with (i.e. predictive of) the presence or absence of certain other triples. In other words, the random variables  $y_{ijk}$  are correlated with each other. There are three main ways to model these correlations:

1. **Latent feature models:** Assume all  $y_{ijk}$  are conditionally independent given latent features associated with the subject, object and relation type and additional parameters.
2. **Graph feature models:** Assume all  $y_{ijk}$  are conditionally independent given observed graph features and additional parameters.
3. **Markov Random Fields:** Assume all  $y_{ijk}$  have local interactions.

The first two model classes predict the existence of a triple  $x_{ijk}$  via a score function  $f(x_{ijk}; \Theta)$

which represents the model's confidence that a triple exists given the parameters  $\Theta$ . The conditional independence assumptions can be written as

$$\Pr [\mathbf{Y} | \mathcal{D}, \Theta] = \prod_{i=1}^{N_e} \prod_{j=1}^{N_e} \prod_{k=1}^{N_r} \text{Ber}(y_{ijk} | \sigma(f(x_{ijk}; \Theta))) \quad (335)$$

where  $\text{Ber}$  is the Bernoulli distribution<sup>101</sup>. Such models will be referred to as *probabilistic models*. We will also discuss *score-based models*, which optimize  $f(\cdot)$  via maximizing the margin between existing and non-existing triples.

**Latent Feature Models.** We assume the variables  $y_{ijk}$  are conditionally independent given a set of global latent features and parameters. All LFM explain triples (observable facts) via latent features of entities<sup>102</sup>. One task of all LFM is to infer these [latent] features automatically from the data.

- **RESCAL:** a bilinear model. Models the score of a triple  $x_{ijk}$  as

$$f_{ijk}^{\text{RESCAL}} := \mathbf{e}_i^T \mathbf{W}_k \mathbf{e}_j = \sum_{a=1}^{H_e} \sum_{b=1}^{H_e} w_{abk} e_{ia} e_{jb} \quad (337)$$

where the entity vectors  $\mathbf{e}_i \in \mathbb{R}^{H_e}$  and  $H_e$  denotes the number of latent features in the model. The parameters of the model are  $\Theta = \{\{\mathbf{e}_i\}_{i=1}^{N_e}, \{\mathbf{W}_k\}_{k=1}^{N_r}\}$ . Note that entities have the same latent representation regardless of whether they occur as subjects or objects in a relationship (shared representation), *thus allowing the model to capture global dependencies in the data*. We can make a connection to **tensor factorization** methods by seeing that the equation above can be written compactly as

$$\mathbf{F}_k = \mathbf{E} \mathbf{W}_k \mathbf{E}^T \quad (338)$$

where  $\mathbf{F}_k \in \mathbb{R}^{N_e \times N_e}$  is the matrix holding all scores for the  $k$ -th relation, and the  $i$ th row of  $\mathbf{E} \in \mathbb{R}^{N_e \times H_e}$  holds the latent representation of  $\mathbf{e}_i$ .

- **Multi-layer perceptrons.** We can rewrite RESCAL as

$$f_{ijk}^{\text{RESCAL}} := \mathbf{w}_k^T \phi_{i,j}^{\text{RESCAL}} \quad (339)$$

$$\phi_{i,j}^{\text{RESCAL}} := \mathbf{e}_j \otimes \mathbf{e}_i \quad (340)$$

where  $\mathbf{w}_k = \text{vec}(\mathbf{W}_k)$  (vector of size  $H_e^2$  obtained by stacking columns of  $\mathbf{W}_k$ ). The

<sup>101</sup>Notation used:

$$\text{Ber}(y | p) = \begin{cases} p & \text{if } y = 1 \\ 1 - p & \text{if } y = 0 \end{cases} \quad (336)$$

<sup>102</sup>It appears that "latent" is being used here synonymously with "not directly observed in the data".

authors extend this to what they call the E-MLP (E for entity) model:

$$f_{ijk}^{E-MLP} := \mathbf{w}_k^T \mathbf{g}(\mathbf{h}_{ijk}^a) \quad (341)$$

$$\mathbf{h}_{ijk}^a := \mathbf{A}_k^T \phi_{ij}^{E-MLP} \quad (342)$$

$$\phi_{ij}^{E-MLP} := [\mathbf{e}_i; \mathbf{e}_j] \quad (343)$$

**Graph Feature Models.** Here we assume that the existence of an edge can be predicted by extracting features from the observed edges in the graph. In contrast to LFM, this kind of reasoning explains triples directly from the observed triples in the KG.

- **Similarity measures for uni-relational data.** Link prediction in graphs that consist only of a single relation (e.g. (Bob, isFriendOf, Sally) for a social network). Various **similarity indices** have been proposed to measure similarity of entities, of which there are three main classes:

1. **Local** similarity indices: Common Neighbors, Adamic-Adar index, Preferential Attachment derive entity similarities from number of common neighbors.
2. **Global** similarity indices: Katz index, Leicht-Holme-Newman index (ensembles of all paths bw entities). Hitting Time, Commute Time, PageRank (random walks).
3. **Quasi-local** similarity indices: Local Katz, Local Random Walks.

- **Path Ranking Algorithm (PRA):** extends the idea of using random walks of bounded lengths for predicting links in multi-relational KGs. Let  $\pi_L(i, j, k, t)$  denote a path of length  $L$  of the form  $e_i \xrightarrow{r_1} e_2 \xrightarrow{r_2} e_3 \cdots \xrightarrow{r_L} e_j$ , where  $t$  represents the sequence of edge types  $t = (r_1, r_2, \dots, r_L)$ . We also require there to be a direct arc  $e_i \xrightarrow{r_k} e_j$ , representing the existence of a relationship of type  $k$  from  $e_i$  to  $e_j$ . Let  $\Pi_L(i, j, k)$  represent the set of all such paths of length  $L$ , ranging over path types  $t$ .

We can compute the probability of following a given path by assuming that at each step we follow an outgoing link uniformly at random. Let  $\Pr[\pi_L(i, j, k, t)]$  be the probability of the path with type  $t$ . *The key idea in PRA is to use these path probabilities as features for predicting the probabilities of missing edges.* More precisely, the feature vector and score function (logistic regression) are as follows:

$$\phi_{ijk}^{PRA} = [\Pr[\pi] : \pi \in \Pi_L(i, j, k)] \quad (344)$$

$$f_{ijk}^{PRA} := \mathbf{w}_k^T \phi_{ijk}^{PRA} \quad (345)$$

**TODO:** Finish...

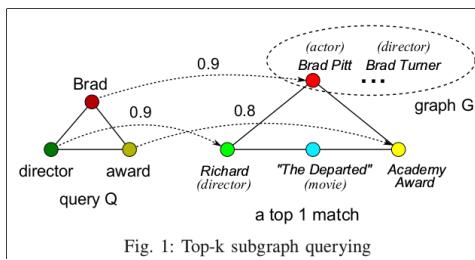
## Fast Top-K Search in Knowledge Graphs

Table of Contents Local

Written by Brandon McKinzie

S. Yang, F. Han, Y. Wu, X. Yan, “Fast Top-K Search in Knowledge Graphs.”

**Introduction.** Task: Given a knowledge graph  $G$ , scoring function  $F$ , and a graph query  $Q$ , top-k subgraph search over  $G$  returns  $k$  answers with the highest matching scores. An example is searching for movie makers (directors) worked with “Brad” and have won awards, illustrated below:



Clearly, it would be extremely inefficient to enumerate all possible matches and then rank them.

### Preliminaries/Terminology.

- **Queries.** A query [graph] is defined as  $Q = (V_Q, E_Q)$ . Each query node in  $Q$  provides information/constraints about an entity, and an edge between two nodes specifies the relationship or the connectivity constraint posed on the two nodes.  $Q^*$  denotes a star-shaped query, which is basically a graph that looks like a star (central node with tree-like structure radially outward).
- **Subgraph Matching.** Given a graph query  $Q$  and a knowledge graph  $G$ , a match  $\phi(Q)$  of  $Q$  in  $G$  is a subgraph of  $G$ , specified by a one-to-one matching function  $\phi$ . It maps each node  $u$  (edge  $e = (u, u')$ ) in  $Q$  to a node match  $\phi(u)$  (edge match  $\phi(e) = (\phi(u), \phi(u'))$ ) in  $\phi(Q)$ .

The matching score between query  $Q$  and its match  $\phi(Q)$  is

$$F(\phi(Q)) = \sum_{v \in V_Q} F_V(v, \phi(v)) + \sum_{e \in E_Q} F_E(e, \phi(e)) \quad (346)$$

$$F_V(v, \phi(v)) = \sum_i \alpha_i f_i(v, \phi(v)) \quad (347)$$

$$F_E(e, \phi(e)) = \sum_j \beta_j f_j(e, \phi(e)) \quad (348)$$

### Star-Based Top-K Matching.

1. **Query decomposition:** Given query  $Q$ , STAR decomposes  $Q$  to a set of star queries  $\mathcal{Q}$ . A star query contains a pivot node and a set of leaves as its neighbors in  $Q$ .
2. **Star querying engine:** Generate a set of top matches for each star query  $\mathcal{Q}$ .
3. **Top-k rank join.** The top matches for multiple star queries are collected and joined to produce top-k complete matches of  $Q$ .

## Dynamic Recurrent Acyclic Graphical Neural Networks (DRAGNN)

Table of Contents Local

Written by Brandon McKinzie

Kong et al., “DRAGNN: A Transition-based Framework for Dynamically Connected Neural Networks,” (2017).

**Transition Systems.** Define a transition system  $\mathcal{T} \triangleq \{\mathcal{S}, \mathcal{A}, t\}$ , where

- $\mathcal{S} = \mathcal{S}(x)$  is a set of states, where that set depends on the input sequence  $x$ .
- A special start state  $s^\dagger \in \mathcal{S}(x)$ .
- A set of allowed decisions  $\mathcal{A}(s, x) \forall s \in \mathcal{S}(x)$ .
- A transition function  $t(s, d, x)$  returning a new state  $s'$  for any decision  $d \in \mathcal{A}(s, x)$ .

The authors then define a **complete structure** as a sequence of state/decision pairs  $(s_1, d_1) \dots (s_n, d_n)$  such that  $s_1 = s^\dagger$ ,  $d_i \in \mathcal{A}(s_i)$  for  $i = 1, \dots, n$  and  $s_{i+1} = t(s_i, d_i)$ , where  $n = n(x)$  is the number of decisions for input  $x$ <sup>103</sup>. We’ll use transition systems to map inputs  $x$  into a sequence of output symbols  $d_1, \dots, d_n$ .

**Transition Based Recurrent Networks.** When combining transition systems with recurrent networks, we will refer to them as **Transition Based Recurrent Units (TBRU)**, which consist of:

- Transition system  $\mathcal{T}$ .
- Input function  $\mathbf{m}(s) : \mathcal{S} \mapsto \mathbb{R}^K$  that maps states to some fixed-size vector representation (e.g. an embedding lookup operation).
- Recurrence function  $\mathbf{r}(s) : \mathcal{S} \mapsto \mathbb{P}\{1, \dots, i-1\}$  that maps states to a set of previous time steps, where  $\mathbb{P}$  is the power set. Note that  $|\mathbf{r}(s)|$  may vary with  $s$ . We use  $\mathbf{r}$  to specify state-dependent recurrent links in the unrolled computation graph.
- The RNN cell  $\mathbf{h}_s \leftarrow \mathbf{RNN}(\mathbf{m}(s), \{\mathbf{h}_i \mid i \in \mathbf{r}(s)\})$ .

**Example: Sequential tagging RNN.** Let  $\mathbf{x} = \{x_1, \dots, x_n\}$  be a sequence of input tokens. Let the  $i$ th output,  $d_i$ , be a tag from some predefined set of tags  $\mathcal{A}$ . Then our model can be defined as:

- Transition system:  $\mathcal{T} = \{ s_i = \mathcal{S}(x_i) = \{1, \dots, d_{i-1}\}, \mathcal{A}, t(s_i, d_i, x_i) = s_{i+1} = s_i + \{d_i\} \}$ .
- Input function:  $\mathbf{m}(s_i) = \text{embed}(x_i)$ .
- Recurrence function:  $\mathbf{r}(s_i) = \{i-1\}$  to connect the network to the previous state.
- RNN cell:  $\mathbf{h}_i \leftarrow \text{LSTM}(\mathbf{m}(s_i) \mid \mathbf{r}(s_i) = \{i-1\})$ .

---

<sup>103</sup>The authors state that we are only concerned with complete structures that have the same number of decisions  $n(x)$  for the same input  $x$ .

### Example: Parsey McParseface.

- Transition system: the **arc-standard** transition system, defined in image below<sup>104</sup>.

<b>Initialization:</b>	$c_s(x = x_1, \dots, x_n) = ([0], [1, \dots, n], \emptyset)$
<b>Terminal:</b>	$C_t = \{c \in C \mid c = ([0], [], A)\}$
<b>Transitions:</b>	$(\sigma, [i \beta], A) \Rightarrow ([\sigma i], \beta, A)$ (SHIFT) $([\sigma i j], B, A) \Rightarrow ([\sigma j], B, A \cup \{(j, l, i)\})^1$ (LEFT-ARC <sub>l</sub> ) $([\sigma i j], B, A) \Rightarrow ([i l], B, A \cup \{(i, l, j)\})$ (RIGHT-ARC <sub>r</sub> )
<sup>1</sup>	Permitted only if $i \neq 0$ .

FIGURE 1. The **arc-standard** stack-based transition system for projective dependency parsing. The notation  $[\sigma|i]$  (for the stack) denotes a right-headed list with head  $i$  and tail  $\sigma$ ; the notation  $[j|\beta]$  (for the buffer) denotes a left-headed list with head  $j$  and tail  $\beta$ .

so the state contains all words and partially built trees (stack) as well as unseen words (buffer).

- Input function:  $\mathbf{m}(s_i)$  is the concatenation of 52 feature embeddings extracted from tokens based on their positions in the stack and the buffer.
- Recurrence function:  $\mathbf{r}(s_i)$  is empty, as this is a feed-forward network.
- RNN cell: a feed-forward MLP (so not an RNN...).

**Inference with TBRUs.** To predict the output sequence  $\{d_1, \dots, d_n\}$  given input sequence  $x = \{x_1, \dots, x_n\}$ , do:

1. Initialize  $s_1 = s^\dagger$ .
2. For  $i = 1, \dots, n$ :
  - (a) Compute  $\mathbf{h}_i = \text{RNN}(\mathbf{m}(s_i), \{\mathbf{h}_j \mid j \in \mathbf{r}(s_i)\})$ .
  - (b) Update transition state:

$$d_i \leftarrow \arg \max_{d \in \mathcal{A}(s_i)} \mathbf{w}_d^T \mathbf{h}_i \quad (349)$$

$$s_{i+1} \leftarrow t(s_i, d_i) \quad (350)$$

**NOTE:** This defines a locally normalized training procedure, whereas Andor et al. of Syntaxnet clearly conclude that their globally normalized model is the preferred choice.

---

<sup>104</sup>Image taken from “Transition-Based Parsing” by Joakim Nivre. Note that “right-headed” means “goes from left to right” or “headed to the right”.

**Combining multiple TBRUs.** We connect multiple TBRUs with different transition systems via  $\mathbf{r}(s)$ .

1. We execute a list of  $T$  TBRU components sequentially, so that each TBRU advances a global step counter.
2. *Each* transition state,  $s^\tau$ , from the  $\tau$ 'th component has access the *terminal* states from every prior transition system, and the recurrence function  $\mathbf{r}(s^\tau)$  for any given component can pull hidden activations from every prior one as well.

**Example: Multi-task bi-directional tagging.** Say we want to do both POS and NER tagging (indices start at 1).

- Left-to-right:  $\mathcal{T} = \text{shift-only}$ ,  $\mathbf{m}(s_i) = \mathbf{x}_i$ ,  $\mathbf{r}(s_i) = \{i - 1\}$ .
- Right-to-left:  $\mathcal{T} = \text{shift-only}$ ,  $\mathbf{m}(s_{n+i}) = \mathbf{x}_{(n-i)+1}$ ,  $\mathbf{r}(s_{n+i}) = \{n + i - 1\}$ .
- POS Tagger:  $\mathcal{T}_{POS} = \text{tagger}$ ,  $\mathbf{m}(s_{2n+i}) = \{\}$ ,  $\mathbf{r}(s_{2n+i}) = \{i, (2n - i) + 1\}$
- NER Tagger:  $\mathcal{T}_{NER} = \text{tagger}$ ,  $\mathbf{m}(s_{3n+i}) = \{\}$ ,  $\mathbf{r}(s_{3n+i}) = \{i, (2n - i) + 1, 2n + i\}$

which illustrates the most important aspect of the TBRU:

*A TBRU can serve as both an encoder for downstream tasks and a decoder for its own task simultaneously.*

For this example, the POS Tagger served as both an encoder for the NER task as well as a decoder for the POS task.

**Training a DRAGNN.** Assume training data consists of examples  $\mathbf{x}$  along with gold decision sequences for a given TBRU in the DRAGNN. Given decisions  $d_1, \dots, d_N$  from prior components  $1, \dots, T - 1$ , the log-likelihood for training the  $T$ 'th TBRU along its gold decision sequence  $d_{N+1}^*, \dots, d_{N+n}^*$  is then:

$$L(\mathbf{x}, d_{N+1:N+n}^*; \theta) = \sum_i \log \Pr [d_{N+i}^* | d_{1:N}, d_{N+1:N+i-1}^*; \theta] \quad (351)$$

During training, the entire input sequence is unrolled and **backpropagation through structure** is used for gradient computation.

---

#### 4.31.1 MORE DETAIL: ARC-STANDARD TRANSITION SYSTEM

---

The arc-standard transition system is mentioned a lot, but with little detail. Here I'll synthesize what I find from external resources. The literature defines the states in a transition system slightly differently than the DRAGNN paper. Here we'll define them as a **configuration**  $c = (\Sigma, B, A)$  triplet, where

- $\Sigma$  is the **stack** of tokens in  $x$  that we've [partially] processed.
- $B$  is the **buffer** of remaining tokens in  $x$  that we need to process.
- $A$  is a set of **arcs**  $(w_i, w_j, \ell)$  that link  $w_i$  to  $w_j$ , and label the arc/link as  $\ell$ .

So, in the arc-standard transition system figure presented with Parsey McParseface earlier,

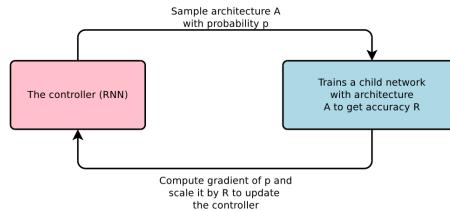
- SHIFT just means “move the head element of the buffer to the tail element of the buffer”.
- Left-arc just means “make a link *from* the tail element of the stack *to* the element before it. Remove the pointed-to element from the stack.”
- Right-arc just means “make a link *from* the element before the tail element in the stack *to* the tail element. Remove the pointed-to element from the stack.”

## Neural Architecture Search with Reinforcement Learning

Table of Contents Local

Written by Brandon McKinzie

B. Zoph and Q. Le, “Neural Architecture Search with Reinforcement Learning,” (2017).



**Controller RNN.** Generates architectures with a predefined number of layers, which is increased manually as training progresses. At convergence, validation accuracy of the generated network is recorded. Then, the controller parameters  $\theta_c$  are optimized to maximize the expected validation accuracy over a batch of generated architectures.

**Reinforcement Learning** to learn the controller parameters  $\theta_c$ . Let  $a_{1:T}$  denote a list of actions taken by the controller<sup>105</sup>, which defines a generated architecture. We denote the resulting validation accuracy by  $R$ , which is the reward signal for our RL task. Concretely, we want our controller to maximize its expected reward,  $J(\theta_c)$ :

$$J(\theta_c) = \mathbb{E}_{P(a_{1:T}; \theta_c)} [R] \quad (352)$$

Since the quantity  $\nabla_{\theta_c} R$  is non-differentiable<sup>106</sup>, we use a **policy gradient** method to iteratively update  $\theta_c$ . All this means is that we instead compute gradients over the softmax outputs (the action probabilities), and use the value of  $R$  as a simple weight factor.

$$\nabla_{\theta_c} J(\theta_c) = R \sum_{t=1}^T \mathbb{E}_{P(a_{1:T}; \theta_c)} [\nabla_{\theta_c} \log P(a_t | a_{1:t-1}; \theta_c)] \quad (353)$$

$$\approx \frac{1}{m} \sum_{k=1}^m R_k \sum_{t=1}^T \nabla_{\theta_c} \log P(a_t | a_{1:t-1}; \theta_c) \quad (354)$$

where the second equation is the empirical approximation (batch-average instead of expectation) over a batch of size  $m$ , an unbiased estimator for our gradient<sup>107</sup>. Also note that we do

<sup>105</sup>Note that  $T$  is not necessarily the number of layers, since a single generated layer can correspond to multiple actions (e.g. stride height, stride width, num filters, etc.).

<sup>106</sup> $R$  is a function of the action sequence  $a_{1:T}$  and the parameters  $\theta_c$ , and implicitly depends on the samples used for the validation set. Clearly, we do not have access to an analytical form of  $R$ , and computing numerical gradients via small perturbations of  $\theta_c$  is computationally intractable.

<sup>107</sup>It is unbiased for the same reason that any average over samples  $x$  drawn from a distribution  $P(x)$  is an unbiased estimator for  $\mathbb{E}_P [x]$ .

have access to the distribution  $P(a_{1:T}; \theta_c)$  since it is *defined* to be the joint softmax probabilities of our controller, given its parameter values  $\theta_c$  (i.e. this is not a  $p_{data}$  vs  $p_{model}$  situation). The approximation 353 is an unbiased estimator for 354, but has high variance. To reduce the variance of our estimator, the authors employ a baseline function  $b$  that does not depend on the current action:

$$\frac{1}{m} \sum_{k=1}^m \sum_{t=1}^T \nabla_{\theta_c} \log P(a_t | a_{1:t-1}; \theta_c) (R_k - b) \quad (355)$$

## Joint Extraction of Events and Entities within a Document Context

Table of Contents Local

Written by Brandon McKinzie

B. Yang and T. Mitchell, “Joint Extraction of Events and Entities within a Document Context,” (2016).

**Introduction.** Two main reasons that state-of-the-art event extraction systems have difficulties:

1. They extract events and entities in separate stages.
2. They extract events independently from each individual sentence, ignoring the rest of the document.

This paper proposes an approach that simultaneously extracts events and entities within a document context. They do this by first decomposing the problem into 3 tractable subproblems:

1. Learning the dependencies between a single event [trigger] and all of its potential arguments.
2. Learning the co-occurrence relations between events across the document.
3. Learning for entity extraction.

and then combine these learned models into a joint optimization framework.

**Learning Within-Event Structures.** For now, assume we have some document  $x$ , a set of candidate event triggers  $\mathcal{T}$ , and a set of candidate entities  $\mathcal{N}$ . Denote the set of entity candidates that are potential arguments for trigger candidate  $i$  as  $\mathcal{N}_i$ . The joint distribution over the possible trigger types, roles, and entities for those roles, is given by

$$\Pr_{\boldsymbol{\theta}} [t_i, \mathbf{r}_i, \mathbf{a} | i, \mathcal{N}_i, x] \propto \quad (356)$$

$$\exp \left( \boldsymbol{\theta}_1^T \mathbf{f}_1(t_i) + \sum_{j \in \mathcal{N}_i} [\boldsymbol{\theta}_2^T \mathbf{f}_2(r_{ij}) + \boldsymbol{\theta}_3^T f_3(t_i, r_{ij}) + \boldsymbol{\theta}_4^T \mathbf{f}_4(a_j) + \boldsymbol{\theta}_5^T f_5(r_{ij}, a_j)] \right) \quad (357)$$

All  $f_i$  also depend on  $i, x$ . In addition, all  $f_i$  except  $f_1$  depend on the current  $j$  in the summation.

where each  $f_i$  is a feature function, and I’ve colored the unary feature functions green. The unary features are tabulated in Table 1 of the original paper. They use simple indicator functions  $1_{t,r}$  and  $1_{r,a}$  for the pairwise features. They train using maximum-likelihood estimates with L2 regularization:

$$\boldsymbol{\theta}^* = \arg \max_{\boldsymbol{\theta}} \mathcal{L}(\boldsymbol{\theta}) - \lambda \|\boldsymbol{\theta}\|_2^2 \quad (358)$$

$$\mathcal{L}(\boldsymbol{\theta}) = \sum_i \log (\Pr_{\boldsymbol{\theta}} [t_i, \mathbf{r}_i, \mathbf{a} | i, \mathcal{N}_i, x]) \quad (359)$$

and use **L-BFGS** to optimize the training objective.

**Learning Event-Event Relations.** A pairwise model of event-event relations in a document. Training data consists of all pair of trigger candidates that co-occur in the same sentence or are connected by a co-referent subject/object if they’re in different sentences. Given a trigger candidate pair  $(i, i')$ , we estimate the probabilities for their event types  $(t_i, t_{i'})$  as

$$\Pr_\phi [t_i, t_{i'} | x, i, i'] \propto \exp (\phi^T g(t_i, t_{i'}, x, i, i')) \quad (360)$$

where  $g$  is a feature function that depends on the trigger candidate pair and their context. In addition to re-using the trigger features in Table 1 of the paper, they also introduce relational trigger features:

1. whether they’re connected by a conjunction dependency relation
2. whether they share a subject or an object
3. whether they have the same head word lemma
4. whether they share a semantic frame based on FrameNet.

As before, they use L-BFGS to compute the maximum-likelihood estimates of the parameters  $\phi$ .

**Entity Extraction.** Trained a standard linear-chain **CRF** using the BIO scheme. Their CRF features:

1. current words and POS tags
2. context words in a window of size 2
3. word type such as all-capitalized, is-capitalized, all-digits
4. Gazetteer-based entity type if the current word matches an entry in the gazetteers collected from Wikipedia.
5. pre-trained word2vec embeddings for each word

**Joint Inference.** Allows information flow among the 3 local models and finds globally-optimal assignments of all variables. Define the following objective:

$$\max_{\mathbf{t}, \mathbf{r}, \mathbf{a}} \sum_{i \in \mathcal{T}} E(t_i, \mathbf{r}_i, \mathbf{a}) + \sum_{i, i' \in \mathcal{T}} R(t_i, t_{i'}) + \sum_{j \in \mathcal{N}} D(a_j) \quad (361)$$

where

- The first term is the sum of confidence scores for individual event mentions from the within-event model.

$$E(t_i, \mathbf{r}_i, \mathbf{a}) = \log p_\theta(t_i) + \sum_{j \in \mathcal{N}_i} [\log p_\theta(r_{ij}) + \log p_\theta(r_{ij}, a_j)] \quad (362)$$

- The second term is the sum of confidence scores for event relations based on the pairwise event model.

$$R(t_i, t_{i'}) = \log p_\phi(t_i, t_{i'} | i, i', x) \quad (363)$$

- The third term is sum of confidence scores for entity mentions, where

$$D(a_j) = \log p_\psi(a_j | j, x) \quad (364)$$

and  $p_\psi(a_j | j, x)$  is the marginal probability derived from the linear-chain CRF.

The optimization is subject to agreement constraints that enforce the overlapping variables among the 3 components to agree on their values. The joint inference problem can be formulated as an **integer linear problem (ILP)**<sup>108</sup>. To solve it efficiently, they find solutions for the relaxation of the problem using a **dual decomposition algorithm AD<sup>3</sup>**.

---

<sup>108</sup>From Wikipedia: An integer linear program in canonical form: maximize  $c^T x$  subject to  $Ax \leq b, x \geq 0, x \in \mathbb{Z}^n$

## Globally Normalized Transition-Based Neural Networks

Table of Contents Local

Written by Brandon McKinzie

D. Andor et al., “Globally Normalized Transition-Based Neural Networks,” (2016).

**Introduction.** Authors demonstrate that simple FF neural networks can achieve comparable or better accuracies than LSTMs, as long as they are **globally normalized**. They don’t use any recurrence, but perform **beam search** for maintaining multiple hypotheses and introduce global normalization with a **CRF** objective to overcome the label bias problem that locally normalized models suffer from.

**Transition System.** Given an input sequence  $x$ , define:

- Set of states  $S(x)$ .
- Start state  $s^\dagger \in S(x)$ .
- Set of decisions  $\mathcal{A}(s, x)$  for all  $s \in S(x)$ .
- Transition function  $t(s, d, x)$  returning new state  $s'$  for any decision  $d \in A(s, x)$ .

The scoring function  $\rho(s, d; \theta)$ , which gives the score for decision  $d$  in state  $s$ , will be defined:

$$\rho(s, d; \theta) = \phi(s; \theta^{(l)}) \cdot \theta^{(d)} \quad (365)$$

which is just the familiar logits computation for decision  $d$ .  $\theta^{(l)}$  are the parameters of the network excluding the parameters at the final layer,  $\theta^{(d)}$ .  $\phi(s; \theta^{(l)})$  gives the representation for state  $s$  computed by the neural network under parameters  $\theta^{(l)}$ .

### Global vs. Local Normalization.

- **Local.** Conditional probabilities  $\Pr[d_j | s_j; \theta]$  are normalized locally over the scores for each possible action  $d_j$  from the current state  $s_j$ .

$$\Pr_L [d_{1:n}] = \prod_{j=1}^n \Pr [d_j | s_j; \theta] = \frac{\exp \left( \sum_{j=1}^n \rho(s_j, d_j; \theta) \right)}{\prod_{j=1}^n Z_L(s_j; \theta)} \quad (366)$$

$$Z_L(s_j; \theta) = \sum_{d' \in \mathcal{A}(s_j)} \exp (\rho(s_j, d'; \theta)) \quad (367)$$

Beam search can be used to attempt to find the action sequence with highest probability.

- **Global.** In contrast, a CRF defines:

$$\Pr_G [d_{1:n}] = \frac{\exp \left( \sum_{j=1}^n \rho(s_j, d_j; \theta) \right)}{Z_G(\theta)} \quad (368)$$

$$Z_G(\theta) = \sum_{d'_{1:n} \in \mathcal{D}_n} \exp \left( \sum_{j=1}^n \rho(s'_j, d'_j; \theta) \right) \quad (369)$$

where  $\mathcal{D}_n$  is the set of all valid sequences of decisions of length  $n$ . The inference problem is now to find

$$\arg \max_{d_{1:n} \in \mathcal{D}_n} \Pr_G [d_{1:n}] = \arg \max_{d_{1:n} \in \mathcal{D}_n} \sum_{j=1}^n \rho(s_j, d_j; \theta) \quad (370)$$

and we can also use beam search to approximately find the argmax.

**Training.** SGD on the NLL of the data under the model. The NLL takes a different form depending on whether we choose a locally normalized model vs a globally normalized model.

- **Local.**

$$L_{local}(d_{1:n}^*; \theta) = -\ln \Pr_L [d_{1:n}^*; \theta] \quad (371)$$

$$= -\sum_{j=1}^n \rho(s_j^*, d_j^*; \theta) + \sum_{j=1}^n \ln Z_L(s_j^*; \theta) \quad (372)$$

- **Global.**

$$L_{global}(d_{1:n}^*; \theta) = -\ln \Pr_G [d_{1:n}^*; \theta] \quad (373)$$

$$= -\sum_{j=1}^n \rho(s_j^*, d_j^*; \theta) + \ln Z_G(\theta) \quad (374)$$

To make learning tractable for the globally normalized model, the authors use **beam search with early updates**, defined as follows. Keep track of the location of the gold path<sup>109</sup> in the beam as the prediction sequence is being constructed. If the gold path is not found in the beam after step  $j$ , run one step of SGD on the following objective:

$$L_{global-beam}(d_{1:j}^*; \theta) = -\sum_{t=1}^j \rho(d_{1:t-1}^*, d_t^*; \theta) - \ln \sum_{d'_{1:j} \in \mathcal{B}_j} \exp \left( \sum_{t=1}^j \rho(d'_{1:t-1}, d'_t; \theta) \right) \quad (375)$$

where  $\mathcal{B}_j$  contains all paths in the beam at step  $j$ , and the gold path prefix  $d^* 1:j$ . If the gold path remains in the beam throughout decoding, a gradient step is performed using  $\mathcal{B}_T$ , the beam at the end of decoding. When training the global model, the authors first pretrain<sup>110</sup> using the local objective function, and then perform additional training steps using the global objective function..

---

<sup>109</sup>The gold path is the predicted sequence that matches the true labeled sequence, up to the current timestep.

**The Label Bias Problem.** Locally normalized models often have a very weak ability to revise earlier decisions. Here we will prove that **globally normalized models are strictly more expressive than locally normalized models**<sup>111</sup>. Let  $\mathcal{P}_L$  denote the set of all possible distributions  $p_L(d_{1:n} | x_{1:n})$  under the local model as the scores  $\rho$  vary. Let  $\mathcal{P}_G$  be the same, but for the global model.

**Theorem 3.1**  $\mathcal{P}_L$  is a strict subset<sup>112</sup> of  $\mathcal{P}_G$ , that is  $\mathcal{P}_L \subsetneq \mathcal{P}_G$ .

In other words, a globally normalized model can model any distribution that a locally normalized one can, but the converse is not true. I've worked through the proof below.

**Proof:**  $\mathcal{P}_L \subsetneq \mathcal{P}_G$

**Proof that  $\mathcal{P}_L \subseteq \mathcal{P}_G$ .** For any locally normalized model with scores  $\rho_L(d_{1:t-1}, d_t, x_{1:t})$ , we can define a corresponding  $p_G$  over scores

$$\rho_G(d_{1:t-1}, d_t, x_{1:t}) = \ln p_L(d_t | d_{1:t-1}, x_{1:t}) \quad (376)$$

By definition, this means that  $p_G(d_{1:t} | x_{1:t}) = p_L(d_{1:t} | x_{1:t})$ .

**Proof that  $\mathcal{P}_G \not\subseteq \mathcal{P}_L$ .** A proof by example. Consider a dataset consisting entirely of one of the following tagged sequences:

$$\mathbf{x} = abc, \quad \mathbf{d} = ABC \quad (377)$$

$$\mathbf{x} = abe, \quad \mathbf{d} = ADE \quad (378)$$

Similar to a typical linear-chain CRF, let  $\mathcal{T}$  denote the set of observed label transitions, and let  $\mathcal{E}$  denote the set of observed  $(x_t, d_t)$  pairs. Let  $\alpha$  be the single scalar parameter of this simple model, where

$$\rho(d_{1:t-1}, d_t, x_{1:t}) = \alpha (\mathbb{1}_{(d_{t-1}, d_t) \in \mathcal{T}} + \mathbb{1}_{(x_t, d_t) \in \mathcal{E}}) \quad (379)$$

for all  $t$ . This results in the following distributions  $p_G$  and  $p_L$ , evaluating on input sequence of length 3

$$p_G(d_{1:3} | x_{1:3}) = \frac{\exp(\alpha \sum_{t=1}^3 (\mathbb{1}_{(d_{t-1}, d_t) \in \mathcal{T}} + \mathbb{1}_{(x_t, d_t) \in \mathcal{E}}))}{Z_G(x_{1:3})} \quad (380)$$

$$p_L(d_{1:3} | x_{1:3}) = p_L(d_1 | x_1)p_L(d_2 | d_1, x_{1:2})p_L(d_3 | d_{1:2}, x_{1:3}) \quad (381)$$

where I've written  $p_L$  as a product over its local CPDs because it reveals the core observation that the proof is based on: *for any given subsequence  $(d_{1:t-1}, x_{1:t})$ , the local CPD is constrained to satisfy  $\sum_{d_t} p_L(d_t | d_{1:t-1}, x_{1:t}) = 1$ .* With this, the following comparison of  $p_G$  and  $p_L$  for large  $\alpha$  completes the proof of  $\mathcal{P}_G \not\subseteq \mathcal{P}_L$ :

$$\lim_{\alpha \rightarrow \infty} p_G(ABC | abc) = \lim_{\alpha \rightarrow \infty} p_G(ADE | abe) = 1 \quad (382)$$

$$p_L(ABC | abc) + p_L(ADE | abe) \leq 1 \quad (\forall \alpha) \quad (383)$$

$\therefore \mathcal{P}_L \subsetneq \mathcal{P}_G$ .

We are assuming that both  $P_L$  and  $P_G$  consist of log-linear distributions of scoring functions  $\rho(d_{1:t-1}, d_t, x_{1:t})$

<sup>111</sup>This is for conditional models only.

<sup>112</sup>Note that  $\subset$  and  $\subsetneq$  mean the same thing. Matter of notational preference/being explicit/etc.

## An Introduction to Conditional Random Fields

Table of Contents   Local

*Written by Brandon McKinzie*

Sutton et al., “An Introduction to Conditional Random Fields,” (2012).

**Graphical Modeling** (2.1). Some notation. Denote factors as  $\psi_a(\mathbf{y}_a)$  where  $1 \leq a \leq A$  and  $A$  is the total number factors.  $\mathbf{y}_a$  is an assignment to the subset  $Y_a \subseteq Y$  of variables associated with  $\psi_a$ . The value returned by  $\psi_a$  is a non-negative scalar that can be thought of as a measure of how compatible the values  $\mathbf{y}_a$  are with each other. Given a collection of subsets  $\{Y_a\}_{a=1}^A$  of  $Y$ , an **undirected graphical model** is the set of all distributions that can be written as

$$p(y) = \frac{1}{Z} \prod_{a=1}^A \psi_a(\mathbf{y}_a) \quad (384)$$

$$Z = \sum_{\mathbf{y}} \prod_{a=1}^A \psi_a(\mathbf{y}_a) \quad (385)$$

for any choice of **factors**  $\mathcal{F} = \{\psi_a\}$  that have  $\psi_a(\mathbf{y}_a) \geq 0$  for all  $\mathbf{y}_a$ . The sum for the **partition function**,  $Z$ , is over all possible assignments  $\mathbf{y}$  of the variables  $Y$ . We’ll use the term **random field** to refer to a particular distribution among those defined by an undirected model<sup>113</sup>.

We can represent the factorization with a **factor graph**: a bipartite graph  $G = (V, F, E)$  in which one set of nodes  $V = \{1, 2, \dots, |Y|\}$  indexes the RVs in the model, and the set of nodes  $F = \{1, 2, \dots, A\}$  indexes the factors. A connection between a variable node  $Y_s$  for  $s \in V$  to some factor node  $\psi_a$  for  $a \in F$  means that  $Y_s$  is one of the arguments of  $\psi_a$ . It is common to draw the factor nodes as squares, and the variable nodes as circles.

**Generative versus Discriminative Models** (2.2). Naive Bayes is generative, while logistic regression (a.k.a maximum entropy) is discriminative. Recall that Naive Bayes and logistic are defined as, respectively,

$$p(y, \mathbf{x}) = p(y) \prod_{k=1}^K p(x_k | y) \quad (386)$$

$$p(y | \mathbf{x}) = \frac{1}{Z(\mathbf{x})} \exp \left( \sum_{k=1}^K \theta_k f_k(y, \mathbf{x}) \right) \quad (387)$$

where the  $f_k$  in the definition of logistic regression denote the feature functions. We could set them, for example, as  $f_{y',j}(y, \mathbf{x}) = 1_{y'=y} x_j$ .

---

<sup>113</sup>i.e. a particular set of factors.

An example generative model for sequence prediction is the **HMM**. Recall that an HMM defines

$$p(\mathbf{y}, \mathbf{x}) = \prod_{t=1}^T p(y_t | y_{t-1})p(x_t | y_t) \quad (388)$$

where we are using the dummy notation of assuming an initial-initial state  $y_0$  clamped to 0 and begins every state sequence, so we can write the initial state distribution as  $p(y_1 | y_0)$ .

We see that the generative models, like naive Bayes and the HMM, define a family of joint distributions that factorizes as  $p(y, x) = p(y)p(x | y)$ . Discriminative models, like logistic regression, define a family of conditional distributions  $p(y | x)$ . The main conceptual difference here is that a conditional distribution  $p(y | x)$  doesn't include a model of  $p(x)$ . The principal advantage of discriminative modeling is that it's better suited to include rich, overlapping features. Discriminative models like CRFs make conditional independence assumptions both (1) among  $y$  and (2) about how the  $y$  can depend on  $x$ , but do *not* make conditional independence assumptions among  $x$ .

The difference between NB and LR is due *only* to the fact that NB is generative and LR is discriminative. Any LR classifier can be converted into a NB classifier with the same decision boundary, and vice versa. In other words, NB defines the same family as LR, if we interpret NB generatively as

$$p(y, \mathbf{x}) = \frac{\exp(\sum_k \theta_k f_k(y, \mathbf{x}))}{\sum_{\tilde{y}, \tilde{\mathbf{x}}} \exp(\sum_k \theta_k f_k(\tilde{y}, \tilde{\mathbf{x}}))} \quad (389)$$

and train it to maximize the conditional likelihood. Similarly, if the LR model is interpreted as above, and trained to maximize the joint likelihood, then we recover the same classifier as NB.

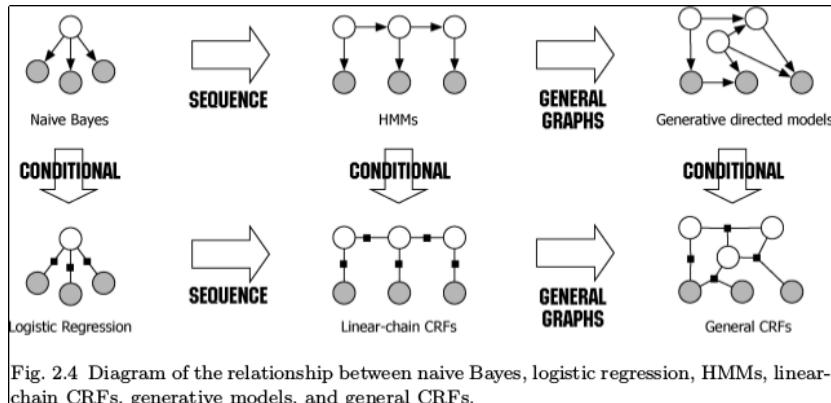


Fig. 2.4 Diagram of the relationship between naive Bayes, logistic regression, HMMs, linear-chain CRFs, generative models, and general CRFs.

**Linear-Chain CRFs** (2.3). Key point: the conditional distribution  $p(\mathbf{y} \mid \mathbf{x})$  that follows from the joint distribution  $p(\mathbf{y}, \mathbf{x})$  of an HMM is in fact a CRF with a particular choice of feature functions. First, we rewrite the HMM joint in a form that's more amenable to generalization:

$$p(\mathbf{y}, \mathbf{x}) = \frac{1}{Z} \prod_{t=1}^T \exp \left( \sum_{i,j \in S} \theta_{i,j} \mathbb{1}_{\{y_t=i, y_{t-1}=j\}} + \sum_{i \in S, o \in O} \mu_{o,i} \mathbb{1}_{\{y_t=i, x_t=o\}} \right) \quad (390) \quad \text{HMM joint distribution}$$

$$= \frac{1}{Z} \prod_{t=1}^T \exp \left( \sum_{k=1}^K \theta_k f_k(y_t, y_{t-1}, x_t) \right) \quad (391)$$

and the latter provides the more compact notation<sup>114</sup>. We can use Bayes rule to then write  $p(\mathbf{y} \mid \mathbf{x})$ , which would give us a particular kind of linear-chain CRF that only includes features for the current word's identity. The general definition of linear-chain CRFs is given below:

Let  $Y, X$  be random vectors,  $\theta = \{\theta_k\} \in \mathbb{R}^K$  be a parameter vector, and  $\mathcal{F} = \{f_k(y, y', \mathbf{x}_t)\}_{k=1}^K$  be a set of real-valued feature functions. Then a **linear-chain conditional random field** is a distribution  $p(\mathbf{y} \mid \mathbf{x})$  that takes the form:

$$p(\mathbf{y} \mid \mathbf{x}) = \frac{1}{Z(\mathbf{x})} \prod_{t=1}^T \exp \left( \sum_k \theta_k f_k(y_t, y_{t-1}, \mathbf{x}_t) \right) \quad (392)$$

$$Z(\mathbf{x}) = \sum_{\mathbf{y}} \prod_{t=1}^T \exp \left( \sum_k \theta_k f_k(y_t, y_{t-1}, \mathbf{x}_t) \right) \quad (393)$$

Notice that a linear chain CRF can be described as a factor graph over  $\mathbf{x}$  and  $\mathbf{y}$ , where each local function (factor)  $\psi_t$  has the special log-linear form:

$$\psi_t(y_t, y_{t-1}, x_t) = \exp \left( \sum_k \theta_k f_k(y_t, y_{t-1}, x_t) \right) \quad (394)$$

**General CRFs** (2.4). Let  $G$  be a factor graph over  $X$  and  $Y$ . Then  $(X, Y)$  is a conditional random field if for any value  $\mathbf{x}$  of  $X$ , the distribution  $p(\mathbf{y} \mid \mathbf{x})$  factorizes according to  $G$ . If  $F = \{\psi_a\}$  is the set of  $A$  factors in  $G$ , then the conditional distribution for a CRF is

$$p(\mathbf{y} \mid \mathbf{x}) = \frac{1}{Z(\mathbf{x})} \prod_{a=1}^A \psi_a(\mathbf{y}_a, \mathbf{x}_a) \quad (395)$$

It is often useful to require that the factors be log-linear over a prespecified set of feature functions, which allows us to write the conditional distribution as

$$p(\mathbf{y} \mid \mathbf{x}) = \frac{1}{Z(\mathbf{x})} \prod_{\psi_a \in F} \exp \left( \sum_{k=1}^{K(a)} \theta_{a,k} f_{a,k}(\mathbf{y}_a, \mathbf{x}_a) \right) \quad (396)$$

---

<sup>114</sup>Note how we collapsed the summations over  $i, j$  and  $i, o$  to simply  $k$ . This is purely notational. Each value  $k$  can be mapped to/from a unique  $i, j$  or  $i, o$  in the first version. Also note that, necessarily, each feature function  $f_k$  in the latter version maps to a specific indicator function  $\mathbb{1}$  in the first.

In addition, most models rely extensively on **parameter tying**<sup>115</sup>. To denote this, we partition the factors of  $G$  into  $\mathcal{C} = \{C_1, C_2, \dots, C_P\}$ , where each  $C_p$  is a **clique template**: a set of factors sharing a set of feature functions  $\{f_{p,k}(\mathbf{x}_c, \mathbf{y}_c)\}_{k=1}^{K(p)}$  and a corresponding set of parameters  $\boldsymbol{\theta}_p \in \mathbb{R}^{K(p)}$ . A CRF that uses clique templates can be written as

$$p(\mathbf{y} | \mathbf{x}) = \frac{1}{Z(\mathbf{x})} \prod_{C_p \in \mathcal{C}} \prod_{\psi_c \in C_p} \psi_c(\mathbf{x}_c, \mathbf{y}_c; \boldsymbol{\theta}_p) \quad (397)$$

$$= \frac{1}{Z(\mathbf{x})} \prod_{C_p \in \mathcal{C}} \prod_{c \in C_p} \exp \left\{ \sum_{k=1}^{K(p)} \theta_{p,k} f_{p,k}(\mathbf{x}_c, \mathbf{y}_c) \right\} \quad (398)$$

In a linear-chain CRF, typically one uses one clique template  $C_0 = \{\psi_t\}_{t=1}^T$ . Again, each factor in a given template *shares the same feature functions and parameters*, so the previous sentence means that we reuse the set of features and parameters for each timestep.

### Feature engineering (2.5).

- **Label-observation features.** When our label variables are discrete, the features  $f_{p,k}$  of a clique template  $C_p$  are ordinarily chosen to have a particular form:

$$f_{p,k}(\mathbf{y}_c, \mathbf{x}_c) = \mathbf{1}_{\{\mathbf{y}_c = \tilde{\mathbf{y}}_c\}} q_{p,k}(\mathbf{x}_c) \quad (399)$$

and we refer to the functions  $q_{p,k}(\mathbf{x}_c)$  as *observation functions*.

- **Unsupported features.** Many observation-label pairs may never occur in our training data (e.g. having the word “with” being associated with label “CITY”). Such features are called unsupported features, and can be useful since often their weights will be driven negative, which can help prevent the model from making predictions in the future that are far from what was seen in the training data.
- **Edge-Observation and Node-Observation features:** the two most common types of label-observation features. Edge-observation features are for the transition factors, while node-observation features are the form introduced for label-observation features above.

$$[\text{edge-obs}] \quad f(y_t, y_{t-1}, \mathbf{x}_t) = q_m(\mathbf{x}_t) \mathbf{1}_{y_t=y, y_{t-1}=y'} \quad (\forall y, y' \in \mathcal{Y}, \forall m) \quad (400)$$

$$[\text{node-obs}] \quad f(y_t, y_{t-1}, \mathbf{x}_t) = \mathbf{1}_{y_t=y, y_{t-1}=y'} \quad (\forall y, y' \in \mathcal{Y}) \quad (401)$$

and both use the same  $f(y_t, \mathbf{x}_t) = q_m(\mathbf{x}_t) \mathbf{1}_{y_t=y}$  ( $\forall y \in \mathcal{Y}, \forall m$ ). Recall that  $m$  is the index into our set of observation features<sup>116</sup>.

- **Feature Induction.** The model begins with a number of base features, and the training procedure adds conjunctions of those features.

---

<sup>115</sup>Note how, for CRFs, the actual parameters  $\theta$  are tightly coupled with the feature functions  $f$ .

<sup>116</sup>In CRFSuite, the observation features are all the attributes we define, and any features that use both label and observation are defined within CRFSuite itself.

---

### 4.35.1 INFERENCE (SEC. 4)

---

There are two inference problems that arise:

1. Wanting to predict the labels of a new input  $\mathbf{x}$  using the most likely labeling  $\mathbf{y}^* = \arg \max_{\mathbf{y}} p(\mathbf{y} | \mathbf{x})$ .
2. Computing marginal distributions (during parameter estimation, for example) such as node marginals  $p(y_t | \mathbf{x})$  and edge marginals  $p(y_t, y_{t-1} | \mathbf{x})$ .

For linear-chain CRFs, the **forward-backward algorithm** is used for computing marginals, and the **Viterbi algorithm** for computing the most probable assignment. We'll first derive these for the case of HMMs, and then generalize to the linear-chain CRF case.

**Forward-backward algorithm (HMMs).** An efficient technique for computing marginals. We begin by writing out  $p(\mathbf{x})$ , and using the distributive law to convert the sum of products to a product of sums:

$$p(\mathbf{x}) = \sum_{\mathbf{y}} p(\mathbf{x}, \mathbf{y}) \quad (402) \quad \text{We define the } \psi_t \text{ as } p(y_t | y_{t-1})p(x_t | y_t)$$

$$= \sum_{\mathbf{y}} \prod_{t=1}^T \psi_t(y_t, y_{t-1}, x_t) \quad (403)$$

$$= \sum_{y_1} p(y_1) p(x_1 | y_1) \sum_{y_2} p(y_2 | y_1) p(x_2 | y_2) \sum_{y_3} \cdots \sum_{y_T} p(y_T | y_{T-1}) p(x_T | y_T) \quad (404)$$

$$= \sum_{y_1} \psi_1(y_1, x_1) \sum_{y_2} \cdots \sum_{y_T} \psi_T(y_T, y_{T-1}, x_T) \quad (405)$$

$$= \sum_T \sum_{T-1} \psi_T(y_T, y_{T-1}, x_T) \sum_{y_{T-2}} \cdots \sum_{y_1} \psi_1(y_1, x_1) \quad (406)$$

We see that we can save an exponential amount of work by caching the inner sums as we go. Let  $M$  denote the number of possible states for the  $y$  variables. We define a set of  $T$  **forward variables**  $\alpha_t$ , each of which is a vector of size  $M$ :

$$\alpha_t(j) \triangleq p(\mathbf{x}_{\langle 1 \dots t \rangle}, y_t = j) \quad (407)$$

$$= \sum_{\mathbf{y}_{\langle 1 \dots t-1 \rangle}} p(\mathbf{x}_{\langle 1 \dots t \rangle}, \mathbf{y}_{\langle 1 \dots t-1 \rangle}, y_t = j) \quad (408)$$

$$= \sum_{\mathbf{y}_{\langle 1 \dots t-1 \rangle}} \psi_t(j, y_{t-1}, x_t) \prod_{t'=1}^{t'-1} \psi_{t'}(y_{t'}, y_{t-1}, x_{t'}) \quad (409)$$

$$= \sum_{i \in S} \psi_t(j, i, x_t) \sum_{\mathbf{y}_{\langle 1 \dots t-2 \rangle}} \psi_{t-1}(y_{t-1}, y_{t-2}, x_{t-1}) \prod_{t'=1}^{t-2} \psi_{t'}(y_{t'}, y_{t-1}, x_{t'}) \quad (410)$$

$$= \sum_{i \in S} \psi_t(j, i, x_t) \alpha_{t-1}(i) \quad (411)$$

where  $S$  is the set of  $M$  possible states. By recognizing that  $p(\mathbf{x}) = \sum_{j \in S} \sum_{\mathbf{y}_{\langle 1 \dots t-1 \rangle}} p(\mathbf{x}_{\langle 1 \dots t \rangle}, \mathbf{y}_{\langle 1 \dots t-1 \rangle}, j)$ , we can rewrite  $p(\mathbf{x})$  as

$$p(\mathbf{x}) = \sum_{j \in S} \alpha_T(j) \quad (412)$$

Notice how in the step from equation 407 to 408, we marginalized over all possible  $y$  subsequences that could've been aligned with  $\mathbf{x}_{\langle 1 \dots t \rangle}$ . We will repeat this pattern to derive the backward recursion for  $\beta_t$ , which is the same idea except now we go from  $T$  backward until some  $t$  (instead of going from 1 *forward* until some  $t$ ).

$$\beta_t(i) \triangleq p(\mathbf{x}_{\langle t+1 \dots T \rangle} \mid y_t = i) \quad (413)$$

$$= \sum_{\mathbf{y}_{\langle t+1 \dots T \rangle}} \psi_{t+1}(y_{t+1}, i, x_{t+1}) \prod_{t'=t+2}^T \psi_{t'}(y_{t'}, y_{t'-1}, x_{t'}) \quad \text{We initialize } \beta_T(i) = 1. \quad (414)$$

$$= \sum_{y_{t+1}} \psi_{t+1}(y_{t+1}, i, x_{t+1}) \beta_{t+1}(y_{t+1}) \quad (415)$$

Similar to how we obtained equation 412, we can rewrite  $p(\mathbf{x})$  in terms of the  $\beta$ :

$$p(\mathbf{x}) = \beta_0(y_0) \triangleq \sum_{y_1} \psi_1(y_1, y_0, x_1) \beta_1(y_1) \quad (416)$$

We can then combine the definition for  $\alpha$  and  $\beta$  to compute marginals of the form  $p(y_{t-1}, y_t \mid \mathbf{x})$ :

$$p(y_{t-1}, y_t \mid \mathbf{x}) = \frac{1}{p(\mathbf{x})} \alpha_{t-1}(y_{t-1}) \psi_t(y_t, y_{t-1}, x_t) \beta_t(y_t) \quad (417)$$

$$\text{where } p(\mathbf{x}) = \sum_{j \in S} \alpha_T(j) = \beta_0(y_0) \quad (418)$$

In summary, the forward-backward algorithm consists of the following steps:

1. Compute  $\alpha_t$  for all  $t$  using equation 412.
2. Compute  $\beta_t$  for all  $t$  using equation 416.
3. Return the marginal distributions computed from equation 417.

**Viterbi algorithm (HMMs).** For computing  $\mathbf{y}^* = \arg \max_{\mathbf{y}} p(\mathbf{y} \mid \mathbf{x})$ . The derivation is nearly the same as how we derived the forward-backward algorithm, except now we've replaced the summations in equation 406 with maximization. The analog of  $\alpha$  for viterbi are defined as:

$$\delta_t(j) = \max_{\mathbf{y}_{(1 \dots t-1)}} \psi_t(j, y_{t-1}, x_t) \prod_{t'=1}^{t-1} \psi_{t'}(y_{t'}, y_{t'-1}, x_{t'}) \quad (419)$$

$$= \max_{i \in S} \psi_t(j, i, x_t) \delta_{t-1}(i) \quad (420)$$

and the maximizing assignment is computed by a backwards recursion,

$$y_T^* = \arg \max_{i \in S} \delta_T(i) \quad (421)$$

$$y_t^* = \arg \max_{i \in S} \psi_t(y_{t+1}^*, i, x_{t+1}) \delta_t(i) \text{ for } t < T \quad (422)$$

Computing the recursions for  $\delta_t$  and  $y_t^*$  together is the *Viterbi algorithm*.

**Forward-backward and Viterbi for linear-chain CRF.** Generalizing to the linear-chain CRF, where now

$$p(\mathbf{y} \mid \mathbf{x}) = \frac{1}{Z(\mathbf{x})} \prod_{t=1}^T \psi_t(y_t, y_{t-1}, x_t) \quad (423)$$

$$\text{where } \psi_t(y_t, y_{t-1}, x_t) = \exp \left\{ \sum_k \theta_k f_k(y_t, y_{t-1}, x_t) \right\} \quad (424)$$

and the results for the forward-backward algorithm become

$$p(y_{t-1}, y_t \mid \mathbf{x}) = \frac{1}{Z(\mathbf{x})} \alpha_{t-1}(y_{t-1}) \psi_t(y_t, y_{t-1}, x_t) \beta_t(y_t) \quad (425)$$

$$p(y_t \mid \mathbf{x}) = \frac{1}{Z(\mathbf{x})} \alpha_t(y_t) \beta_t(y_t) \quad (426)$$

$$\text{where } Z(\mathbf{x}) = \sum_{j \in S} \alpha_T(j) = \beta_0(y_0) \quad (427)$$

Note that the interpretation is also slightly different. The  $\alpha$ ,  $\beta$ , and  $\delta$  variables should only be interpreted with the factorization formulas, and *not* as probabilities. Specifically, use

$$\alpha_t(j) = \sum_{\mathbf{y}_{(1 \dots t-1)}} \exp \left\{ \sum_k \theta_k f_k(j, y_{t-1}, x_t) \right\} \prod_{t'=1}^{t-1} \exp \left\{ \sum_k \theta_k f_k(y_{t'}, y_{t'-1}, x_{t'}) \right\} \quad (428)$$

$$\beta_t(i) = \sum_{\mathbf{y}_{(t+1 \dots T)}} \exp \left\{ \sum_k \theta_k f_k(y_{t+1}, i, x_{t+1}) \right\} \prod_{t'=t+2}^T \exp \left\{ \sum_k \theta_k f_k(y_{t'}, y_{t'-1}, x_{t'}) \right\} \quad (429)$$

$$\delta_t(j) = \max_{\mathbf{y}_{(1 \dots t-1)}} \exp \left\{ \sum_k \theta_k f_k(j, y_{t-1}, x_t) \right\} \prod_{t'=1}^{t-1} \exp \left\{ \sum_k \theta_k f_k(y_{t'}, y_{t'-1}, x_{t'}) \right\} \quad (430)$$

**Markov Chain Monte Carlo** (MCMC). The two most popular classes of approximate inference algorithms are **Monte Carlo** algorithms and **variational** algorithms. In what follows, we drop the CRF-specific notation and refer to the more general joint distribution

$$p(\mathbf{y}) = Z^{-1} \prod_{a \in F} \psi_a(\mathbf{y}_a) \quad (431)$$

that factorizes according to some factor graph  $G = (V, F)$ . MCMC methods construct a Markov chain whose state space is the same as that of  $Y$ , and sample from this chain to approximate, e.g., the expectation of some function  $f(\mathbf{y})$  over the distribution  $p(\mathbf{y})$ . MCMC algorithms aren't commonly applied in the context of CRFs, since parameter estimation by maximum likelihood requires calculating marginals many times.

---

#### 4.35.2 PARAMETER ESTIMATION (SEC. 5)

---

**Maximum Likelihood for Linear-Chain CRFs.** Since we're modeling the conditional distribution with CRFs, we use the **conditional log likelihood**  $\ell(\boldsymbol{\theta})$  with l2-regularization:

$$\ell(\boldsymbol{\theta}) = \sum_{i=1}^N \log p(\mathbf{y}^{(i)} | \mathbf{x}^{(i)}; \boldsymbol{\theta}) - \frac{1}{2\sigma^2} \sum_{k=1}^K \theta_k^2 \quad (432)$$

$$= \sum_{i=1}^N \sum_{t=1}^T \sum_{k=1}^K \theta_k f_k(y_t^{(i)}, y_{t-1}^{(i)}, \mathbf{x}_t^{(i)}) - \sum_{i=1}^N \log Z(\mathbf{x}^{(i)}) - \frac{1}{2\sigma^2} \sum_{k=1}^K \theta_k^2 \quad (433)$$

with regularization parameter  $1/2\sigma^2$ . The partial derivatives are

$$\frac{\partial \ell}{\partial \theta_k} = \sum_{i,t} f_k(y_t^{(i)}, y_{t-1}^{(i)}, \mathbf{x}_t^{(i)}) - \sum_{i,t,y,y'} f_k(y, y', \mathbf{x}_t^{(i)}) p(y, y' | \mathbf{x}^{(i)}) - \frac{\theta_k}{\sigma^2} \quad (434)$$

and the derivation for the partial derivative of  $\log(Z(x))$  is

$$\frac{\partial \log Z(x)}{\partial \theta_k} = \frac{1}{Z(x)} \frac{\partial}{\partial \theta_k} \left[ \sum_{\mathbf{y}_{\langle 1 \dots T \rangle}} \prod_{t=1}^T \exp \left\{ \sum_k \theta_k f_k(y_t, y_{t-1}, x_t) \right\} \right] \quad (435)$$

$$= \frac{1}{Z(x)} \sum_{\mathbf{y}_{\langle 1 \dots T \rangle}} \frac{\partial}{\partial \theta_k} \exp \left\{ \sum_t \sum_k \theta_k f_k(y_t, y_{t-1}, x_t) \right\} \quad (436)$$

$$= \frac{1}{Z(x)} \sum_{\mathbf{y}_{\langle 1 \dots T \rangle}} \sum_t f_k(y_t, y_{t-1}, x_t) \exp \left\{ \sum_t \sum_k \theta_k f_k(y_t, y_{t-1}, x_t) \right\} \quad (437)$$

$$= \sum_t \sum_{y_t} \sum_{y_{t-1}} f_k(y_t, y_{t-1}, x_t) \left[ \sum_{\mathbf{y}_{\langle 1 \dots t-2 \rangle}} \sum_{\mathbf{y}_{\langle t+1 \dots T \rangle}} \frac{1}{Z(x)} \exp \left\{ \sum_t \sum_k \theta_k f_k(y_t, y_{t-1}, x_t) \right\} \right] \quad (438)$$

$$= \sum_t \sum_y \sum_{y'} f_k(y, y', x_t) p(y_t = y, y_{t-1} = y' | x) \quad (439)$$

We can rewrite this in the form of expectations. For now, let  $\tilde{p}(\mathbf{y}, \mathbf{x})$  denote the *empirical distribution*, and let  $\hat{p}(\mathbf{y} | \mathbf{x}; \theta) \tilde{p}(\mathbf{x})$  denote the *model distribution*.

$$\frac{\partial \ell}{\partial \theta_k} = \mathbb{E}_{\mathbf{x}, \mathbf{y} \sim \tilde{p}(\mathbf{y}, \mathbf{x})} \left[ \sum_t f_k(y_t, y_{t-1}, x_t) \right] - \mathbb{E}_{\mathbf{x}, \mathbf{y} \sim \hat{p}(\mathbf{y}, \mathbf{x})} \left[ \sum_t f_k(y_t, y_{t-1}, x_t) \right] \quad (440)$$

### Procedure: Training Linear-Chain CRFs

Here I summarize the main steps involved during a parameter update.

**Inference.** We need to compute the log probabilities for each instance in the dataset, under the current parameters. We will need them when evaluating  $Z(\mathbf{x})$  and the marginals  $p(y, y' | \mathbf{x})$  when computing gradients.

1. Initialize  $\alpha_1(j) = \exp\{\sum_k \theta_k f_k(j, y_0, x_1)\}$  ( $y_0$  is the fixed initial state) and  $\beta_T(i) = 1$ .
2. For all  $t$ , compute:

$$\alpha_t(j) = \sum_{i \in S} \exp \left\{ \sum_k \theta_k f_k(j, i, x_t) \right\} \alpha_{t-1}(i) \quad (441)$$

$$\beta_t(j) = \sum_{i \in S} \exp \left\{ \sum_k \theta_k f_k(i, j, x_{t+1}) \right\} \beta_{t+1}(i) \quad (442)$$

3. Compute the marginals:

$$p(y_t, y_{t-1}, | \mathbf{x}) = \frac{1}{Z(\mathbf{x})} \alpha_{t-1}(y_{t-1}) \psi_t(y_t, y_{t-1}, x_t) \beta_t(y_t) \quad (443)$$

$$p(y_t | \mathbf{x}) = \frac{1}{Z(\mathbf{x})} \alpha_t(y_t) \beta_t(y_t) \quad (444)$$

$$\text{where } Z(\mathbf{x}) = \sum_{j \in S} \alpha_T(j) = \beta_0(y_0) \quad (445)$$

**Gradients.** For each parameter  $\theta_k$ , compute:

$$\frac{\partial \ell}{\partial \theta_k} = \sum_{i,t} f_k(y_t^{(i)}, y_{t-1}^{(i)}, \mathbf{x}_t^{(i)}) - \sum_{i,t,y,y'} f_k(y, y', \mathbf{x}_t^{(i)}) p(y, y' | \mathbf{x}^{(i)}) - \frac{\theta_k}{\sigma^2} \quad (446)$$

**CRF with latent variables.** Now we have additional latent variables  $\mathbf{z}$ :

$$p(\mathbf{y}, \mathbf{z} | \mathbf{x}) = \frac{1}{Z(\mathbf{x})} \prod_t \psi_t(y_t, y_{t-1}, \mathbf{z}_t, \mathbf{x}_t) \quad (447)$$

Since we observe  $\mathbf{y}$  during training, what if we instead treat this as a CRF with both  $\mathbf{x}$  and  $\mathbf{y}$  observed?

$$p(\mathbf{z} | \mathbf{y}, \mathbf{x}) = \frac{1}{Z(\mathbf{y}, \mathbf{x})} \prod_t \psi_t(y_t, y_{t-1}, \mathbf{z}_t, \mathbf{x}_t) \quad (448)$$

$$Z(\mathbf{y}, \mathbf{x}) = \sum_{\mathbf{z}} \prod_t \psi_t(y_t, y_{t-1}, \mathbf{z}_t, \mathbf{x}_t) \quad (449)$$

We can use the same inference algorithms as usual to compute  $Z(\mathbf{y}, \mathbf{x})$ , and the key result is that we can now write

$$p(\mathbf{y} \mid \mathbf{x}) = \frac{1}{Z(\mathbf{x})} \sum_{\mathbf{z}} \prod_t \psi_t(y_t, y_{t-1}, \mathbf{z}_t, \mathbf{x}_t) = \frac{Z(\mathbf{y}, \mathbf{x})}{Z(\mathbf{x})} \quad (450)$$

Finally, we can write the gradient as<sup>117</sup>

$$\frac{\partial \ell}{\partial \theta_k} = \sum_{\mathbf{z}} p(\mathbf{z} \mid \mathbf{y}, \mathbf{x}) \frac{\partial}{\partial \theta_k} [\log p(\mathbf{y}, \mathbf{z} \mid \mathbf{x})] \quad (451)$$

$$= \sum_t \sum_{\mathbf{z}_t} \left[ p(\mathbf{z}_t \mid \mathbf{y}, \mathbf{x}) f_k(y_t, y_{t-1}, \mathbf{z}_t, \mathbf{x}_t) - \sum_{y, y'} p(\mathbf{z}_t, y, y' \mid \mathbf{x}_t) f_k(y_t, y_{t-1}, \mathbf{z}_t, \mathbf{x}_t) \right] \quad (452)$$

where I've assumed there are no connections between any  $\mathbf{z}_t$  and  $\mathbf{z}_{t' \neq t}$ .

**Stochastic Gradient Methods.** Now we compute gradients for a single example, and for a linear-chain CRF:

$$\frac{\partial \ell_i}{\partial \theta_k} = \sum_t f_k(y_t^{(i)}, y_{t-1}^{(i)}, \mathbf{x}_t^{(i)}) - \sum_{t, y, y'} f_k(y, y', \mathbf{x}_t^{(i)}) p(y, y' \mid \mathbf{x}^{(i)}) - \frac{\theta_k}{N\sigma^2} \quad (453)$$

which corresponds to parameter update (remember: we are using the LL, not the negative LL):

$$\theta^{(m)} = \theta^{(m-1)} + \alpha_m \nabla \ell_i(\theta^{(m-1)}) \quad (454)$$

where  $m$  denotes this is the  $m$ th update of the training process.

---

<sup>117</sup>This uses the trick

$$\frac{df}{d\theta} = f(\theta) \frac{d \log f}{d\theta}$$

---

### 4.35.3 RELATED WORK AND FUTURE DIRECTIONS (SEC. 6)

---

**MEMMs.** Maximum-entropy Markov models. Essentially a Markov model in which the transition probabilities are given by logistic regression. Formally, a MEMM is defined by

$$p_{MEMM}(\mathbf{y} \mid \mathbf{x}) = \prod_{t=1}^T p(y_t \mid y_{t-1}, \mathbf{x}) \quad (455)$$

$$p(y_t \mid y_{t-1}, \mathbf{x}) = \frac{\exp \left\{ \sum_{k=1}^K \theta_k f_k(y_t, y_{t-1}, \mathbf{x}_t) \right\}}{Z_t(y_{t-1}, \mathbf{x})} \quad (456)$$

$$Z_t(y_{t-1}, \mathbf{x}) = \sum_{y'} \exp \left\{ \sum_{k=1}^K \theta_k f_k(y_t, y_{t-1}, \mathbf{x}_t) \right\} \quad (457)$$

which has some important differences compared to the linear-chain CRF. Notice how maximum-likelihood training of MEMMs does *not* require performance inference over full output sequences  $\mathbf{y}$ , because  $Z_t$  is a simple sum over the labels at a single position. MEMMs, however, suffer from *label bias*, while CRFs do not.

**Bayesian CRFs.** Instead of predicting the optimal labeling  $y_{ML}^*$  for input sequence  $x$  with maximum likelihood (ML), we can instead use a fully Bayesian (B) approach, both of which are shown below for comparison:

$$y_{ML}^* \leftarrow \arg \max_y p(y \mid x; \hat{\theta}) \quad (458)$$

$$y_B^* \leftarrow \arg \max_y \mathbb{E}_{\theta \sim p(\theta|x)}(p(y \mid x; \theta)) \quad (459)$$

Unfortunately, computing the exact integral (the expectation) is usually intractable, and we must resort to approximate methods like MCMC.

## Co-sampling: Training Robust Networks for Extremely Noisy Supervision

Table of Contents Local

Written by Brandon McKinzie

Han et al., “Co-sampling: Training Robust Networks for Extremely Noisy Supervision,” (2018).

**Introduction.** The authors state that current methodologies [for training networks under noisy labels] involves estimating the **noise transition matrix** (which they don’t define). Patrini et al. (2017) define the matrix as follows:

Denote by  $T \in [0, 1]^{c \times c}$  the noise transition matrix specifying the probability of one label being flipped to another, so that  $\forall i, j \quad T_{ij} \triangleq \Pr[\tilde{y} = e^j | y = e^i]$ . The matrix is row-stochastic<sup>118</sup> and not necessarily symmetric across the classes.

**Algorithm.** Authors propose a learning paradigm called **Co-sampling**. They maintain two networks  $f_{w_1}$  and  $f_{w_2}$  simultaneously. For each mini-batch data  $\tilde{\mathcal{D}}$ , each network selects  $\mathcal{R}_T$  small-loss instances as a “clean” mini-batch  $\hat{\mathcal{D}}_1$  and  $\hat{\mathcal{D}}_2$ , respectively. Each of the two networks then uses the clean mini-batch data to update the parameters  $w_2$  ( $w_1$ ) of its *peer* network.

- Why small-loss instances? Because deep networks tend to fit clean instances first, then noisy/harder instances progressively after.
- Why two networks? Because if we just trained a single network on clean instances, we would not be robust in extremely high-noise rates, since the training error would accumulate if the selected instances are not “fully clean.”

The Co-sampling paradigm algorithm is shown below.

```

Input:  $w_1$  and  $w_2$ ; learning rate  $\eta$ ; fixed  $\mathfrak{D}$ ; epoch  $T_k$  and  $T_{max}$ ; iteration  $N_{max}$ ;
for  $T = 1, 2, \dots, T_{max}$  do
    Shuffle: training set  $\tilde{\mathcal{D}}$                                 //noisy dataset
    for  $N = 1, \dots, N_{max}$  do
        Draw: mini-batch  $\hat{\mathcal{D}}$  from  $\tilde{\mathcal{D}}$ 
        Sample:  $\hat{\mathcal{D}}_1 = \arg \min_{\hat{\mathcal{D}}} \ell(f_{w_1}, \hat{\mathcal{D}}, \mathfrak{R}_T)$       //sample  $\mathfrak{R}_T$  small-loss instances
        Sample:  $\hat{\mathcal{D}}_2 = \arg \min_{\hat{\mathcal{D}}} \ell(f_{w_2}, \hat{\mathcal{D}}, \mathfrak{R}_T)$       //sample  $\mathfrak{R}_T$  small-loss instances
        Update:  $w_1 = w_1 - \eta \nabla f_{w_1}(\hat{\mathcal{D}}_2)$                       //update  $w_1$  by  $\hat{\mathcal{D}}_2$ 
        Update:  $w_2 = w_2 - \eta \nabla f_{w_2}(\hat{\mathcal{D}}_1)$                       //update  $w_2$  by  $\hat{\mathcal{D}}_1$ 
    end
    Update:  $\mathfrak{R}_T = 1 - \min\left\{\frac{T}{T_k} \mathfrak{D}, \mathfrak{D}\right\}$ 
end
Output:  $f_{w_1}$  and  $f_{w_2}$ 
```

<sup>118</sup>Each row sums to 1.

## Hidden-Unit Conditional Random Fields

Table of Contents Local

Written by Brandon McKinzie

Maaten et al., “Hidden-Unit Conditional Random Fields,” (2011).

**Introduction.** Three key advantages of CRFs over HMMs:

1. CRFs don’t assume that the observations are conditionally independent given the hidden (or target if linear-chain CRF) states.
2. CRFs don’t suffer from the label bias problems of models that do local probability normalization<sup>119</sup>.
3. For certain choices of factors, the negative conditional L.L. is convex.

The hidden-unit CRF (HUCRF), similar to discriminative restricted Boltzmann machines (RBMs), has binary stochastic hidden units that are conditionally independent given the data and the label sequence. By exploiting the conditional independence properties, we can efficiently compute:

1. The exact gradient of the C.L.L.
2. The most likely label sequence.
3. The marginal distributions over label sequences.

**Hidden-Unit CRFs.** At each time step  $t$ , the HUCRF employs  $H$  binary stochastic hidden units  $\mathbf{z}_t$ . It models the conditional distribution as

$$p(\mathbf{y} \mid \mathbf{x}) = \frac{1}{Z(\mathbf{x})} \sum_{\mathbf{z}} \exp(E(\mathbf{x}, \mathbf{z}, \mathbf{y})) \quad (460)$$

$$\begin{aligned} E(\mathbf{x}, \mathbf{z}, \mathbf{y}) = & \sum_{t=2}^T [\mathbf{y}_{t-1}^T \mathbf{A} \mathbf{y}_t] + \sum_{t=1}^T [\mathbf{x}_t^T \mathbf{W} \mathbf{z}_t + \mathbf{z}_t^T \mathbf{V} \mathbf{y}_t + \mathbf{b}^T \mathbf{z}_t + \mathbf{c}^T \mathbf{y}_t] \\ & + \mathbf{y}_1^T \boldsymbol{\pi} + \mathbf{y}_T^T \boldsymbol{\tau} \end{aligned} \quad (461)$$

Since the hidden units are conditionally independent given the data and labels, the hidden units can be marginalized out one-by-one. This, along with the nice property that the hidden units have binary elements, allows us to write  $p(\mathbf{y} \mid \mathbf{x})$  without writing any  $\mathbf{z}_t$  explicitly, as

---

<sup>119</sup>See the introduction in my CRF notes for recap of label-bias.

shown below:

$$p(\mathbf{y} \mid \mathbf{x}) = \frac{\exp\{\mathbf{y}_1^T \boldsymbol{\pi} + \mathbf{y}_T^T \boldsymbol{\tau}\}}{Z(\mathbf{x})} \prod_{t=1}^T \left[ \exp\{\mathbf{c}^T \mathbf{y}_t + \mathbf{y}_{t-1}^T \mathbf{A} \mathbf{y}_t\} \right. \\ \left. \prod_{h=1}^H \sum_{z_h \in \{0,1\}} \exp\{z_h b_h + z_h \mathbf{w}_h^T \mathbf{x}_t + z_h \mathbf{v}_h^T \mathbf{y}_t\} \right] \quad (462)$$

$$= \frac{\exp\{\mathbf{y}_1^T \boldsymbol{\pi} + \mathbf{y}_T^T \boldsymbol{\tau}\}}{Z(\mathbf{x})} \prod_{t=1}^T \left[ \exp\{\mathbf{c}^T \mathbf{y}_t + \mathbf{y}_{t-1}^T \mathbf{A} \mathbf{y}_t\} \right. \\ \left. \prod_{h=1}^H \left( 1 + \exp\{b_h + \mathbf{w}_h^T \mathbf{x}_t + \mathbf{v}_h^T \mathbf{y}_t\} \right) \right] \quad (463)$$

For inference, we'll need an algorithm for computing the marginals  $p(y_t \mid \mathbf{x})$  and  $p(y_t, y_{t-1} \mid \mathbf{x})$ . The equations are essentially the same as the forward-backward formulas for the linear-chain CRF, but with summations over  $\mathbf{z}$ :

$$p(y_t, y_{t-1} \mid \mathbf{x}) \propto \alpha_{t-1}(y_{t-1}) \left[ \sum_{\mathbf{z}_t} \Psi_t(\mathbf{x}_t, \mathbf{z}_t, y_t, y_{t-1}) \right] \beta_t(y_t) \quad (464)$$

$$p(y_t \mid \mathbf{x}) \propto \alpha_t(y_t) \beta_t(y_t) \quad (465)$$

$$(466)$$

$$\alpha_t(j) = \sum_{i \in \mathcal{Y}} \sum_{\mathbf{z}_t} \Psi_t(\mathbf{x}_t, \mathbf{z}_t, j, i) \alpha_{t-1}(i) \quad (467)$$

$$\beta_t(j) = \sum_{i \in \mathcal{Y}} \sum_{\mathbf{z}_{t+1}} \Psi_{t+1}(\mathbf{x}_{t+1}, \mathbf{z}_{t+1}, i, j) \beta_{t+1}(i) \quad (468)$$

**Training.** The conditional log likelihood for a single example  $(\mathbf{x}, \mathbf{y})$  is (bias and initial-state terms omitted)

$$\mathcal{L} = \log p(\mathbf{y} \mid \mathbf{x}) \quad (469)$$

$$= \sum_{t=1}^T \log \left( \sum_{\mathbf{z}_t} \Psi_t(\mathbf{x}_t, \mathbf{z}_t, y_{t-1}, y_t) \right) - \log Z(\mathbf{x}) \quad (470)$$

$$\text{where } \Psi_t := \exp\{y_{t-1}^T \mathbf{A} \mathbf{y}_t + \mathbf{x}_t^T \mathbf{W} \mathbf{z}_t + \mathbf{z}_t^T \mathbf{V} \mathbf{y}_t\} \quad (471)$$

Let  $\Upsilon = \{\mathbf{W}, \mathbf{V}, \mathbf{b}, \mathbf{c}, \}$  be the set of model parameters. The gradient w.r.t. the data-dependent<sup>120</sup> parameters  $v \in \Upsilon$  is given by

$$\frac{\partial \mathcal{L}}{\partial v} = \sum_{t=1}^T \left[ \sum_{k \in \mathcal{Y}} \left( (\mathbb{1}_{y_t=k} - p(y_t=k \mid \mathbf{x})) \sum_{h=1}^H \sigma(o_{hk}(\mathbf{x}_t)) \frac{\partial o_{hk}(\mathbf{x}_t)}{\partial v} \right) \right] \quad (472)$$

$$\text{where } o_{hk}(\mathbf{x}_t) = b_h + c_k + V_{hk} + \mathbf{w}_h^T \mathbf{x}_t \quad (473)$$

Unfortunately, the negative CLL is **non-convex**, and so we are only guaranteed to converge to a *local* maximum of the CLL.

---

<sup>120</sup>The data-dependent parameters are each individual element of the elements of  $\Upsilon$ . “Data” here means  $(\mathbf{x}, \mathbf{y})$ . Notice that  $\Upsilon$  does not include  $\mathbf{A}$ ,  $\boldsymbol{\pi}$ , or  $\boldsymbol{\tau}$ .

---

#### 4.37.1 DETAILED DERIVATIONS

---

Unfortunately, the paper leaves out a lot of details regarding derivations and implementations. I'm going to work through them here. First, a recap of the main equations, and with all biases/initial states included. Not leaving anything out<sup>121</sup>

$$p(\mathbf{y} \mid \mathbf{x}) = \frac{\exp\{\mathbf{y}_1^T \boldsymbol{\pi} + \mathbf{y}_T^T \boldsymbol{\tau}\}}{Z(\mathbf{x})} \prod_{t=1}^T \left[ \exp\{\mathbf{c}^T \mathbf{y}_t + \mathbf{y}_{t-1}^T \mathbf{A} \mathbf{y}_t\} \prod_{h=1}^H (1 + \exp\{b_h + \mathbf{w}_h^T \mathbf{x}_t + \mathbf{v}_h^T \mathbf{y}_t\}) \right] \quad (474)$$

$$NLL = - \sum_{i=1}^N \log p(\mathbf{y}^{(i)} \mid \mathbf{x}^{(i)}) \quad (475)$$

$$= - \sum_{i=1}^N \left[ \sum_{t=1}^T \log \left( \sum_{\mathbf{z}_t} \psi_t(\mathbf{x}_t, \mathbf{z}_t, \mathbf{y}_{t-1}, \mathbf{y}_t) \right) - \log Z(\mathbf{x}^{(i)}) \right] \quad (476)$$

The above formula for  $p(\mathbf{y} \mid \mathbf{x})$  implies something that will be very useful:

$$\sum_{\mathbf{z}_t} \psi_t(\mathbf{x}_t, \mathbf{z}_t, \mathbf{y}_{t-1}, \mathbf{y}_t) = \exp\{\mathbf{c}^T \mathbf{y}_t + \mathbf{y}_{t-1}^T \mathbf{A} \mathbf{y}_t\} \prod_{h=1}^H (1 + \exp\{b_h + \mathbf{w}_h^T \mathbf{x}_t + \mathbf{v}_h^T \mathbf{y}_t\}) \quad (477)$$

Using the generalization of the product rule for derivatives over N products, we can derive that

$$\frac{\partial}{\partial v} \prod_h (1 + \exp\{o(h)\}) = \left[ \prod_h (1 + \exp\{o(h)\}) \right] \left[ \sum_h \sigma(o(h)) \frac{\partial o(h)}{\partial v} \right] \quad (478)$$

Which means the derivatives of  $\sum_z \psi$  for the data-dependent params  $v_{dat}$  and transition params  $v_{tr}$ , are:

$$\frac{\partial \sum_{\mathbf{z}_t} \psi_t}{\partial v_{dat}} = \left[ \sum_h \sigma(o(h, y_t)) \frac{\partial o(h, y_t)}{\partial v_{dat}} \right] \sum_{\mathbf{z}_t} \psi_t \quad (479)$$

$$\frac{\partial \sum_{\mathbf{z}_t} \psi_t}{\partial v_{tr}} = \left[ \frac{\partial}{\partial v_{tr}} \mathbf{c}^T \mathbf{y}_t + \mathbf{y}_{t-1}^T \mathbf{A} \mathbf{y}_t \right] \sum_{\mathbf{z}_t} \psi_t \quad (480)$$

which also conveniently means that

$$\frac{\partial}{\partial v_{dat}} \log \left( \sum_{\mathbf{z}_t} \psi_t \right) = \sum_h \sigma(o(h, y_t)) \frac{\partial o(h, y_t)}{\partial v_{dat}} \quad (481)$$

$$\frac{\partial}{\partial v_{tr}} \log \left( \sum_{\mathbf{z}_t} \psi_t \right) = \frac{\partial}{\partial v_{tr}} \mathbf{c}^T \mathbf{y}_t + \mathbf{y}_{t-1}^T \mathbf{A} \mathbf{y}_t \quad (482)$$

I'll now proceed to derive the gradients of negative (conditional) log-likelihood for the main parameters. We can save some time by getting the base formula for any of the gradients with

---

<sup>121</sup>The equation for  $p(\mathbf{y} \mid \mathbf{x})$  from the paper, and thus here, is technically incorrect. The term  $\exp\{\mathbf{c}^T \mathbf{y}_t + \mathbf{y}_{t-1}^T \mathbf{A} \mathbf{y}_t\}$  should not be included in the product over  $t$  for  $t = 1$ .

respect to a specific single parameter  $v$ :

$$\frac{\partial NLL}{\partial v} = - \sum_{i=1}^N \left[ \sum_{t=1}^T \frac{\partial}{\partial v} \log \left( \sum_{\mathbf{z}_t} \psi_t(\mathbf{x}_t^{(i)}, \mathbf{z}_t, \mathbf{y}_{t-1}^{(i)}, \mathbf{y}_t^{(i)}) \right) - \frac{\partial}{\partial v} \log Z(\mathbf{x}^{(i)}) \right] \quad (483)$$

$$\frac{\partial \log Z(\mathbf{x}^{(i)})}{\partial v} = \frac{1}{Z(\mathbf{x}^{(i)})} \frac{\partial}{\partial v} \sum_{\mathbf{y}_{\langle 1 \dots T \rangle}} \tilde{p}(\mathbf{y}_1, \dots, \mathbf{y}_T \mid \mathbf{x}^{(i)}) \quad (484)$$

$$\frac{\partial}{\partial v} \tilde{p}(\mathbf{y}_1, \dots, \mathbf{y}_T \mid \mathbf{x}^{(i)}) = \frac{\partial}{\partial v} \prod_t \sum_{\mathbf{z}_t} \psi_t \quad (485)$$

$$= \left[ \prod_t \sum_{\mathbf{z}_t} \psi_t \right] \left[ \sum_t \frac{\frac{\partial}{\partial v} \sum_{\mathbf{z}_t} \psi_t}{\sum_{\mathbf{z}_t} \psi_t} \right] \quad (486)$$

where I've done some regrouping on the last line to be more gradient-friendly.

### Data-dependent parameters

All params  $v$  that are not transition params.

$$\frac{\partial NLL}{\partial v_{dat}} = - \sum_{i=1}^N \left[ \sum_{t=1}^T \frac{\partial}{\partial v_{dat}} \log \left( \sum_{\mathbf{z}_t} \psi_t \right) - \frac{\partial}{\partial v} \log Z(\mathbf{x}^{(i)}) \right] \quad (487)$$

$$= - \sum_{i=1}^N \left[ \sum_t \sum_h \sigma(o(h, y_t)) \frac{\partial o(h, y_t)}{\partial v_{dat}} - \frac{1}{Z(\mathbf{x}^{(i)})} \sum_{\mathbf{y}_{\langle 1 \dots T \rangle}} \left[ \prod_t \sum_{\mathbf{z}_t} \psi_t \right] \left[ \sum_t \frac{\frac{\partial}{\partial v} \sum_{\mathbf{z}_t} \psi_t}{\sum_{\mathbf{z}_t} \psi_t} \right] \right] \quad (488)$$

$$= - \sum_{i=1}^N \left[ \sum_t \sum_h \sigma(o(h, y_t)) \frac{\partial o(h, y_t)}{\partial v_{dat}} - \frac{1}{Z(\mathbf{x}^{(i)})} \sum_{\mathbf{y}_{\langle 1 \dots T \rangle}} \left[ \prod_t \sum_{\mathbf{z}_t} \psi_t \right] \left[ \sum_t \sum_h \sigma(o(h, y_t)) \frac{\partial o(h, y_t)}{\partial v_{dat}} \right] \right] \quad (489)$$

$$= - \sum_{i=1}^N \left[ \sum_t \sum_h \sigma(o(h, y_t)) \frac{\partial o(h, y_t)}{\partial v_{dat}} - \sum_t \sum_y \sum_{y'} \left[ \sum_h \sigma(o(h, y_t)) \frac{\partial o(h, y_t)}{\partial v_{dat}} \right] \xi_{t,y,y'} \right] \quad (490)$$

$$= - \sum_{i=1}^N \left[ \sum_t \sum_h \sigma(o(h, y_t)) \frac{\partial o(h, y_t)}{\partial v_{dat}} - \sum_t \sum_y \left[ \sum_h \sigma(o(h, y_t)) \frac{\partial o(h, y_t)}{\partial v_{dat}} \right] \gamma_{t,y} \right] \quad (491)$$

$$= - \sum_{i=1}^N \left[ \sum_t \left( \sum_h \sigma(o(h, y_t)) \frac{\partial o(h, y_t)}{\partial v_{dat}} - \sum_y \left[ \sum_h \sigma(o(h, y_t)) \frac{\partial o(h, y_t)}{\partial v_{dat}} \right] \gamma_{t,y} \right) \right] \quad (492)$$

$$= - \sum_{i=1}^N \left[ \sum_t \sum_y \left( (\mathbb{1}_{y_t=y} - \gamma_{t,y}) \sum_h \sigma(o(h, y_t)) \frac{\partial o(h, y_t)}{\partial v_{dat}} \right) \right] \quad (493)$$

NOTE: Although I haven't thoroughly checked the last few steps, they are required to be true in order to match the paper's results.

## Transition parameters

$$\frac{\partial NLL}{\partial v_{tr}} = - \sum_{i=1}^N \left[ \sum_{t=1}^T \frac{\partial}{\partial v_{tr}} \log \left( \sum_{\mathbf{z}_t} \psi_t \right) - \frac{\partial}{\partial v_{tr}} \log Z(\mathbf{x}^{(i)}) \right] \quad (494)$$

$$= - \sum_i^N \left[ \sum_t \frac{\partial}{\partial v_{tr}} [\mathbf{c}^T \mathbf{y}_t + \mathbf{y}_{t-1} \mathbf{A} \mathbf{y}_t] - \frac{\partial}{\partial v_{tr}} \log Z(\mathbf{x}^{(i)}) \right] \quad (495)$$

$$= - \sum_{i=1}^N \left[ \sum_t \frac{\partial}{\partial v_{tr}} [\mathbf{c}^T \mathbf{y}_t + \mathbf{y}_{t-1} \mathbf{A} \mathbf{y}_t] - \frac{1}{Z(\mathbf{x}^{(i)})} \sum_{\mathbf{y}_{(1 \dots T)}} \left[ \prod_t \sum_{\mathbf{z}_t} \psi_t \right] \left[ \sum_t \frac{\partial}{\partial v_{tr}} [\mathbf{c}^T \mathbf{y}_t + \mathbf{y}_{t-1} \mathbf{A} \mathbf{y}_t] \right] \right] \quad (496)$$

$$= - \sum_{i=1}^N \left[ \sum_t \sum_y \left( (\mathbb{1}_{y_t=y} - \gamma_{t,y}) \frac{\partial}{\partial v_{tr}} [\mathbf{c}^T \mathbf{y}_t + \mathbf{y}_{t-1} \mathbf{A} \mathbf{y}_t] \right) \right] \quad (497)$$

## Boundary parameters

$$\frac{\partial NLL}{\partial \pi_\ell} = - \sum_{i=1}^N [\mathbb{1}_{y_1=\ell} - \gamma_{1,\ell}] \quad (498)$$

$$\frac{\partial NLL}{\partial \tau_\ell} = - \sum_{i=1}^N [\mathbb{1}_{y_T=\ell} - \gamma_{T,\ell}] \quad (499)$$

The results of each of the boxes above are summarized below, for the case of  $N = 1$  to save space.

$$\frac{\partial NLL}{\partial v_{dat}} = - \sum_t \sum_y \left( (\mathbb{1}_{y_t=y} - \gamma_{t,y}) \sum_h \sigma(o(h,y)) \frac{\partial o(h,y)}{\partial v_{dat}} \right) \quad (500)$$

$$\frac{\partial NLL}{\partial v_{tr}} = - \sum_t \sum_y \left( (\mathbb{1}_{y_t=y} - \gamma_{t,y}) \frac{\partial}{\partial v_{tr}} [\mathbf{c}^T \mathbf{y}_t + \mathbf{y}_{t-1} \mathbf{A} \mathbf{y}_t] \right) \quad (501)$$

$$\frac{\partial NLL}{\partial \pi_\ell} = - \sum_{i=1}^N [\mathbb{1}_{y_1=\ell} - \gamma_{1,\ell}] \quad (502)$$

$$\frac{\partial NLL}{\partial \tau_\ell} = - \sum_{i=1}^N [\mathbb{1}_{y_T=\ell} - \gamma_{T,\ell}] \quad (503)$$

Now I'll further go through and show how the equations simplify for each type of data-

dependent parameter.

$$\frac{\partial NLL}{\partial W_{c,h}} = - \sum_t \sum_y \left( (\mathbb{1}_{y_t=y} - \gamma_{t,y}) \sum_{h'} \sigma(o(h', y)) \frac{\partial}{\partial W_{c,h}} (b_{h'} + c_y + V_{h',y} + \mathbf{w}_{h'}^T \mathbf{x}_t) \right) \quad (504)$$

$$= - \sum_t \sum_y \left( (\mathbb{1}_{y_t=y} - \gamma_{t,y}) \sum_{h'} \sigma(o(h, y)) \mathbb{1}_{h=h'} \mathbb{1}_{c \in \mathbf{x}_t} \right) \quad (505)$$

$$= - \sum_t \sum_y (\mathbb{1}_{y_t=y} - \gamma_{t,y}) \sigma(o(h, y)) \mathbb{1}_{c \in \mathbf{x}_t} \quad (506)$$

$$\frac{\partial NLL}{\partial V_{h,y}} = - \sum_t (\mathbb{1}_{y_t=y} - \gamma_{t,y}) \sigma(o(h, y)) \quad (507)$$

$$\frac{\partial NLL}{\partial b_h} = - \sum_t \sum_y (\mathbb{1}_{y_t=y} - \gamma_{t,y}) \sigma(o(h, y)) \quad (508)$$

$$\frac{\partial NLL}{\partial c_y} = - \sum_t (\mathbb{1}_{y_t=y} - \gamma_{t,y}) \sum_h \sigma(o(h, y)) \quad (509)$$

$$(510)$$

**Alternative Approach.** The above was a bit more cumbersome than needed. I'll now derive it in an easier way.

$$p(\mathbf{y} \mid \mathbf{x}) = \frac{\exp\{\mathbf{y}_1^T \boldsymbol{\pi} + \mathbf{y}_T^T \boldsymbol{\tau}\}}{Z(\mathbf{x})} \prod_{t=1}^T \left[ \exp\{\mathbf{c}^T \mathbf{y}_t + \mathbf{y}_{t-1}^T \mathbf{A} \mathbf{y}_t\} \prod_{h=1}^H (1 + \exp\{b_h + \mathbf{w}_h^T \mathbf{x}_t + \mathbf{v}_h^T \mathbf{y}_t\}) \right] \quad (511)$$

$$= \frac{\exp\{I+T\}}{Z(\mathbf{x})} \prod_{t=1}^T \prod_{h=1}^H (1 + \exp\{b_h + \mathbf{w}_h^T \mathbf{x}_t + \mathbf{v}_h^T \mathbf{y}_t\}) \quad (512)$$

$$NLL = - \sum_{i=1}^N \log p(\mathbf{y}^{(i)} \mid \mathbf{x}^{(i)}) \quad (513)$$

$$= - \sum_{i=1}^N \left[ I + T + \sum_{t=1}^T \sum_{h=1}^H [\log(1 + \exp\{b_h + \mathbf{w}_h^T \mathbf{x}_t + \mathbf{v}_h^T \mathbf{y}_t\})] - \log Z(\mathbf{x}^{(i)}) \right] \quad (514)$$

Now, focusing on the regular log-likelihood for a single example, we have

$$\frac{\partial \mathcal{L}_i}{\partial v} = \frac{\partial}{\partial v} \log p(\mathbf{y}^{(i)}, \mathbf{x}^{(i)}) \quad (515)$$

$$= \frac{\partial}{\partial v} \left[ I + T + \sum_{t,h} \log(1 + \exp\{b_h + \mathbf{w}_h^T \mathbf{x}_t + \mathbf{v}_h^T \mathbf{y}_t\}) - \log Z(\mathbf{x}^{(i)}) \right] \quad (516)$$

$$\frac{\partial \log Z(\mathbf{x}^{(i)})}{\partial v} = \frac{1}{Z(\mathbf{x}^{(i)})} \sum_{\mathbf{y}_{\langle 1 \dots T \rangle}} \frac{\partial}{\partial v} \tilde{p}(\mathbf{y} \mid \mathbf{x}^{(i)}) \quad (517)$$

$$= \frac{1}{Z(\mathbf{x}^{(i)})} \sum_{\mathbf{y}_{\langle 1 \dots T \rangle}} \tilde{p}(\mathbf{y} \mid \mathbf{x}^{(i)}) \frac{\partial}{\partial v} \log \tilde{p}(\mathbf{y} \mid \mathbf{x}^{(i)}) \quad (518)$$

as our base formula for partial derivatives.

### Transition parameters

$$\frac{\partial \mathcal{L}_i}{\partial A_{i,j}} = \frac{\partial}{\partial A_{i,j}} \left[ I + T + \sum_{t,h} \log \left( 1 + \exp \left\{ b_h + \mathbf{w}_h^T \mathbf{x}_t + \mathbf{v}_h^T \mathbf{y}_t \right\} \right) - \log Z(\mathbf{x}^{(i)}) \right] \quad (519)$$

$$= \frac{\partial}{\partial A_{i,j}} [I + T] - \sum_{\mathbf{y}_{\langle 1 \dots T \rangle}} p(\mathbf{y} \mid \mathbf{x}^{(i)}) \frac{\partial}{\partial A_{i,j}} \left[ y_1^T \pi + y_T^T \tau + \sum_t y_{t-1} A y_t \right] \quad (520)$$

$$= \sum_t \mathbb{1}_{y_{t-1}^{(i)}=i} \mathbb{1}_{y_t^{(i)}=j} - \sum_{\mathbf{y}_{\langle 1 \dots T \rangle}} \frac{1}{Z(\mathbf{x}^{(i)})} \tilde{p}(\mathbf{y} \mid \mathbf{x}^{(i)}) \sum_{t=1}^T \mathbb{1}_{y_{t-1}=i} \mathbb{1}_{y_t=j} \quad (521)$$

$$= \sum_t \mathbb{1}_{y_{t-1}^{(i)}=i} \mathbb{1}_{y_t^{(i)}=j} - \sum_t \sum_{y_t} \sum_{y_{t-1}} \mathbb{1}_{y_{t-1}=i} \mathbb{1}_{y_t=j} \sum_{\mathbf{y}_{\langle 1 \dots t-2 \rangle}} \sum_{\mathbf{y}_{\langle t+1 \dots T \rangle}} \frac{1}{Z(\mathbf{x}^{(i)})} \tilde{p}(\mathbf{y} \mid \mathbf{x}^{(i)}) \quad (522)$$

## Pre-training of Hidden-Unit CRFs

Table of Contents Local

Written by Brandon McKinzie

Kim et al., “Pre-training of Hidden-Unit CRFs,” (2018).

**Model Definition.** The Hidden-Unit CRF (HUCRF) accepts the usual observation sequence  $\mathbf{x} = x_1, \dots, x_n$ , and associated label sequence  $\mathbf{y} = y_1, \dots, y_n$  for training. The HUCRF also has a hidden layer of binary-valued  $\mathbf{z} = z_1 \dots z_n$ . It defines the joint probability

$$p_{\theta, \gamma}(\mathbf{y}, \mathbf{z} \mid \mathbf{x}) = \frac{\exp(\boldsymbol{\theta}^T \Phi(\mathbf{x}, \mathbf{z}) + \boldsymbol{\gamma}^T \Psi(\mathbf{z}, \mathbf{y}))}{\sum_{\mathbf{z}', \mathbf{y}' \in \mathcal{Y}(\mathbf{x}, \mathbf{z}')} \exp(\boldsymbol{\theta}^T \Phi(\mathbf{x}, \mathbf{z}') + \boldsymbol{\gamma}^T \Psi(\mathbf{z}', \mathbf{y}'))} \quad (523)$$

where

- $\mathcal{Y}(\mathbf{x}, \mathbf{z})$  is the set of all possible label sequences for  $\mathbf{x}$  and  $\mathbf{z}$ .
- $\Phi(\mathbf{x}, \mathbf{z}) = \sum_{j=1}^n \phi(x_j, z_j)$
- $\Psi(\mathbf{z}, \mathbf{y}) = \sum_{j=1}^n \psi(z_j, y_{j-1}, y_j)$ .

Also note that we model  $(z_i \perp z_{i \neq i} \mid \mathbf{x}, \mathbf{y})$ .

**Pre-training HUCRFs.** Since the objective for HUCRFs is non-convex, we should choose a better initialization method than random initialization. This is where pre-training comes in, a simple 2-step approach:

1. Cluster observed tokens from  $M$  unlabeled sequences and treat the clusters as labels to train an intermediate HUCRF. Let  $C(u^{(i)})$  be the sequence of cluster assignments/labels for the unlabeled sequence  $u^{(i)}$ . We compute:

$$(\theta_1, \gamma_1) \approx \arg \max_{\theta, \gamma} \sum_{i=1}^M \log p_{\theta, \gamma}(C(u^{(i)}) \mid u^{(i)}) \quad (524)$$

2. Train a final model on the labeled data  $\{(x^{(i)}, y^{(i)})\}_{i=1}^N$ , using  $\theta_1$  as an initialization point:

$$(\theta_2, \gamma_2) \approx \arg \max_{\theta, \gamma} \sum_{i=1}^N \log p_{\theta, \gamma}(y^{(i)} \mid x^{(i)}) \quad (525)$$

Note that pre-training only defines the initialization for  $\theta$ , the parameters between  $\mathbf{x}$  and  $\mathbf{z}$ . We still train  $\gamma$ , the parameters from  $\mathbf{z}$  to  $\mathbf{y}$ , from scratch.

**Canonical Correlation Analysis** (CCA). A general technique that we will need to understand as a prerequisite for the multi-sense clustering approach (defined in the next section). Given  $n$  samples of the form  $(x^{(i)}, y^{(i)})$ , where each  $x^{(i)} \in \{0, 1\}^d$  and  $y^{(i)} \in \{0, 1\}^{d'}$ , CCA returns *projection matrices*  $A \in \mathbb{R}^{d \times k}$  and  $B \in \mathbb{R}^{d' \times k}$  that we can use to project the samples to  $k$  dimensions:

$$x \longrightarrow A^T x \quad (526)$$

$$y \longrightarrow B^T y \quad (527)$$

The CCA algorithm is outlined below.

#### Algorithm: CCA

1. Calculate  $D \in \mathbb{R}^{d \times d'}$ ,  $u \in \mathbb{R}^d$ , and  $v \in \mathbb{R}^{d'}$  as follows:

$$D_{i,j} = \sum_{l=1}^n \mathbb{1}_{x_i^{(l)}=1} \mathbb{1}_{y_j^{(l)}=1} \quad (528)$$

$$u_i = \sum_{l=1}^n \mathbb{1}_{x_i^{(l)}=1} \quad (529)$$

$$v_i = \sum_{l=1}^n \mathbb{1}_{y_i^{(l)}=1} \quad (530)$$

2. Define  $\hat{\Omega} = \text{diag}(u)^{-1/2} D \text{ diag}(v)^{-1/2}$ .
3. Calculate rank- $k$  SVD  $\hat{\Omega}$ . Let  $U \in \mathbb{R}^{d \times k}$  and  $V \in \mathbb{R}^{d' \times k}$  contain the left and right, respectively, singular vectors for the largest  $k$  singular values.
4. Return  $A = \text{diag}(u)^{-1/2} U$  and  $B = \text{diag}(v)^{-1/2} V$ .

**Multi-sense clustering.** For each word type, use CCA to create a set of context embeddings corresponding to all occurrences of that word type. Then, cluster these embeddings with  $k$ -means. Set the number of word senses  $k$  to the number of label types occurring in the labeled data.

**TODO:** finish this note

## Structured Attention Networks

[Table of Contents](#)   [Local](#)

*Written by Brandon McKinzie*

Kim et al., “Structured Attention Networks,” (2017).

**Background: Attention Networks.** Goal: produce a context  $c$  based on input sequence  $x$  and query  $q$ . We assume we have an attention distribution<sup>122</sup>  $z \sim p(z | x, q)$ . Interpret  $z$  as a categorical latent variable over  $T$  categories, where  $T = \text{len}(x)$  is the length of the input sequence. We can then compute the context,  $c = \mathbb{E}_{z \sim p(z|x,q)} [f(x, z)]$ , where  $f(x, z)$  is an annotation function<sup>123</sup>.

*We interpret the attention mechanism as taking the expectation of an annotation function  $f(x, z)$  with respect to a latent variable  $z \sim p$ , where  $p$  is parameterized to be a function of  $x$  and  $q$ .*

For comparisons later on with the traditional attention mechanism, here it is:

$$c = \sum_t^T p(z = t | x, q) \mathbf{x}_t \quad (531)$$

$$p(z = t | x, q) = \text{softmax}(\theta_t) \quad (532)$$

where usually  $x$  is the sequence of hidden state of the encoder RNN,  $q$  is the hidden state of the decoder RNN at the most recent time step,  $z$  gives the source position to be attended to, and  $\theta_t = \text{score}(x_t, q)$ .

**Structured Attention.** In a **structured attention model**,  $z$  is now a *vector* of discrete random variables  $z_1, \dots, z_m$  and the attention distribution  $p(z | x, q)$  is now modeled as a *conditional random field*, specifying the structure of the  $z$  variables. We also assume now that the annotation function  $f$  factors into clique annotation functions  $f(x, z) = \sum_C f_C(x, z_C)$ , where the summation is over the  $C$  factors,  $\psi_C$ , of the CRF. Our context vector takes the form:

$$c = \mathbb{E}_{z \sim p(z|x,q)} [f(x, z)] = \sum_C \mathbb{E}_{z \sim p(z_C|x,q)} [f_C(x, z_C)] \quad (533)$$

$$p(z | x, q) = \frac{1}{Z(x, q)} \prod_C \psi_C(z_C) \quad (534)$$

---

<sup>122</sup>Also called the “alignments”. It is the output of the softmax layer of attention scores in the majority of cases.

<sup>123</sup>In all applications I’ve seen,  $f(x, z) = x_z$ .

**Example 1: Subsequence Selection.** Let  $m = T$ , and let each  $z_t \in \{0, 1\}$  be a binary R.V. Let  $f(x, z) = \sum_t^T f_t(x, z_t) = \sum_t^T \mathbf{x}_t \mathbb{1}_{z_t=1}$ <sup>124</sup>. This yields the context vector,

$$c = \mathbb{E}_{z_1, \dots, z_T} [f(x, z)] = \sum_t^T p(z_t = 1 | x, q) \mathbf{x}_t \quad (535)$$

Although this looks similar to equation 531, we haven't yet revealed the functional form for  $p(z | x, q)$ . Two possible choices:

**Linear-Chain CRF:**  $p(z_1, \dots, z_T | x, q) = \frac{1}{Z(x, q)} \prod_t^T \psi_t(z_t, z_{t-1}) \quad (536)$

**Bernoulli:**  $p(z_1, \dots, z_T | x, q) = \prod_t^T p(z_t = 1 | x, q) = \prod_t^T \sigma(\psi_t(z_t)) \quad (537)$

The factor  $\psi$  for the CRF is NOT the same as the factor for the Bernoulli!

These show why equation 535 is fundamentally different than equation 531:

- It allows for multiple inputs (or no inputs) to be selected for a given query.
- We can incorporate structural dependencies across the  $z_t$ 's.

Also note that all methods can use potentials from the same neural network or RNN that takes  $x$  and  $q$  as input. By this we mean, for example, that we can take the same parameters we'd use when computing the scores in our attention layer, and reinterpret them as e.g. CRF parameters. Then, we can compute the marginals  $p(z_t | x)$  using the forward-backward algorithm<sup>125</sup>.

*Crucially this generalization from vector softmax to forward-backward is just a series of differentiable steps, and we can compute gradients of its output (marginals) with respect to its input (potentials), allowing the structured attention model to be trained end-to-end as part of a deep model.*

---

<sup>124</sup>Ok, so equivalently,  $z^T x$ , i.e. the indicator function can just be replaced by  $z_t$  here

<sup>125</sup>This is different than the simple softmax we usually use in an attention layer, which does not model any interdependencies between the  $z_t$ . The marginals we end up with when using the CRF originate from a joint distribution over the entire sequence  $z_1, \dots, z_T$ . This seems potentially incredibly powerful. Need to analyze in depth.

## Neural Conditional Random Fields

Table of Contents Local

Written by Brandon McKinzie

Do and Artieres, “Neural Conditional Random Fields,” (2010).

**Neural CRFs.** Essentially, we feed the input sequence  $\mathbf{x}$  through a feed-forward network whose output layer has a linear activation function. The output layer is then connected with the target variable sequence  $\mathbf{Y}$ . In other words, instead of feeding instances  $\mathbf{x}$  of the observation variables  $\mathbf{X}$ , we feed the hidden layer activations of the NN. This results in the conditional probability

$$p(\mathbf{y} \mid \mathbf{x}) \propto \prod_{c \in C} e^{-E_c(\mathbf{x}, \mathbf{y}_c, \mathbf{w})} = \prod_{c \in C} e^{\langle \mathbf{w}_c^{\mathbf{y}_c}, \Phi_c(\mathbf{x}, \mathbf{w}_{NN}) \rangle} \quad (538)$$

We can set  $\Phi_c = \Phi$  for a shared-weights approach

where

- $\mathbf{w}_{NN}$  are the weights for the NN.
- $\mathbf{w}_c^{\mathbf{y}_c}$  are the weights (for clique  $c$ ) for the CRF.
- $\Phi_c(\mathbf{x}, \mathbf{w}_{NN})$  is the output of the NN. It symbolizes the high-level feature representation of the input  $\mathbf{x}$  at clique  $c$  computed by the NN.

The authors refer to the linear output layer (containing the CRF weights) as the *energy outputs*. For the sake of writing this in more familiar notation for the linear-chain CRF case, here is the above equation translated for the case where each clique corresponds to a timestep  $t$  of the input sequence and is either a label-label clique or a state-label clique.

$$p(\mathbf{y} \mid \mathbf{x}) = \frac{1}{Z(\mathbf{x})} \prod_t^T \exp \{-E_t(\mathbf{x}, \mathbf{y}_t, \mathbf{y}_{t-1}, \mathbf{w})\} \quad (539)$$

$$= \frac{1}{Z(\mathbf{x})} \prod_t^T \exp \{-E_{loc}(\mathbf{x}, t, y_t, \mathbf{w}) - E_{trans}(\mathbf{x}, t, y_{t-1}, y_t, \mathbf{w})\} \quad (540)$$

$$(541)$$

where the authors are using a blanket  $\mathbf{w}$  to denote all model parameters<sup>126</sup>.

---

<sup>126</sup>Also note that the authors allow for utilizing the input sequence  $\mathbf{x}$  in the transition energy function,  $E_{trans}$ , although usually we implement  $E_{trans}$  using only  $y_{t-1}$  and  $y_t$ .

**Initialization & Fine-Tuning.** The hidden layers of the NN are initialized layer-by-layer in an unsupervised manner using RBMs. It's important to note that the hidden layers of the NN consist of binary units. Then, using the pre-trained hidden layers, the CRF layer is initialized by training it in the usual way, and keeping the pretrained NN weights fixed.

Next, fine-tuning is used to learn all parameters globally.

$$\frac{\partial L(\mathbf{w})}{\partial \mathbf{w}} = \frac{1}{n} \sum_{i=1}^n \frac{\partial L_i(\mathbf{w})}{\partial \mathbf{w}} \quad (542)$$

$$\frac{\partial L_i(\mathbf{w})}{\partial \mathbf{w}} = \frac{\partial L_i(\mathbf{w})}{\partial \mathbf{E}(\mathbf{x}^{(i)})} \frac{\partial \mathbf{E}(\mathbf{x}^{(i)})}{\partial \mathbf{w}} \quad (543)$$

$$[\mathbf{E}(\mathbf{x}^{(i)})]_t = E_{loc}(\mathbf{x}, t, y_t, \mathbf{w}) + E_{trans}(\mathbf{x}, t, y_{t-1}, y_t, \mathbf{w}) \quad (544)$$

$$(545)$$

where  $\frac{\partial \mathbf{E}_i}{\partial \mathbf{w}}$  is the Jacobian matrix of the NN outputs for input sequence  $\mathbf{x}^{(i)}$  w.r.t. weights  $\mathbf{w}$ . By setting  $\frac{\partial L_i(\mathbf{w})}{\partial \mathbf{E}_i}$  as backprop errors of the NN output units, we can backpropagate and get  $\frac{\partial L_i(\mathbf{w})}{\partial \mathbf{w}}$  using the chain rule over the hidden layers.

## Bidirectional LSTM-CRF Models for Sequence Tagging

Table of Contents Local

Written by Brandon McKinzie

Huang et al., “Bidirectional LSTM-CRF Models for Sequence Tagging,” (2015).

**BI-LSTM-CRF Networks.** Consider the matrix of scores  $f_\theta(\mathbf{x})$  for input sentence  $\mathbf{x}$ . The element  $[f_\theta]_{\ell,t}$  gives the score for label  $\ell$  with the  $t$ -th word. This is output by the LSTM network parameterized by  $\theta$ . We let  $[A]_{i,j}$  denote the transition score from label  $i$  to label  $j$  within the CRF. The total set of parameters is denoted  $\tilde{\theta} = \theta \cup \mathbf{A}$ . The total score for input sentence  $\mathbf{x}$  and predicted label sequence  $y$  is then

$$s(\mathbf{x}, \mathbf{y}, \tilde{\theta}) = \sum_t^T (A_{y_{t-1}, y_t} + [f_\theta]_{y_t, t}) \quad (546)$$

**Features.** The authors incorporate 3 distinct types of input features:

- **Spelling features.** Various standard lexical/syntactical features. One-hot encoded as usual.
- **Context features.** Unigram, bigram, and sometimes tri-gram features. One-hot encoded.
- **Word embeddings.** Distinct from the word features. Use a pretrained embedding for each word.

Although it's not entirely clear, it appears they concatenate all of the aforementioned features together as input to the BI-LSTM. This necessarily means they are learning an embedding for the one-hot encoded spelling and word features. They also add direct connections from the input to the CRF for the spelling and word features.

EDIT: they may actually replace the one-hot encoded word features with the word embeddings.  
Unclear.

## Relation Extraction: A Survey

Table of Contents Local

Written by Brandon McKinzie

Pawar et al., “Relation Extraction: A Survey,” (December 2017).

**Feature-based Methods.** For each entity pair  $(e_1, e_2)$ , generate a set of features and train a classifier to predict the relation<sup>127</sup>. Some useful features are shown in the figure below.

Feature Types	Example
<b>Words:</b> Words of both the mentions and all the words in between	M11_leaders, M21_Venice; B1_of, B2_Italy, B3_’s, B4_left-wing, B5_government, B6_were, B7_in
<b>Entity Types:</b> Entity types of both the mentions	E1_PERSON, E2_GPE
<b>Mention Level:</b> Mention types (NAME, NOMINAL or PRONOUN) of both the mentions	M1_NOMINAL, M2_NAME
<b>Overlap:</b> #words separating the two mentions, #other mentions in between, flags indicating whether the two mentions are in the same NP, VP or PP	7_Words_Apart, 2_Mentions_In_Between (Italy & government), Not_Same_NP, Not_Same_VP, Not_Same_PP
<b>Dependency:</b> Words, POS and chunk labels of words on which the mentions are dependent in the dependency tree, #links traversed in dependency tree to go from one mentions to another	M1W_were, M1P_VBD, M1C_VP, M2W_in, M2P_IN, M2C_PP, DepLinks_3
<b>Parse Tree:</b> Path of non-terminals connecting the two mentions in the parse tree, and the path annotated with head words	PERSON-NP-S-VP-PP-GPE, PERSON-NP:leaders-S -VP:were-PP:in-GPE

Authors found that SVMs outperform MaxEnt (logistic reg) classifiers for this task.

**Kernel methods.** Instead of explicit feature engineering, we can design kernel functions for computing similarities between representations of two relation instances<sup>128</sup> (a relation instance is a triplet of the form  $(e_1, e_2)$ ), and SVM for the classification.

---

<sup>127</sup>If there are  $N$  unique relations for our data, it is common to train the classifier to predict  $2N$  total relations, to handle both possible orderings of relation arguments.

<sup>128</sup>Recall that kernel methods are for the general optimization problem

$$\min_w \sum_{i=1}^n \text{loss} \left( \sum_{j=1}^n \alpha_j \phi(x_j)^T \phi(x_i), y_i \right) + \lambda \sum_{j=1}^n \sum_{k=1}^n \alpha_j \alpha_k \phi(x_j)^T \phi(x_k) \quad (547)$$

which changes the direct focus from feature engineering to “similarity” engineering.

One approach is the **sequence kernel**. We represent each relation instance as a sequence of feature vectors:

$$(e_1, e_2) \rightarrow (\mathbf{f}_1, \dots, \mathbf{f}_N)$$

where  $N$  might be e.g. the number of words between the two entities, and the dimension of each  $f$  is the same, and could correspond to e.g. POS tag, NER tag, etc. More formally, define the *generalized subsequence kernel*,  $K_n(s, t, \lambda)$ , that computes some number of weighted subsequences  $u$  such that

- There exist index sequences  $ii := (i_1, \dots, i_n)$  and  $jj := (j_1, \dots, j_n)$  of length  $n$  such that

$$u_i \in \mathbf{s}_i \quad \forall i \in ii \tag{548}$$

$$u_j \in \mathbf{t}_j \quad \forall j \in jj \tag{549}$$

$$\tag{550}$$

- The weight of  $u$  is  $\lambda^{l(ii)+l(jj)}$ , where  $l(x) = \max(x) - \min(x)$  and  $0 < \lambda < 1$ . Sparser (more spaced out) subsequences get lower weight.

The authors then provide the recursion formulas for  $K$ , and describe some extensions of sequence kernels for relation extraction.

**Syntactic Tree Kernels.** Structural properties of a sentence are encoded by its constituent parse tree. The tree defines the syntax of the sentence in terms of constituents such as noun phrases (NP), verb phrases (VP), prepositional phrases (PP), POS tags (NN, VB, IN, etc.) as non-terminals and actual words as leaves. The syntax is usually governed by Context Free Grammar (CFG). Constructing a constituent parse tree for a given sentence is called *parsing*. The **Convolution Parse Tree Kernel**  $K_T$  can be used for computing similarity between two syntactic trees.

Syntactic Tree Kernels

**Dependency Tree Kernels.** For grammatical relations between words in a sentence. Words are the nodes and dependency relations are the edges (in the tree), typically from dependent to parent. In the **relation instance representation**, we use the smallest subtree containing the entity pair of a given sentence. Each node is augmented with additional features like POS, chunk, entity level (name, nominal, pronoun), hypernyms, relation argument, etc. Formally, an *augmented dependency tree* is defined as a tree  $T$  where

- Each node  $t_i$  has features  $\phi(t_i) = \{v_1, \dots, v_d\}$ .
- Let  $t_i[c]$  denote all children of  $t_i$ , and let  $Pa(t_i)$  denote its parent.
- For comparison of two nodes we use:
  - **Matching function**  $m(t_i, t_j)$ : equal to 1 if some important features are shared between  $t_i$  and  $t_j$ , else 0.
  - **Similarity function**  $s(t_i, t_j)$ : returns a positive real similarity score, and defined as

Dependency Tree  
Kernels

$$s(t_i, t_j) = \sum_{v_q \in \phi(t_i)} \sum_{v_r \in \phi(t_j)} \text{Compat}(v_q, v_r) \tag{551}$$

over some compatibility function between two feature values.

Finally, we can define the overall dependency tree kernel  $K(T_1, T_2)$  for similarity between trees  $T_1$  and  $T_2$  as follows. Let  $r_i$  denote the root node of tree  $T_i$ .

$$K(T_1, T_2) = \begin{cases} 0 & \text{if } m(r_1, r_2) = 0 \\ s(r_1, r_2) + K_c(r_1[\mathbf{c}], r_2[\mathbf{c}]) & \text{otherwise} \end{cases} \quad (552)$$

$$K_c(t_i[\mathbf{c}], t_j[\mathbf{c}]) = \sum_{\mathbf{a}, \mathbf{b}, l(\mathbf{a})=l(\mathbf{b})} \lambda^{d(\mathbf{a})+d(\mathbf{b})} K(t_i[\mathbf{a}], t_j[\mathbf{b}]) \quad (553)$$

$$a_1 \leq a_2 \leq \dots \leq a_n$$

$$d(\mathbf{a}) \triangleq a_n - a_1 + 1$$

$$0 < \lambda < 1$$

The interpretation is that, whenever a pair of matching nodes is found, *all* possible matching subsequences<sup>129</sup> of their children are found. Two subsequences  $\mathbf{a}$  and  $\mathbf{b}$  are said to “match” if  $m(a_i, b_1) = 1 (\forall i < n)$ . Similar to the sequence kernel seen earlier,  $\lambda$  is a decay factor that penalizes sparser subsequences.

---

<sup>129</sup>Note that a summation over subsequences of a sequence  $\mathbf{a}$ , denoted here as  $\sum_{\mathbf{a}}$ , expands to  $\{a_1, \dots, a_n, a_1a_2, a_1a_3, \dots, a_1a_n, a_1a_2a_3, \dots, a_2a_5a_6, \dots\}$  and so on and so forth.

## Neural Relation Extraction with Selective Attention over Instances

Table of Contents Local

Written by Brandon McKinzie

Lin et al., “Neural Relation Extraction with Selective Attention over Instances,” (2016).

**Introduction.** A common distant supervision approach for RE is aligning a KB with text. For any  $(e_1, r, e_2)$  in the KB, it assumes that all text mentions of  $(e_1, e_2)$  express the relation  $r$ . Of course, this assumption will often not be true. This motivates the notion of **multi-instance learning**, wherein we predict whether a set of instances  $\{x_1, \dots, x_n\}$  (each of which contain mention(s) of  $(e_1, e_2)$ ) imply the existence of  $(e_1, r, e_2)$  being true.

**Input Representation.** Each instance sentence  $x$  is tokenized into a sequence of words. Each word  $w_i$  is transformed into a concatenation  $\mathbf{w}_i \in \mathbb{R}^d$  ( $d = d_w + 2d_{pos}$ ),

$$\mathbf{w}_i := [\text{word2Vec}(w_i); \text{dist}(w_i, e_1); \text{dist}(w_i, e_2)] \quad (554)$$

where  $\text{dist}(a, b)$  returns the [embedded] relative distance (num tokens) between a and b in the given sentence (positive integer)<sup>130</sup>.

**Convolutional Network.** We use a CNN to encode a sentence of embeddings  $\{\mathbf{w}_1, \dots, \mathbf{w}_T\}$  into a single sentence vector representation  $\mathbf{x}$ . Denote the kernel/filter/window size as  $\ell$  and the number of words in the given sentence as  $T$ . Let  $\mathbf{q}_i \in \mathbb{R}^{\ell \cdot d}$  denote the vector for the  $i$ th window,

$$\mathbf{q}_i = \mathbf{w}_{i-\ell+1:i} \quad (1 \leq i \leq T + \ell - 1) \quad (555)$$

and let  $\mathbf{Q} \in \mathbb{R}^{(T+\ell-1) \times \ell \cdot d}$  be defined such that row  $\mathbf{Q}_i = \mathbf{q}_i^T$ . It follows that, for convolution matrix  $\mathbf{W} \in \mathbb{R}^{K \times (\ell \cdot d)}$ , the output of the  $k$ th filter, and subsequent max-pooling, is

$$\mathbf{p}_k = [\mathbf{W}\mathbf{Q}^T]_k + \mathbf{b} \quad (556)$$

$$[\mathbf{x}]_k = [\max(\mathbf{p}_{k1}); \max(\mathbf{p}_{k2}); \max(\mathbf{p}_{k3})] \quad (557)$$

where we’ve divided  $\mathbf{p}_k$  into three segments, corresponding to before entity 1, middle, and after entity 2 of the given sentence. The sentence vector  $\mathbf{x} \in \mathbb{R}^{3K}$  is the concatenation of all of these, after feeding through a non-linear activation function like a ReLU.

---

<sup>130</sup>More specifically, it is an embedded representation of the relative distance. To actually implement it, you’d first shift the relative distances such that they begin at 0, and learn  $2 * \text{window\_size} + 1$  embedding vectors for each of the possible position offsets. Anything outside the window is embedded into the zero vector.

**Selective Attention over Instances.** An attention mechanism is employed over all  $n$  sentence instances  $x_i$  for some candidate entity pair  $(e_1, e_2)$ . The output is a **set vector**  $\mathbf{s}$ , a real-valued vector representation of the set of instances, where

$$\mathbf{s} = \sum_i \alpha_i \mathbf{x}_i \quad (558)$$

$$\alpha_i = \text{softmax}(\mathbf{x}_i \mathbf{A} \mathbf{r}) \quad (559)$$

is an attention-weighted sum over the instance embeddings. Note that they constrain  $\mathbf{A}$  to be diagonal. Finally, the predictive distribution is defined as

$$p(r | \mathcal{S}, (e_1, e_2) \theta) = \text{softmax}(\mathbf{M} \mathbf{s} + \mathbf{d})_r \quad (560)$$

where  $\mathcal{S}$  is the set of  $n$  sentences for the given entity pair  $(e_1, e_2)$ .

## On Herding and the Perceptron Cycling Theorem

[Table of Contents](#)   [Local](#)

*Written by Brandon McKinzie*

Gelfand et al., “On Herding and the Perceptron Cycling Theorem,” (2010).

**Introduction.** Begin with the familiar learning rule of Rosenblatt’s perceptron, after some *incorrect* prediction  $\hat{y}_i = \text{sgn}(\mathbf{w}^T \mathbf{x}_i)$ ,

$$\mathbf{w} \leftarrow \mathbf{w} + \mathbf{x}_i(y_i - \hat{y}_i) \quad (561)$$

which has the effect that a subsequent prediction on  $\mathbf{x}_i$  will (before taking the sign) be  $\|\mathbf{x}_i\|_2^2$  closer to the correct side of the hyperplane. The **perceptron cycling theorem** (PCT) states that if the data is *not* linearly separable, the weights will still remain bounded and won’t diverge to infinity. **This paper shows that the PCT implies that certain moments are conserved on average.** Formally, their result says that, for some  $N$  number of iterations over samples selected from training data (with replacement)<sup>131</sup>,

$$\left\| \frac{1}{N} \sum_i^N \mathbf{x}_i y_i - \frac{1}{N} \sum_i^N \mathbf{x}_i \hat{y}_i \right\| \sim \mathcal{O}\left(\frac{1}{N}\right) \quad (562)$$

where it’s important to remember that  $\hat{y}_i$  here is the prediction for  $\mathbf{x}_i$  when it was encountered at that training iteration. This result shows that perceptron learning generate predictions that correlate with the input attributes the same way as the true labels do, and [the correlations] converge to the sample mean with a rate of  $1/N$ . This also hints at why averaged perceptron algorithms (using the average of weights across training) makes sense, as opposed to just selected the best weights. This paper also shows that *supervised perceptron algorithms and unsupervised herding algorithms can all be derived from the PCT*.

Below are some theorems that will be used throughout the paper. Let  $\{\mathbf{w}_t\}$  be a sequence of vectors  $\mathbf{w}_t \in \mathbb{R}^D$ , each generated according to iterative updates  $\mathbf{w}_{t+1} = \mathbf{w}_t + \mathbf{v}_t$ , where  $\mathbf{v}_t$  is an element of a *finite* set  $\mathbf{V}$ , and the norm of  $\mathbf{v}_t$  is bounded:  $\max \|\mathbf{v}_t\| = R < \infty$ .

**PCT:**  $\forall t \geq 0$ : If  $\mathbf{w}_t^T \mathbf{v}_t \leq 0$ ,  $\exists M > 0$  s.t.  $\|\mathbf{w}_t - \mathbf{w}_0\| < M$ .

**Convergence Thm:** If PCT holds, then  $\|\frac{1}{T} \sum_{t=1}^T \mathbf{v}_t\| \sim \mathcal{O}\frac{1}{T}$ .

---

<sup>131</sup>Unclear whether this is only for samples that corresponded to an update, or just all samples during training.

**Herding.** Consider a fully observed Markov Random Field (MRF) over  $m$  variables, each of which can take on an integer value in the range  $[1, K]$ . In herding, our energy function and weight updates for observation  $\mathbf{x}$  (over all  $m$  variables in  $\mathcal{X}$ ),

$$E(\mathbf{x}) = -\mathbf{w}^T \phi(\mathbf{x}) \quad (563)$$

$$\mathbf{w}_{t+1} = \mathbf{w}_t + \bar{\phi} - \phi(\mathbf{x}_t^*) \quad (564)$$

$$\text{where } \bar{\phi} = \mathbb{E}_{\mathbf{x}^{(i)} \sim p_{data}} [\phi(\mathbf{x}^{(i)})] \quad (565)$$

$$\text{and } \mathbf{x}_t^* = \arg \max_{\mathbf{x}} \mathbf{w}_t^T \phi(\mathbf{x}) \quad (566)$$

What if we consider more complicated features that depend on the weights  $\mathbf{w}$ ? This situation may arise in e.g. models with hidden units  $\mathbf{z}$ , where our feature function would take the form  $\phi(\mathbf{x}, \mathbf{z})$ , and we always select  $\mathbf{z}$  via

$$\mathbf{z}(\mathbf{x}, \mathbf{w}) = \arg \max_{\mathbf{z}'} \mathbf{w}^T \phi(\mathbf{x}, \mathbf{z}') \quad (567)$$

and therefore our feature function  $\phi$  depends on weights  $\mathbf{w}$  through  $\mathbf{z}$ . In this case, our herding update terms from above take the form

$$\bar{\phi} = \mathbb{E}_{\mathbf{x}^{(i)} \sim p_{data}} [\phi(\mathbf{x}^{(i)}, \mathbf{z}(\mathbf{x}^{(i)}, \mathbf{w}))] \quad (568)$$

$$\mathbf{x}_t^*, \mathbf{z}_t^* = \arg \max_{\mathbf{x}, \mathbf{z}} \mathbf{w}_t^T \phi(\mathbf{x}, \mathbf{z}) \quad (569)$$

**Conditional Herding.** Main contribution of this paper. It's basically identical to regular herding, but now we decompose  $\mathbf{x}$  into inputs and outputs  $(\mathbf{x}, \mathbf{y})$  for interpreting in a discriminative setting. In the paper, they express  $\mathbf{w}^T \phi(\mathbf{x}, \mathbf{y}, \mathbf{z})$  identically as a discriminative RBM. The parameter update for mini-batch  $\mathcal{D}_t$  is given by

$$\mathbf{w}_{t+1} = \mathbf{w}_t + \frac{\eta}{|\mathcal{D}_t|} \sum_{(\mathbf{x}^{(i)}, \mathbf{y}^{(i)}) \in \mathcal{D}_t} (\phi(\mathbf{x}^{(i)}, \mathbf{y}^{(i)}, \mathbf{z}) - \phi(\mathbf{x}^{(i)}, \mathbf{y}^*, \mathbf{z}^*)) \quad (570)$$

## Non-Convex Optimization for Machine Learning

Table of Contents Local

Written by Brandon McKinzie

P. Jain and P. Kar, “Non-convex Optimization for Machine Learning,” (2017).

**Convex Analysis** (2.1). First, let’s summarize some definitions.

### Convex Combination

A **convex combination** of a set of  $n$  vectors  $\mathbf{x}_i \in \mathbb{R}^p$ ,  $i = 1 \dots n$  is a vector  $\mathbf{x}_\theta := \sum_{i=1}^n \theta_i \mathbf{x}_i$ , where  $\theta_i \geq 0$  and  $\sum_{i=1}^n \theta_i = 1$ .

My interp: A weighted average where the weights can be interpreted as probability mass associated with each vector.

**Convex Set.** Sets that contain all [points in] line segments that join any 2 points in the set.

A set  $\mathcal{C}$  is called a **convex set** if  $\forall \mathbf{x}, \mathbf{y} \in \mathcal{C}$  and  $\lambda \in [0, 1]$ , we have that  $(1 - \lambda)\mathbf{x} + \lambda\mathbf{y} \in \mathcal{C}$  as well.

#### Proving conv. comb. of 3 vectors $\in \mathcal{C}$ too.

After reading the definition of a convex set above, it seemed intuitive that any convex combination of points  $\in \mathcal{C}$  should also be in it as well (i.e. generalizing the pairwise definition). Let  $x, y, z \in \mathcal{C}$ . How can we prove that  $\theta_1x + \theta_2y + \theta_3z$  (where  $\theta_i$  satisfy the constraints of a convex comb.) is also in  $\mathcal{C}$ ? Here is how I ended up doing it:

- If we can prove that  $\theta_1x + \theta_2y = (1 - \theta_3)v$  for some  $v \in \mathcal{C}$ , then our work is done. This is pretty easy to show via simple arithmetic.
- Case 1: assume  $\theta_3 < 1$ , so that we can divide both sides by  $1 - \theta_3$ :

$$v = \frac{\theta_1}{1 - \theta_3}x + \frac{\theta_2}{1 - \theta_3}$$

Clearly, the two coefficients here sum 1 and satisfy the constraints of a convex combination, and therefore we know that  $v \in \mathcal{C}$ , and this case is done.

- Case 2: assume  $\theta_3 = 1$ . Well, that means  $\theta_1 = \theta_2 = 0$ . Trivially,  $z \in \mathcal{C}$  and this case is done.

### Convex Function

A continuously differentiable function  $f : \mathbb{R}^p \mapsto \mathbb{R}$  is a **convex function** if  $\forall \mathbf{x}, \mathbf{y} \in \mathbb{R}^p$ , we have that

$$f(\mathbf{y}) \geq f(\mathbf{x}) + \langle \nabla f(\mathbf{x}), \mathbf{y} - \mathbf{x} \rangle \quad (571)$$

While thinking about how to gain intuition for the above, I came across chapter 3 of “Convex Optimization” which describes this in much more detail. It’s crucial to recognize that the RHS of the inequality is the 1st-order Taylor expansion of the function  $f$  about  $\mathbf{x}$ , evaluated at  $\mathbf{y}$ . In other words, the first-order Taylor approximation is a **global underestimator** of any convex function  $f$ <sup>132</sup>.

---

<sup>132</sup>Consider what this implies about all the 1st-order gradient-based optimizers we use.

**Strongly Convex/Smooth Function.** Informally, strong convexity ensures a convex function doesn't grow too *slow*, while strong smoothness ensures a ~~convex~~<sup>133</sup> function doesn't grow too *fast*. Formally,

A continuously differentiable function is considered  $\alpha$ -strongly convex (SC) and  $\beta$ -strongly smooth (SS) if  $\forall \mathbf{x}, \mathbf{y} \in \mathbb{R}^p$  we have

$$\frac{\alpha}{2} \|\mathbf{x} - \mathbf{y}\|_2^2 \leq f(\mathbf{y}) - f(\mathbf{x}) - \langle \nabla f(\mathbf{x}), \mathbf{y} - \mathbf{x} \rangle \leq \frac{\beta}{2} \|\mathbf{x} - \mathbf{y}\|_2^2 \quad (572)$$

Considering the aforementioned 1st-order Taylor approximation interpretation, we see that  $\alpha$  determines just how much larger  $f(\mathbf{y})$  must be compared to its linear approximation. Conversely,  $\beta$  determines the upper bound for how large this discrepancy is allowed to be<sup>134</sup>.

### Exercise 2.1: SS does not imply convexity

Construct a non-convex function  $f : \mathbb{R}^p \mapsto \mathbb{R}$  that is 1-SS.

We need to find a function whose linear approximation is always more than  $\frac{1}{2}$  times the magnitude of the difference in inputs **squared**, compared to the true value. Intuitively, I'd expect any *concave* function to satisfy this, since its linear approximation is a global *overestimator* of the true value. So, for example,  $f(\mathbf{y}) = -\|\mathbf{y}\|_2^2$  would satisfy 1-SS while being non-convex.

### Lipschitz Function

A function  $f$  is **B-Lipschitz** if  $\forall \mathbf{x}, \mathbf{y} \in \mathbb{R}^p$ ,

$$|f(\mathbf{x}) - f(\mathbf{y})| \leq B \cdot \|\mathbf{x} - \mathbf{y}\|_2 \quad (573)$$

**Jensen's Inequality.** Generalizes behavior of convex functions on convex combinations<sup>135</sup>.

If  $X$  is a R.V. taking values in the domain of a convex function  $f$ , then

$$\mathbb{E}[f(X)] \geq f(\mathbb{E}[X]) \quad (574)$$

---

<sup>133</sup>Strong smoothness alone does not imply convexity.

<sup>134</sup>Notice that SC and SS are *quadratic* lower and upper bounds, respectively. This means that the allowed deltas grow as a function of the distance between  $\mathbf{x}$  and  $\mathbf{y}$ , whereas things like Lipschitzness grow linearly.

<sup>135</sup>It should be obvious that expectations are convex combinations.

**Convex Projections** (2.2). Given any closed set  $\mathcal{C} \in \mathbb{R}^p$ , the projection operator  $\Pi_{\mathcal{C}}(\cdot)$  is defined as

$$\Pi_{\mathcal{C}}(\mathbf{z}) := \arg \min_{\mathbf{x} \in \mathcal{C}} \|\mathbf{x} - \mathbf{z}\|_2 \quad (575)$$

If  $\mathcal{C}$  is a convex set, then the above reduces to a convex optimization problem. Projections onto convex sets have three particularly interesting properties. For each of them, the setup is: “For any convex set  $\mathcal{C} \subset \mathbb{R}^p$ , and any  $\mathbf{z} \in \mathbb{R}^p$ , let  $\hat{\mathbf{z}} := \Pi_{\mathcal{C}}(\mathbf{z})$ . Then  $\forall \mathbf{x} \in \mathcal{C}, \dots$ ”

- **Property-O**<sup>136</sup>:  $\|\hat{\mathbf{z}} - \mathbf{z}\|_2 \leq \|\mathbf{x} - \mathbf{z}\|_2$ . Informally: “the projection results in the point  $\hat{\mathbf{z}}$  in  $\mathcal{C}$  that is closest to the original  $\mathbf{z}$ ”. This basically just restates the optimization problem.
- **Property-I**.  $\langle \mathbf{x} - \hat{\mathbf{z}}, \mathbf{z} - \hat{\mathbf{z}} \rangle \leq 0$ . Informally: “from the perspective of  $\hat{\mathbf{z}}$ , all points  $\mathbf{x} \in \mathcal{C}$  are in the (informally) opposite direction of  $\mathbf{z}$ ”
- **Property-II**.  $\|\hat{\mathbf{z}} - \mathbf{x}\|_2 \leq \|\mathbf{z} - \mathbf{x}\|_2$ . Informally: “the projection brings the point closer to all points in  $\mathcal{C}$  compared to its original location.”

### Proving Property-I

A proof by contradiction.

1. Assume that  $\exists \mathbf{x} \in \mathcal{C}$  s.t.  $\langle \mathbf{x} - \hat{\mathbf{z}}, \mathbf{z} - \hat{\mathbf{z}} \rangle > 0$ .
2. We know that  $\hat{\mathbf{z}}$  is also in  $\mathcal{C}$ , and since  $\mathcal{C}$  is convex, then for any  $\lambda \in [0, 1]$ ,

$$\mathbf{x}_\lambda := \lambda \mathbf{x} + (1 - \lambda) \hat{\mathbf{z}} \quad (576)$$

must also be in  $\mathcal{C}$ .

3. If we can show that some value of  $\lambda$  guarantees that  $\|\mathbf{z} - \mathbf{x}_\lambda\|_2 < \|\mathbf{z} - \hat{\mathbf{z}}\|_2$ , this would directly contradict property-O, implying  $\hat{\mathbf{z}}$  is not the closest member of  $\mathcal{C}$  to  $\mathbf{z}$ . I’m not sure how to actually derive the range of  $\lambda$  values that satisfy this, though.

**(Convex) Projected Gradient Descent** (2.3). Our optimization problem is

$$\min_{\mathbf{x} \in \mathbb{R}^p} f(\mathbf{x}) \quad \text{s.t.} \quad \mathbf{x} \in \mathcal{C} \quad (577)$$

where  $\mathcal{C} \subset \mathbb{R}^p$  is a convex constraint set, and  $f : \mathbb{R}^p \mapsto \mathbb{R}$  is a convex objective function. Projected gradient descent iteratively updates the value of  $\mathbf{x}$  that minimizes  $f$  as usual, but additionally projects the current iterate (value of best  $\mathbf{x}$ ) onto  $\mathcal{C}$  at the end of each iteration. That’s the only difference.

---

<sup>136</sup>In this case only,  $\mathcal{C}$  need not be convex

---

#### 4.45.1 NON-CONVEX PROJECTED GRADIENT DESCENT (3)

---

**Non-Convex Projections** (3.1). Here we look at a few special cases where projecting onto a non-convex set can still be carried out efficiently.

- **Projecting into sparse vectors.** The set of  $s$ -sparse vectors (vectors with at most  $s$  nonzero elements) is denoted as

$$\mathcal{B}_0(s) \triangleq \{\mathbf{x} \in \mathbb{R}^p \mid \|\mathbf{x}\|_0 \leq s\} \quad (578)$$

It turns out that  $\hat{\mathbf{z}} := \Pi_{\mathcal{B}_0(s)}(\mathbf{z})$  is obtained by setting all except the top- $s$  elements of  $\mathbf{z}$  to zero.

- **Projecting into low-rank matrices.** The set of  $m \times n$  matrices with rank at most  $r$  is denoted as

$$\mathcal{B}_{rank}(r) \triangleq \{A \in \mathbb{R}^{m \times n} \mid \text{rank}(A) \leq r\} \quad (579)$$

and we want to project some matrix  $A$  onto this set,

$$\Pi_{\mathcal{B}_{rank}(r)}(A) := \arg \min_{X \in \mathcal{B}_{rank}(r)} \|A - X\|_F \quad (580)$$

This can be done efficiently via SVD on  $A$  and retaining the top  $r$  singular values and vectors.

#### Restricted Strong Convexity and Smoothness (3.2).

##### Restricted Convexity

A continuously differentiable function  $f : \mathbb{R}^p \mapsto \mathbb{R}$  is said to satisfy restricted convexity over a (possibly non-convex) region  $\mathcal{C} \subseteq \mathbb{R}^p$  if  $\forall \mathbf{x}, \mathbf{y} \in \mathcal{C}$ , we have that

$$f(\mathbf{y}) \geq f(\mathbf{x}) + \langle \nabla f(\mathbf{x}), \mathbf{y} - \mathbf{x} \rangle \quad (581)$$

and a similar rephrasing for restricted strong convexity (RSC) and restricted strong smoothness (RSS).

## Improving Language Understanding by Generative Pre-Training

Table of Contents Local

Written by Brandon McKinzie

Radford et al., “Improving Language Understanding by Generative Pre-Training,” (2018).

**Unsupervised Pre-Training** (3.1). Given unsupervised corpus of tokens  $\mathcal{U} = \{u_1, \dots, u_n\}$ , train with a standard LM objective:

$$L_1(\mathcal{U}) = \sum_i^n \log P(u_i | u_{i-k}, \dots, u_{i-1}; \Theta) \quad (582)$$

The authors use a **Transformer decoder**, i.e. literally just the decoder part of the Transformer in “Attention is all you need.”

**Supervised Fine-Tuning** (3.2). Now we have a labeled corpus  $\mathcal{C}$ , where each instance consists of a sequence of input tokens  $x^1, \dots, x^m$ , along with a label  $y$ . They just feed the inputs through the transformer until they obtain the final transformer block’s activation  $h_l^m$ , and linearly project it to output space:

$$P(y | x^1, \dots, x^m) = \text{softmax}(h_l^m W_y) \quad (583)$$

$$L_2(\mathcal{C}) = \sum_{(x,y)} \log P(y | x^1, \dots, x^m) \quad (584)$$

They also found that including a language modeling auxiliary objective helped learning,

$$L_3(\mathcal{C}) = L_2(\mathcal{C}) + \lambda L_1(\mathcal{C}) \quad (585)$$

...that’s it. Extremely simple, yet somehow effective.

## Deep Contextualized Word Representations

Table of Contents Local

Written by Brandon McKinzie

Peters et al., “Deep Contextualized Word Representations,” (2018).

**Bidirectional Language Models** (3.1). Given a sequence of  $N$  tokens, a forward LM computes the probability of the sequence via

$$p(t_1, \dots, t_N) = \prod_{k=1}^N p(t_k | t_1, \dots, t_{k-1}) \quad (586)$$

A common approach is learning context-independent token representations  $\mathbf{x}_k$  and passing these through  $L$  layers of forward LSTMs. The top layer LSTM output at step  $k$ ,  $\vec{\mathbf{h}}_{k,L}$ , is used to predict  $t_{k+1}$  with a softmax layer. The authors’ biLM combines a forward and backward LM to jointly maximize

$$\begin{aligned} \sum_{k=1}^N & \left[ \log p(t_k | t_1, \dots, t_{k-1}; \Theta_x, \vec{\Theta}_{LSTM}, \Theta_s) \right. \\ & \left. + \log p(t_k | t_{k+1}, \dots, t_N; \Theta_x, \overleftarrow{\Theta}_{LSTM}, \Theta_s) \right] \end{aligned} \quad (587)$$

and it’s important to note the shared parameters  $\Theta_x$  (token representation) and  $\Theta_s$  (output softmax).

**ELMo** (3.2). A task-specific linear combination of the intermediate representations.

$$\mathbf{ELMo}_k^{task} = \gamma^{task} \sum_{j=0}^L s_j^{task} \mathbf{h}_{k,j}^{LM} \quad (588)$$

where  $s^{task}$  are softmax-normalized weights (so the combination is convex). The authors also mention that, in some cases, it helped to apply **layer normalization** to each biLM layer before weighting.

**Using biLMs for Supervised NLP** (3.3). Given a pretrained biLM and a supervised architecture, we can learn the ELMo representations (jointly with the given supervised task) as follows.

1. Freeze the weights of the [pretrained] biLM.
2. Concatenate the token representations (e.g. GloVe) with the ELMo representation.
3. Pass the concatenated representation into the supervised architecture.

The authors found it beneficial to some dropout to ELMo, and in some cases add L2-regularization on the ELMo weights.

**Pretrained biLM Architecture** (3.4). In addition to the biLM we introduced earlier, the authors make the following changes/specifications for their pretrained biLMs:

- **residual connections** between LSTM layers<sup>137</sup>.
- Halved all embedding and hidden dimensions from the CNN-BIG-LSTM model in *Exploring the Limits of Language Modeling*.
- The  $x_k$  token representations use 2048 character n-gram convolutional filters followed by two **highway layers**.

---

<sup>137</sup>So the output of some layer, instead of being  $\text{LSTM}(x)$ , becomes  $(x + \text{LSTM}(x))$

## Exploring the Limits of Language Modeling

Table of Contents Local

Written by Brandon McKinzie

Josefina et al., “Exploring the Limits of Language Modeling,” (2016).

**NCE and Importance Sampling** (3.1). In this section, assume any  $p(w)$  is shorthand for  $p(w | \{w_{prev}\})$ .

- **Noise Contrastive Estimation** (NCE). Train a classifier to discriminate between true data (from distribution  $p_d$ ) or samples coming from some arbitrary noise distribution  $p_n$ . If these distributions were known, we could compute

$$p(Y=\text{true} | w) = \frac{p(w | \text{true})p(\text{true})}{p(w)} \quad (589)$$

$$= \frac{p_d(w)p(\text{true})}{p(w, \text{true}) + p(w, \text{false})} \quad (590)$$

$$= \frac{p_d(w)p(\text{true})}{p_d(w)p(\text{true}) + p_n(w)p(\text{false})} \quad (591)$$

$$= \frac{p_d(w)}{p_d(w) + kp_n(w)} \quad (592)$$

where  $k$  is the number of negative samples per positive word. The idea is to train a logistic classifier  $p(Y=\text{true} | w) = \sigma(\log p_{model} - \log kp_n(w))$ , then  $\text{softmax}(\log p_{model})$  is a good approx of  $p_d(w)$ .

- **Importance Sampling**. Estimates the partition function. Consider that now we have a set of  $k+1$  words  $W = \{w_1, \dots, w_{k+1}\}$ , where  $w_1$  is the word coming from the true data, and the rest are from the noise distribution. We train a multinomial logistic regression over  $k+1$  classes,

$$p(Y=i | W) = \frac{p_d(w_i)}{p_n(w_i)} \frac{1}{\sum_{i'=1}^{k+1} p_d(w_{i'})/p_n(w_{i'})} \quad (593)$$

$$\propto_Y \frac{p_d(w_i)}{p_n(w_i)} \quad (594)$$

and we end up seeing that IS is the same as NCE, except in the multiclass setting and with cross entropy loss instead of logistic loss.

**CNN Softmax** (3.2). Typically the logit for word  $w$  is given by  $z_w = h^T e_w$ , where  $h$  is often the output state of an LSTM, and  $e_w$  is a vector of parameters that could be interpreted as the word embedding for  $w$ . Instead of this, the authors propose what they call *CNN Softmax*, where we compute  $e_w = CNN(chars_w)$ . Although this makes the function mapping from  $w$  to  $e_w$  much smoother (due to the tied weights), it ends up having a hard time distinguishing between similarly spelled words that may have entirely different meanings. The authors use a correction factor, learned for each word, such that

$$z_w = h^T CNN(chars_w) + h^T M corr_w \quad (595)$$

where  $M$  projects low-dimensional  $corr_w$  back up to the dimensionality of the LSTM state  $h$ .

**Char LSTM Predictions** (3.3). To reduce the computational burden of the partition function, the authors feed the word-level LSTM state  $h$  through a character-level LSTM that predicts the target word one character at a time.

## Connectionist Temporal Classification

Table of Contents Local

Written by Brandon McKinzie

Graves et al., “Connectionist Temporal Classification: Labelling Unsegmented Sequence Data with Recurrent Neural Networks,” (2006).

**Temporal Classification.**

- **Input space:** Let  $\mathcal{X} = (\mathbb{R}^m)^*$  be the set of all sequences of  $m$  dimensional real-valued vectors.
- **Output space:** Let  $\mathcal{Z} = L^*$  be the set of all sequences of a finite vocabulary of  $L$  labels.
- **Data distribution:** Denote by  $\mathcal{D}_{\mathcal{X} \times \mathcal{Z}}$  the probability distribution over samples  $(\mathbf{x}, \mathbf{z})$ . Let  $S$  denote a set of training examples drawn from this distribution.

**From Network Outputs to Labellings** (3.1). Let  $L' = L \cup \{\epsilon\}$  denote the set of unique labels combined with the blank token  $\epsilon$ . We refer to the alignment sequences of length  $T$  (same length as  $\mathbf{x}$ ), i.e. elements of the set  $(L')^T$ , as **paths** and denote them  $\pi$ .

Now that we have alignment sequences  $\pi$ , we need to convert them to label sequences by (1) merging repeated contiguous labels, and then (2) removing blank tokens. We denote this procedure as a many-to-one map  $\mathcal{B} : L'^T \mapsto L^{\leq T}$ . We can then write the conditional posterior over possible output sequences  $\ell$ :

$$p(\ell | \mathbf{x}) = \sum_{\pi \in \mathcal{B}^{-1}(\ell)} p(\pi | \mathbf{x}) \quad (596)$$

**Constructing the Classifier** (3.2). There is no tractable algorithm for exact decoding, i.e. computing

$$h(\mathbf{x}) \triangleq \arg \max_{\ell \in L^{\leq T}} p(\ell | \mathbf{x}) \quad (597)$$

However, the following two approximate methods work well in practice:

1. **Best Path Decoding.**  $h(\mathbf{x}) \approx \mathcal{B}(\pi^*)$  where  $\pi^* = \arg \max_{\pi \in L'^T} p(\pi | \mathbf{x})$ .
2. **Prefix Search Decoding.**

**The CTC Forward-Backward Algorithm** (4.1). Define the probability of obtaining the first  $s$  output labels,  $\ell_{\langle 1 \dots s \rangle}$ , at time  $t$  as

$$\alpha_t(s) \triangleq \sum_{\substack{\pi \in L'^T \\ \mathcal{B}(\pi_{\langle 1 \dots t \rangle}) = \ell_{\langle 1 \dots s \rangle}}} \prod_{t'=1}^t y_{\pi_{t'}}^{t'} \quad (598)$$

Note that the summation here could contain duplicate  $\pi_{\langle 1 \dots t \rangle}$  that differ only in their elements beyond  $t$ .

We insert a blank token at the beginning and end of  $\ell$  and between each pair of labels, and call this augmented sequence  $\ell'$ . We have the following rules for initializing  $\alpha$  at the first output step  $t=1$ , followed by the recursion rule:

$$\alpha_1(s) = \begin{cases} y_\epsilon^1 & s=1 \\ y_{\ell_1}^1 & s=2 \\ 0 & s > 2 \end{cases} \quad (599)$$

$$\alpha_t(s) = \begin{cases} \bar{\alpha}_t(s)y_{\ell'_s}^t & \ell'_s = b \text{ or } \ell'_{s-2} \\ (\bar{\alpha}_t(s) + \alpha_{t-1}(s-2))y_{\ell'_s}^t & \text{otherwise} \end{cases} \quad (600)$$

$$\bar{\alpha}_t(s) \triangleq \alpha_{t-1}(s) + \alpha_{t-1}(s-1) \quad (601)$$

It's worth emphasizing how to interpret these, given we've imposed this weird augmented label sequence. In as-verbose-as-possible terms,

$\alpha_t(s)$  is the probability, after running our RNN for  $t$  time steps to produce the path  $\pi_{\langle 1 \dots t \rangle}$ , that  $\mathcal{B}(\pi_{\langle 1 \dots t \rangle}) == \ell_{\langle 1 \dots \frac{s-1}{2} \rangle}$  for which, after inserting  $\epsilon$  between all elements of  $\ell_{\langle 1 \dots \frac{s-1}{2} \rangle}$ , we obtain the augmented labeling  $\ell'_{\langle 1 \dots s \rangle}$ .

The way you should think about the different possible cases here is that, at time step  $t$ , in order for there to be nonzero probability that we can merge the sequence of  $t$  RNN outputs into the augmented label sequence  $\ell'_{\langle 1 \dots s \rangle}$ , it must be true that:

- We emit the token  $\ell'_s$  at time  $t$  from the RNN.
- At the previous timestep,  $t-1$ , we emitted a token consistent with our rules for merging combined with the fact that we've inserted  $\epsilon$  between every pair of tokens in the final output labeling  $\ell$ , in order to produce  $\ell'$ .

The weird case (in my opinion) to consider is realizing that we can emit, for example, the label  $a$  at time  $t-1$ , then the label  $b$  at time  $t$ , and this would eventually get mapped to a portion of the augmented label sequence,  $[a, \epsilon, b]$ .

## BERT

Table of Contents Local

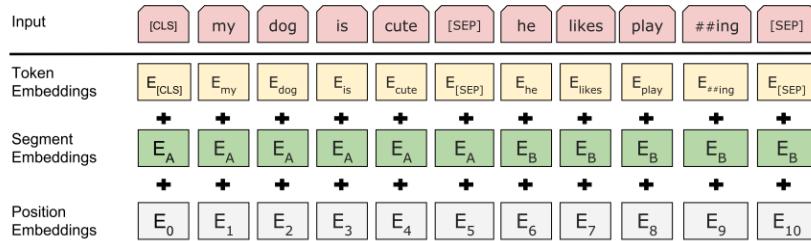
Written by Brandon McKinzie

Devlin et al., “BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding,” Google AI Language (Oct 2018).

**TL;DR.** Bidirectional Encoder Representations from Transformers. Pretrained by jointly conditioning on left and right context, and can be fine-tuned with one additional non-task-specific output layer. Authors claim the following contributions:

- Demonstrate importance of bidirectional pre-training for language representations. Ok, congrats.
- Show that pre-trained representations eliminate needs of task-specific architectures. We already knew this. Seriously, how is this news?
- Advances SOTA for eleven NLP tasks.

**BERT.** Instead of using the unidirectional transformer *decoder*, they use the bidirectional *encoder* architecture for their pre-trained model<sup>138</sup>.



The input representation is shown in the figure above. The input is a sentence pair, as commonly seen in tasks like QA/NLI.

<sup>138</sup>Is this seriously paper-worthy?? I'm taking notes so I can easily refer back on popular approaches, but I don't see what's so special here.

### Pre-training Tasks (3.3).

1. **Masked LM:** Mask 15% of all tokens, and try to predict only those masked tokens<sup>139</sup> Furthermore, at training time, the mask tokens are either fed through as (a) the special [MASK] token 80% of the time, (b) a random word 10% of the time, and (c) the original word unchanged 10% of the time. Now *this* is just hackery.
2. **Next Sentence Prediction:** Given two input sentences A and B, train a binary classifier to predict whether sentence B came directly after sentence A.

They do the pretraining jointly, using a loss function that's the sum of the mean masked LM likelihood and mean next sentence prediction likelihood.

---

<sup>139</sup>This is the only “novel” idea I’ve seen in the paper. Seems hacky-ish but also reasonable.

## Wasserstein is all you need

Table of Contents Local

Written by Brandon McKinzie

Singh et al., “Wasserstein is all you need,” EPFL Switzerland (August 2018).

**TL;DR.** Unsupervised representations of objects/entities via distributional + point estimate. Made possible by **optimal transport**.

**Optimal Transport** (3). First, notation. Let…

- $\Omega$  denote a space of possible outcomes.
- $\mu$  denote an **empirical** probability measure, defined as some convex combination  $\mu(\mathbf{x}) = \sum_{i=1}^n a_i \delta(x_i)$ , where  $x_i \in \Omega$ .
- $\nu$  denote a similar measure, also a convex combination,  $\nu(\mathbf{y}) = \sum_j b_j \delta(y_j)$ .
- $M_{ij}$  denote the ground cost of moving from point  $x_i$  to  $y_j$ .

Intuition break: recognize that  $\mu$  and  $\nu$  are just a formal description of probability densities via normalized “counts”  $a_i$  and/or  $b_j$ . Those weights are basically probability mass. The **Optimal Transport** distance between  $\mu$  and  $\nu$  is the following **linear program**:

$$\text{OT}(\mu, \nu; \mathbf{M}) = \min_{\mathbf{T} \in \mathbb{R}_+^{n \times m}} \sum_{ij} T_{ij} M_{ij} \quad \text{s.t.} \quad (\forall i) \sum_j T_{ij} = a_i, \quad (\forall j) \sum_i T_{ij} = b_j \quad (602)$$

where  $\mathbf{T}$  is called the **transportation matrix**. Informally, the constraints are simply enforcing bijection to/from  $\mu$  and  $\nu$ , in that “all the mass sent from element  $i$  must be exactly  $a_i$ , and all mass sent to element  $j$  must be exactly  $b_j$ ”. A particular case called the **p-Wasserstein distance**, where  $\Omega = \mathbb{R}^d$  and  $M_{ij}$  is a distance metric over  $\mathbb{R}^d$ , is defined as

$$W_p(\mu, \nu) \triangleq \text{OT}(\mu, \nu; D_\Omega^p)^{1/p} \quad (603)$$

where  $D$  is just a distance metric, e.g. for  $p = 2$  it could be euclidean distance.

**Distributional Estimate** (4.1). Let  $\mathcal{C} \triangleq \{c\}_i$  be the set of possible contexts, where each context  $c_i$  can be a word/phrase/sentence/etc. For a given word  $w$  and our set of observed contexts for that word, we essentially want to embed its context histogram into a space  $\Omega$  (where typically  $\Omega = \mathbb{R}^d$ ). Let  $\mathbf{V}$  denote a matrix of context embeddings, such that  $V_{i,:} = \mathbf{c}_i \in \mathbb{R}^d$ , the embedding for context  $c_i$  in what the authors call the *ground space*. Combining the histogram  $H^w$  containing observed context counts for word  $w$  with  $\mathbf{V}$ , the **distributional estimate** of the word  $w$  is defined as

$$P_{\mathbf{V}}^w \triangleq \sum_{c \in \mathcal{C}} (H^w)_c \delta(\mathbf{v}_c) \quad (604)$$

Also, the *point estimate* is just  $\mathbf{v}_w$ , i.e. the embedding of the word  $w$  when viewed as a context.

**Distance** (4.2). Given some distance metric  $D_\Omega$  in ground space  $\Omega = \mathbb{R}^d$ , the distance between words  $w_i$  and  $w_j$  is the solution to the following OT problem<sup>140</sup>:

$$\text{OT}(P_{\mathbf{V}}^{w_i}, P_{\mathbf{V}}^{w_j}; D_\Omega^p) := W_p^\lambda(P_{\mathbf{V}}^{w_i}, P_{\mathbf{V}}^{w_j})^p \quad (605)$$

**Concrete Framework** (5). The authors make use of the **shifted positive pointwise mutual information** (SPPMI),  $\mathbf{S}_s^\alpha$ , for computing the word histograms:

$$(H^w)_c := \frac{\mathbf{S}_s^\alpha(w, c)}{\sum_{c' \in \mathcal{C}} \mathbf{S}_s^\alpha(w, c')} \quad (606)$$

$$\mathbf{S}_s^\alpha(w, c) \triangleq \max \left[ \log \left( \frac{\text{Count}(w, c) \sum_{c'} \text{Count}(c')^\alpha}{\text{Count}(w) \text{Count}(c)^\alpha} \right) - \log(s), 0 \right] \quad (607)$$

---

<sup>140</sup>I'm not sure whether the rightmost exponent of  $p$  is a typo in the paper, but that is how it is written.

## Noise Contrastive Estimation

Table of Contents Local

Written by Brandon McKinzie

M. Gutman and A. Hyvärinen, “Noise contrastive estimation: A new estimation principle for unnormalized statistical models,” University of Helsinki (2010).

**TL;DR:** A few ways of thinking about NCE:

- Instead of directly modeling a normalized word distribution  $p_d(w)$ , we can just model the unnormalized  $\tilde{p}_d$  distribution (and an additional parameter  $c = -\ln Z$ ) by training our model to distinguish between true samples from  $p_d$  and noise samples from some distribution  $p_n$  that we choose.
- Instead of modeling  $p(w | c)$ , we model  $p(D | w, c)$ , where  $D$  is binary RV indicating whether  $w, c$  are from the true data distribution  $p_d$  or the noise distribution  $p_n$ .

**Introduction.** Setup & notation:

- We observe  $\mathbf{x} \sim p_d(\cdot)$  but  $p_d(\cdot)$  itself is unknown.
- We model  $p_d$  by some model  $p_m(\cdot; \alpha)$  parameterized by  $\alpha^{141}$ .

So, can we get away with modeling the *unnormalized* density  $\tilde{p}_m$  instead of requiring the normalization constraint to be baked in to our optimization problem? Similar to approaches like **contrastive divergence** and **score matching**, **noise-contrastive estimation** (NCE) aims to address this question.

**Noise-contrastive estimation** (2.1). Let  $c$  be an estimator for  $-\ln Z(\alpha)$ , and let  $\theta = \{\alpha, c\}$  denote all of our parameters. Given observed data  $X = \{\mathbf{x}_1, \dots, \mathbf{x}_T\}$ , and noise  $Y = \{\mathbf{y}_1, \dots, \mathbf{y}_T\}$ , we seek parameters  $\hat{\theta}_T$  that maximize  $J_T(\theta)^{142}$ :

$$J_T(\theta) = \frac{1}{2T} \sum_t \ln [h(\mathbf{x}_t)] + \ln [1 - h(\mathbf{y}_t)] \quad (608)$$

$$h(\mathbf{u}) = \sigma(G(\mathbf{u})) \quad (609)$$

$$G(\mathbf{u}) = \ln p_m(\mathbf{u}) - \ln p_n(\mathbf{u}) \quad (610)$$

$$\ln p_m(\cdot; \theta) := \ln \tilde{p}_m(\cdot; \alpha) + c \quad (611)$$

where for compactness reasons I’ve removed the explicit dependence of all functions above (except  $p_n$ ) on  $\theta$ . Notice how this fixes the issue of the model just setting  $c$  arbitrarily high to obtain a high likelihood<sup>143</sup>.

<sup>141</sup>The implicit assumption here is that  $\exists \alpha^*$  such that  $p_d(\cdot) = p_m(\cdot; \alpha^*)$ .

<sup>142</sup>This all assumes of course that  $p_n$  is fully defined.

<sup>143</sup>The primary reason why MLE is traditionally unable to parameterize the partition function.

**Connection to supervised learning** (2.2). The NCE objective can be interpreted as binary logistic regression that discriminates whether a point belongs to the data ( $p_d$ ) or to the noise ( $p_n$ ).

$$P(C=1 \mid \mathbf{u} \in X \cup Y) = \frac{P(C=1)p(\mathbf{u} \mid C=1)}{p(\mathbf{u})} \quad (612)$$

$$= \frac{p_m(\mathbf{u})}{p_m(\mathbf{u}) + p_n(\mathbf{u})} \quad (613)$$

$$\equiv h(\mathbf{u}; \theta) \quad (614)$$

We model with a uniform prior:  
 $P(C=1) = P(C=0) = 1/2$

where we're now using the union of  $X$  and  $Y$ ,  $U := \{\mathbf{u}_1, \dots, \mathbf{u}_{2T}\}$ . The log-likelihood of the data under the parameters  $\theta$  is

$$\ell(\theta) = \sum_t^{2T} [C_t \ln P(C_t=1 \mid \mathbf{u}_t; \theta) + (1 - C_t) \ln P(C_t=0 \mid \mathbf{u}_t; \theta)] \quad (615)$$

$$= \sum_t^{2T} [C_t \ln h(\mathbf{u}_t; \theta) + (1 - C_t) \ln [1 - h(\mathbf{u}_t; \theta)]] \quad (616)$$

$$= \sum_t^T [\ln h(\mathbf{x}_t; \theta) + \ln [1 - h(\mathbf{y}_t; \theta)]] \quad (617)$$

which is (up to a constant factor) the same as our NCE objective.

**Properties of the estimator** (2.3). As  $T \rightarrow \infty$ ,  $J_T(\theta) \rightarrow J(\theta)$ , where

$$J(\theta) \triangleq \lim_{T \rightarrow \infty} J_T(\theta) = \frac{1}{2} \mathbb{E} [\ln h(\mathbf{x}; \theta) + \ln [1 - h(\mathbf{y}; \theta)]] \quad (618)$$

$$\tilde{J}(f) \triangleq \frac{1}{2} \mathbb{E} \left[ \ln \left[ \sigma \left( f(\mathbf{x}) - \ln p_n(\mathbf{x}) \right) \right] + \ln \left[ 1 - \sigma \left( f(\mathbf{y}) - \ln p_n(\mathbf{y}) \right) \right] \right] \quad (619)$$

### Theorem 1

$\tilde{J}$  attains exactly one maximum, located at  $f(\cdot) = \ln p_d(\cdot)$ , provided  $p_d(\cdot) \neq 0 \implies p_n(\cdot) \neq 0$ .

---

#### 4.52.1 SELF-NORMALIZED NCE

---

Notes from A. Mnih and Y. Teh, “A fast and simple algorithm for training neural probabilistic language models” (2012).

**Maximum likelihood learning.** Let  $P_\theta^h(w)$  denote the probability of observing word  $w$  given context  $h$ . For neural LMs, we assume this is the softmax of a scoring function  $s_\theta(w, h)$  (logits). In what follows, I’ll drop the explicit  $\theta$  and  $h$  subscript/superscript notation for brevity.

$$\frac{\partial \log P(w)}{\partial \theta} = \frac{\partial}{\partial \theta} s(w, h) - \frac{\partial}{\partial \theta} \log \left[ \sum_{w'} e^{s(w', h)} \right] \quad (620)$$

$$= \frac{\partial}{\partial \theta} s(w, h) - \sum_{w'} P(w') \frac{\partial}{\partial \theta} s(w', h) \quad (621)$$

$$= \frac{\partial}{\partial \theta} s(w, h) - \mathbb{E}_{w \sim P_\theta^h} \left[ \frac{\partial}{\partial \theta} s(w, h) \right] \quad (622)$$

where the expectation (in red) is expensive due to requiring  $s(w, h)$  evaluated for all words in the vocabulary. One approach is **importance sampling** where we sample a subset of  $k$  words from the vocab and compute the probabilities from that approximation:

$$\frac{\partial \log P(w)}{\partial \theta} = \frac{\partial}{\partial \theta} s(w, h) - \sum_{w'} P(w') \frac{\partial}{\partial \theta} s(w', h) \quad (623)$$

$$\approx \frac{\partial}{\partial \theta} s(w, h) - \frac{1}{V} \sum_{j=1}^k v(x_j) \frac{\partial}{\partial \theta} s(w', h) \quad (624)$$

$$\text{where } v(x) = \frac{e^{s_\theta(x, h)}}{Q^h(w=x)} \quad (625)$$

and we refer to  $v$  as the **importance weights**. In NLP, we typically set  $Q$  to the **Zipfian distribution**<sup>144</sup>

---

<sup>144</sup>TensorFlow seems to define this as

$$P_{Zipf}(w) = \frac{\log(w+2) - \log(w+1)}{\log(V+1)} \quad (626)$$

where  $V$  is the vocabulary size. I can’t seem to find this definition anywhere else though. A more common form seems to be

$$P(w) = \frac{\frac{1}{w}}{\sum_{w'} \frac{1}{w'^s}} \quad (627)$$

I plotted both on WolframAlpha (link here) and they do indeed look basically the same, especially for any reasonably large  $V$ .

**NCE.** In NCE, we introduce [unigram] noise distribution  $P_n(w)$  and impose a prior that noise samples are  $k$  times more frequent than data samples from  $P_d^h(w)$ , resulting in the joint distribution,

$$P^h(D, w) = P(D=1)P_d^h(w) + P(D=0)P_n(w) \quad (628)$$

$$= \frac{1}{k+1}P_d^h(w) + \frac{k}{k+1}P_n(w) \quad (629)$$

Our goal is to learn the posterior distribution  $P^h(D=1 | w)$  (so we replace  $P_d$  with  $P_\theta$ ):

$$P^h(D=1 | w, \theta) = \frac{P_\theta^h(w)}{P_\theta^h(w) + kP_n(w)} \quad (630) \quad \text{NCE posterior}$$

In NCE, we re-parameterize  $P_\theta$  by treating  $-\log Z$  as a parameter itself,  $c^{h145}$ .

$$P_\theta^h(w) := P_{\theta^0}^h \exp(c^h) \quad (631)$$

where  $P_{\theta^0}^h$  denotes the unnormalized distribution. It turns out that, in practice, we can impose that  $\exp(c^h)=1$  and use the unnormalized  $P_{\theta^0}^h$  in place of the true probabilities in all that follows. *Critically, note that this means we rewrite equation 630 using the unnormalized probabilities in place of  $P_\theta^h$ .* The NCE objective is to find<sup>146</sup> as follows, where I've shown each step of the derivation:

$$\theta^* = \arg \max_{\theta} J^h(\theta) \quad (632)$$

$$J^h(\theta) = \mathbb{E}_{(D, w) \sim P^h} [\log P(D | w, \theta)] \quad (633)$$

$$= \sum_{D=0}^1 \sum_w P^h(D, w) \log P^h(D | w, \theta) \quad (634)$$

$$= \sum_w P^h(0, w) \log P^h(0 | w) + \sum_w P^h(1, w) \log P^h(1 | w) \quad (635)$$

$$= \frac{1}{k+1} \sum_w [kP_n(w) \log P^h(0 | w) + P_d^h(w) \log P^h(1 | w)] \quad (636)$$

$$= \frac{1}{k+1} [k\mathbb{E}_{P_n} [\log P^h(0 | w)] + \mathbb{E}_{P_d^h} [\log P^h(1 | w)]] \quad (637)$$

$$\propto k\mathbb{E}_{P_n} [\log P^h(0 | w)] + \mathbb{E}_{P_d^h} [\log P^h(1 | w)] \quad (638)$$

The gradient of the NCE objective is thus

$$\frac{\partial}{\partial \theta} J^h(\theta) = \sum_w \frac{kP_n(w)}{P_\theta^h(w) + kP_n(w)} (P_d^h(w) - P_\theta^h(w)) \frac{\partial}{\partial \theta} \log P_\theta^h(w) \quad (639)$$

**TODO:** incorporate more info from Chris Dyer's excellent notes.

---

<sup>145</sup>Reminder that the  $h$  is a reminder that  $Z$  is a function of the context  $h$ .

<sup>146</sup>I'll drop off  $\theta$  dependence wherever obvious for the sake of compactness.

## Neural Ordinary Differential Equations

Table of Contents Local

Written by Brandon McKinzie

R. Chen, Y. Rubanova, J. Bettencourt, and D. Duvenaud, “Neural Ordinary Differential Equations,” University of Toronto (Oct 2018).

**Introduction** (1). Let  $T$  denote the number of layers of our network. In the limit of  $T \rightarrow \infty$  and small  $\delta\mathbf{h}(t)$  between each “layer”, we can parameterize the dynamics via an ODE:

$$\frac{d\mathbf{h}(t)}{dt} = f(\mathbf{h}(t), t, \theta) \quad (640)$$

Benefits of defining models using ODE solvers:

- **Memory.** Constant as a function of depth, since we don’t need to store intermediate values from the forward pass.
- **Adaptive computation.** Modern ODE solvers adapt evaluation strategy on the fly.
- **Parameter efficiency.** Nearby “layers” have shared parameters.

## Review: ODE

Remember the basic idea with ODEs like the one shown above. Our goal is to solve for  $\mathbf{h}(t)$ .

$$d\mathbf{h}(t) = f(\mathbf{h}(t), t, \theta) dt \quad (641)$$

$$\int d\mathbf{h}(t) = \int f(\mathbf{h}(t), t, \theta) dt \quad (642)$$

$$\mathbf{h}(t) + c_1 = \int f(\mathbf{h}(t), t, \theta) dt \quad (643)$$

$$(644)$$

and so the solution of an ODE is often represented as an integral.

**Reverse-mode automatic differentiation of ODE solutions** (2). Our goal is to optimize

$$L(\mathbf{z}(t_1)) = L\left(\int_{t_0}^{t_1} f(\mathbf{z}(t), t, \theta) dt\right) \quad (645)$$

Given our starting definition (eq 640), we can say

$$\mathbf{z}(t + \epsilon) = \mathbf{z}(t) + \int_t^{t+\epsilon} f(\mathbf{z}(t), t, \theta) dt := T_\epsilon(\mathbf{z}(t), t) \quad (646)$$

which we can use to define the **adjoint**  $a(t)$ :

$$a(t) \triangleq -\frac{\partial L}{\partial \mathbf{z}(t)} = -\frac{\partial L}{\partial \mathbf{z}(t+\epsilon)} \frac{d\mathbf{z}(t+\epsilon)}{d\mathbf{z}(t)} \quad (647)$$

$$= a(t+\epsilon) \frac{\partial T_\epsilon(\mathbf{z}(t), t)}{\partial \mathbf{z}(t)} \quad (648)$$

$$\frac{da(t)}{dt} = -a(t)^T \frac{\partial f(\mathbf{z}(t), t, \theta)}{\partial \mathbf{z}} \quad (649)$$

where  $\frac{da(t)}{dt}$  can be derived using the limit definition of a derivative. We'll now outline the algorithm for computing gradients. We use a black box ODE solver as a subroutine that solves a first-order ODE initial value problem. As such, it accepts an initial state, its first derivative, the start time, the stop time, and parameters  $\theta$  as arguments.

### Reverse-mode derivative

*Given start time  $t_0$ , stop time  $t_1$ , final state  $\mathbf{z}(t_1)$ , parameters  $\theta$ , and gradient  $\frac{\partial L}{\partial \mathbf{z}(t_1)}$  compute all gradients of  $L$ .*

1. Compute  $t_1$  gradient:

$$\frac{\partial L}{\partial t_1} = \frac{\partial L}{\partial \mathbf{z}(t_1)}^T \frac{\partial \mathbf{z}(t_1)}{\partial t_1} = \frac{\partial L}{\partial \mathbf{z}(t_1)}^T f(\mathbf{z}(t_1), t_1, \theta) \quad (650)$$

2. Initialize the augmented state:

$$s_0 := \left[ \mathbf{z}(t_1), \mathbf{a}(t_1), \mathbf{0}, -\frac{\partial L}{\partial t_1} \right] \quad (651)$$

3. Define augmented state dynamics:

$$\frac{ds}{dt} \triangleq \left[ f(\mathbf{z}(t), t, \theta), -\mathbf{a}(t)^T \frac{\partial f}{\partial \mathbf{z}}, -\mathbf{a}(t)^T \frac{\partial f}{\partial \theta}, -\mathbf{a}(t)^T \frac{\partial f}{\partial t} \right] \quad (652)$$

4. Solve **reverse-time**<sup>a</sup> ODE:

$$\left[ \mathbf{z}(t_0), \frac{\partial L}{\partial \mathbf{z}(t_0)}, \frac{\partial L}{\partial \theta}, \frac{\partial L}{\partial t_0} \right] = \text{ODESolve}\left(s_0, \frac{ds}{dt}, t_1, t_0, \theta\right) \quad (653)$$

5. Return  $\frac{\partial L}{\partial \mathbf{z}(t_0)}$ ,  $\frac{\partial L}{\partial \theta}$ ,  $\frac{\partial L}{\partial t_0}$ ,  $\frac{\partial L}{\partial t_1}$

---

<sup>a</sup>Notice how our “initial state” actually corresponds to  $t_1$ , and we pass in  $t_1$  and  $t_0$  in the opposite order we typically do.

## On the Dimensionality of Word Embedding

Table of Contents Local

Written by Brandon McKinzie

Z. Yin and Y. Shen, “On the Dimensionality of Word Embedding,” Stanford University (Dec 2018).

**Unitary Invariance of Word Embeddings** (2.1). Authors interpret result any unitary transformation (e.g. a rotation) on embedding matrix as equivalent to the original.

**Word Embeddings from Explicit Matrix Factorization** (2.2). Let  $M$  be the  $V \times V$  co-occurrence counts matrix. One way of getting embeddings is doing a truncated SVD on  $M = UDV^T$ . If we want  $k$ -dimensional embedding vectors, we can do

$$\mathbf{E} = \mathbf{U}_{1:k} \mathbf{D}_{1:k, 1:k}^\alpha \quad (654)$$

for some  $\alpha \in [0, 1]$ .

**PIP Loss** (3). Given embedding matrix  $E \in \mathbb{R}^{V \times d}$ , define its **Pairwise Inner Product** (PIP) matrix to be

$$\text{PIP}(\mathbf{E}) \triangleq \mathbf{E}\mathbf{E}^T \quad (655)$$

Notice that  $\text{PIP}(\mathbf{E})_{i,j} = \langle \mathbf{w}_i, \mathbf{w}_j \rangle$ . Define the **PIP loss** for comparing two embeddings  $\mathbf{E}_1$  and  $\mathbf{E}_2$  for a common vocab of  $V$  words:

$$\|\text{PIP}(\mathbf{E}_1) - \text{PIP}(\mathbf{E}_2)\|_F = \sqrt{\sum_{i,j} \left( \langle \mathbf{w}_i^{(1)}, \mathbf{w}_j^{(1)} \rangle - \langle \mathbf{w}_i^{(2)}, \mathbf{w}_j^{(2)} \rangle \right)^2} \quad (656)$$

## Generative Adversarial Nets

Table of Contents Local

Written by Brandon McKinzie

Goodfellow et al., “Generative Adversarial Nets,” (June 2014)

**TL;DR.** The abstract is actually quite good:

... we simultaneously train two models: a generative model  $\mathbf{G}$  that captures the data distribution, and a discriminative model  $\mathbf{D}$  that estimates the probability that a sample came from the training data rather than  $\mathbf{G}$ . The training procedure for  $\mathbf{G}$  is to maximize the probability of  $\mathbf{D}$  making a mistake. This framework corresponds to a minimax two-player game.

**Adversarial Nets** (3). As usual, first go over notation:

- Generator produces data samples<sup>147</sup>,  $\mathbf{x} := \mathbf{G}(\mathbf{z}; \theta_g)$ , where  $\mathbf{z} \sim p_n$  (noise distribution prior).
- Discriminator,  $\mathbf{D}(\mathbf{x}; \theta_d)$ , outputs probability that  $\mathbf{x}$  came from (true)  $p_{data}$  instead of  $G$ .

Our two-player minimax optimization problem can be written as:

$$\min_G \max_D V(\mathbf{D}, \mathbf{G}) = \mathbb{E}_{\mathbf{x} \sim p_{data}} [\log \mathbf{D}(\mathbf{x})] + \mathbb{E}_{\mathbf{z} \sim p_n} [\log (1 - \mathbf{D}(\mathbf{G}(\mathbf{z})))] \quad (657)$$

**Theoretical Results** (4). Below is the training algorithm.

### SGD with GANs

Repeat the following for each training iteration.

1. Train  $\mathbf{D}$ . For  $k$  steps, repeat:
  - (a) Sample  $m$  noise samples  $\{\mathbf{z}_1, \dots, \mathbf{z}_m\}$  from noise prior  $p_n$ .
  - (b) Sample  $m$  data samples  $\{\mathbf{x}_1, \dots, \mathbf{x}_m\}$  from data distribution  $p_{data}$ .
  - (c) Update discriminator by ascending  $\nabla_{\theta_d} V(\mathbf{D}, \mathbf{G})$ .
2. Train  $\mathbf{G}$ : Sample another  $m$  noise samples  $\{\mathbf{z}_1, \dots, \mathbf{z}_m\}$  and descend on  $\nabla_{\theta_g} V(\mathbf{D}, \mathbf{G})$ .

---

<sup>147</sup>Note that  $\mathbf{G}$  outputs samples  $\mathbf{x}$ , not probabilities. By doing this, it *implicitly* defines a probability distribution  $p_g(\mathbf{x})$ . This is what the authors say.

**Global Optimality of  $p_g \equiv p_{data}$**  (4.1).

### Proposition 1

For fixed  $G$ , the optimal  $D$  is

$$D_G^*(\mathbf{x}) = \frac{p_{data}(\mathbf{x})}{p_{data}(\mathbf{x}) + p_g(\mathbf{x})} \quad (658)$$

**Derivation of  $D_G^*(\mathbf{x})$ .**

**Aside: Law of the unconscious statistician** (LotUS). The distribution  $p_g(\mathbf{x})$  should be read as “the probability that the output of  $G$  yields the value  $\mathbf{x}$ .” Take a step back and recognize that  $G$  is simply a function of a random variable  $\mathbf{z}$ . As such, we can apply familiar rules like

$$\mathbb{E}[G(\mathbf{z})] = \mathbb{E}_{\mathbf{z} \sim p_n}[G(\mathbf{z})] \quad (659)$$

$$= \int_{\mathbf{z}} p_n(\mathbf{z}) G(\mathbf{z}) d\mathbf{z} \quad (660)$$

However, recall that functions of random variables can themselves be interpreted as random variables. In other words, we can also use the interpretation that  $G$  evaluates to some output  $\mathbf{x}$  with probability  $p_g(\mathbf{x})$ .

$$\mathbb{E}[G] = \mathbb{E}_{\mathbf{x} \sim p_g}[\mathbf{x}] \quad (661)$$

$$= \int_{\mathbf{x}} p_g(\mathbf{x}) \mathbf{x} d\mathbf{x} \quad (662)$$

As this blog post details, this equivalence is NOT due to a change of variables, but rather by the **Law of the unconscious statistician**.

**The Proof:** We can directly use LotUS to rewrite  $V(G, D)$ :

$$V(G, D) = \mathbb{E}_{\mathbf{x} \sim p_{data}} [\log D(\mathbf{x})] + \mathbb{E}_{\mathbf{z} \sim p_n} [\log (1 - D(G(\mathbf{z})))] \quad (663)$$

$$= \mathbb{E}_{\mathbf{x} \sim p_{data}} [\log D(\mathbf{x})] + \mathbb{E}_{\mathbf{x} \sim p_g} [\log (1 - D(\mathbf{x}))] \quad (664)$$

$$= \int_{\mathbf{x}} [p_{data}(\mathbf{x}) \log(D(\mathbf{x})) + p_g(\mathbf{x}) \log(1 - D(\mathbf{x}))] d\mathbf{x} \quad (665)$$

LotUS allowed us to express  $V(G, D)$  as a continuous function over  $\mathbf{x}$ . More importantly, it means we can evaluate  $\frac{\partial V}{\partial D}$  and take the derivative inside the integral<sup>a</sup>. Setting the derivative to zero and solving for  $D$  yields  $D_G^*$ , the form that maximizes  $V$ .

---

<sup>a</sup>Also remember that  $D(\cdot) \in [0, 1]$  since it is a probability distribution.

The authors use this proposition to define the virtual training criterion  $C(G) \triangleq V(G, D_G^*)$ :

$$C(G) = \mathbb{E}_{\mathbf{x} \sim p_{data}} \left[ \log \frac{p_{data}(\mathbf{x})}{p_{data}(\mathbf{x}) + p_g(\mathbf{x})} \right] + \mathbb{E}_{\mathbf{x} \sim p_g} \left[ \log \frac{p_g(\mathbf{x})}{p_{data}(\mathbf{x}) + p_g(\mathbf{x})} \right] \quad (666)$$

### Theorem 1.

*The global minimum of  $C(G)$  is achieved IFF  $p_g = p_{data}$ . At that point  $C(G) = -\log 4$ .*

#### Proof: Theorem 1

The authors subtract  $V(D_G^*, G; p_g=p_{data})$  from both sides of 666, do some substitutions, and find that

$$C(G) = 2 \cdot JSD(p_{data} || p_g) - \log 4 \quad (667)$$

where  $JSD$  is the **Jensen-Shannon divergence**<sup>a</sup>. Since  $0 \leq JSD(p||q)$  always, with equivalence only if  $p \equiv q$ , this proves Theorem 1 above.

---

<sup>a</sup>Recall that the JSD represents the divergence of each distribution from the mean of the two

## A Framework for Intelligence and Cortical Function

Table of Contents   Local

*Written by Brandon McKinzie*

J. Hawkins et al., “A Framework for Intelligence and Cortical Function Based on Grid Cells in the NeoCortex,” Numata Inc. (Oct 2018).

**Introduction.** Authors propose new framework based on location processing that provides supporting evidence to the theory that **all regions of the neocortex are fundamentally the same**. We’ve known that **grid cells** exist in the hippocampal complex of mammals, but only recently have seen evidence that they may be present in the neocortex.

**How Grid Cells Represent Location.** Grid cells in the **entorhinal cortex**<sup>148</sup> represent space and location. The main concepts, in order such that they build on one another, are as follows:

- A **single grid cell** is a neuron that fires [when the agent is] at [one of many] multiple locations in a physical environment<sup>149</sup>.
- A **grid cell module** is a set of grid cells that activate with the *same lattice spacing and orientation* but at shifted locations within an environment.
- **Multiple grid cell modules** that differ in tile spacing and/or orientation can provide *unique location information*<sup>150</sup>.

Crucially, the number of unique locations that can be represented by a set of grid cell modules **scales exponentially** with the number of modules. Every learned environment is associated with a set of unique locations (firing patterns of the grid cells).

**Grid Cells in the Neocortex.** The authors propose that we learn the structures of objects (like pencils, cups, etc) via grid cells in the *neocortex*. Specifically, they propose:

1. Every cortical column has neurons that perform a function similar to grid cells.
2. Cortical columns learn models of *objects* similarly to how grid/place cells learn models of *environments*.

<sup>148</sup>The entorhinal cortex is located in the medial temporal lobe and functions as a hub in a widespread network for memory, navigation and the perception of time.

<sup>149</sup>For example, there may be a grid cell in my brain that fires when I’m at certain locations inside my room. Those locations tend to form a lattice of sorts.

<sup>150</sup>A single module alone cannot, because it repeats periodically. In other words, it can only provide relative location information.

## Large-Scale Study of Curiosity Driven Learning

Table of Contents Local

Written by Brandon McKinzie

Burda et al., "Large-Scale Study of Curiosity Driven Learning," OpenAI and UC Berkeley (Aug 2018).

An agent sees observation  $x_t$ , takes action  $a_t$ , and transitions to the next state with observation  $x_{t+1}$ . Goal: incentivize agent with reward  $r_t$  relating to how informative the transition was. The main components in what follows are:

- Observation embedding  $\phi(x)$ .
- Forward dynamics network for prediction  $P(\phi(x_{t+1} | x_t, a_t))$ .
- Exploration reward (*surprisal*):

$$r_t = -\log p(\phi(x_{t+1}) | x_t, a_t) \quad (668)$$

The authors choose to model the next state embedding with a Gaussian,

$$\phi(x_{t+1}) | x_t, a_t \sim \mathcal{N}(f(x_t, a_t), \epsilon) \quad (669)$$

$$r_t = \|f(x_t, a_t) - \phi(x_{t+1})\|_2^2 \quad (670)$$

where  $f$  is the learned dynamics model.

**Feature spaces** (2.1). Some possible ways to define  $\phi$ :

- **Pixels:**  $\phi(x) = x$ .
- **Random Features:** Literally feeding  $\phi(x) = ConvNet(x)$  where ConvNet is *randomly initialized* and fixed.
- **VAE:** Use the mapping to the mean [of the approximated distribution] as the embedding network  $\phi$ .

**Interpretation.** It seems that this works because after awhile, it is boring and predictable to take actions that result in losing a game. The most surprising actions seem to be those that advance us forward, to new and uncharted territory. However, these experiments are all on games that have a very "linear" uni-directional-like sequence of success. I wonder how successful this would be in a game like rocket league, where there is no tight coupling of direction with success and novelties (e.g. moving forward in mario bros).

## Universal Language Model Fine-Tuning for Text Classification

Table of Contents Local

Written by Brandon McKinzie

J. Howard and S. Ruder, “Universal Language Model Fine-Tuning for Text Classification,” (May 2018).

**TL;DR.** ULMFiT is a transfer learning method. They introduce techniques for fine-tuning language models.

**Universal Language Model Fine-tuning (3).** Define the general **inductive transfer learning** setting for NLP:

*Given a source task  $\mathcal{T}_s$  and target task  $\mathcal{T}_T \neq \mathcal{T}_s$ , improve performance on  $\mathcal{T}_T$ .*

ULMFiT is defined by the following three steps:

1. General-domain LM pretraining.
2. Target task LM fine-tuning.
3. Target task classifier fine-tuning.

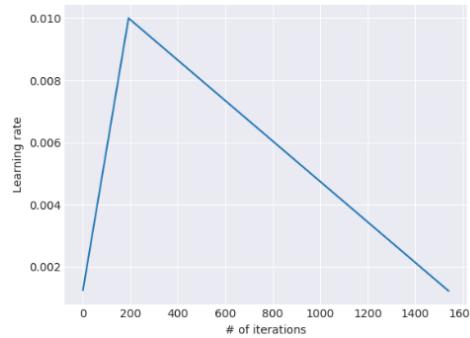
**Target Task LM Fine-tuning (3.2).** For step 2, the authors propose what they call **discriminative fine-tuning** and **slanted triangular learning rates**.

- **Discriminative fine-tuning.** Tune each layer with different learning rates:

$$\theta_t^\ell = \theta_{t-1}^\ell - \eta^\ell \cdot \nabla_{\theta^\ell} J(\theta) \quad (671)$$

The authors suggest setting  $\eta^{\ell-1} = \eta^\ell / 2.6$ .

- **Slanted triangular learning rates.** A type of learning rate schedule that looks like the picture below.



First, we define the following hyperparameters:

- $T$ : total number of training iterations<sup>151</sup>
- $cfrac$ : fraction of  $T$  (in num iterations) where we’re *increasing* the learning rate.
- $cut$ :  $\lfloor T \cdot cfrac \rfloor$ . Iteration where we switch from increasing the LR to decreasing it.
- $ratio$ :  $\eta_{max}/\eta_{min}$ . We of course must also define  $\eta_{max}$ .

We can now compute the learning rate for a given iteration  $t$ :

$$\eta_t = \eta_{max} \cdot \frac{1 + p \cdot (ratio - 1)}{ratio} \quad (672)$$

$$p = \begin{cases} \frac{t}{cut} & \text{if } t < cut \\ 1 - \frac{t-cut}{cut \cdot (1/cfrac-1)} & \text{otherwise} \end{cases} \quad (673)$$

Suggested:  
 $cfrac=0.1$   
 $ratio=32$   
 $\eta_{max}=0.01$

**Target Task Classifier Fine-tuning** (3.3). Augment the LM with two fully-connected layers. The first with ReLU activation and the second with softmax. Each uses batch normalization and dropout. The first is fed the output hidden state of the LM concatenated with the max- and mean-pooled hidden states over all timesteps<sup>152</sup>:

$$\mathbf{h}_c = [\mathbf{h}_t, \text{maxpool}(\mathbf{H}), \text{meanpool}(\mathbf{H})] \quad (674)$$

In addition to DF-T and STLR from above, they also employ the following techniques:

- **Gradual Unfreezing**: first unfreeze the last layer and fine-tune it alone with all other layers frozen *for one epoch*. Then, unfreeze the next layer and fine-tune the last-two layers only *for the next epoch*. Continue until the entire network is being trained, at which time we just train until convergence.
- **BPTT for Text Classification**. Divide documents into fixed-length “batches”<sup>153</sup> of size  $b$ . They initialize the  $i$ th section with the final state of the model run on section  $i - 1$ .

---

<sup>151</sup>Steps-per-epoch times number of epochs.

<sup>152</sup>Or just as much as we can fit into GPU memory.

<sup>153</sup>Not a fan of how they overloaded this term here.

## The Marginal Value of Adaptive Gradient Methods in Machine Learning

Table of Contents Local

Written by Brandon McKinzie

Wilson et al., “The Marginal Value of Adaptive Gradient Methods in Machine Learning,” (May 2018).

**TL;DR.** For simple overparameterized problems, adaptive methods (a) often find drastically different solutions than SGD, and (b) tend to give undue influence to spurious features that have no effect on out-of-sample generalization. They also found that tuning the initial learning and decay scheme for Adam yields significant improvements over its default settings in all cases.

**Background (2).** The gradient updates for general stochastic gradient, stochastic momentum, and adaptive gradient methods, respectively, can be formalized as follows<sup>154</sup>

$$[\text{regular}] \quad \Delta w_{k+1} = -\alpha_k \tilde{\nabla} f(w_k) \quad (675)$$

$$[\text{momentum}] \quad \Delta w_{k+1} = -\alpha_k \tilde{\nabla} f(w_k + \gamma_k \Delta w_k) + \beta_k \Delta w_k \quad (676)$$

$$[\text{adaptive}] \quad \Delta w_{k+1} = -\alpha_k \mathbf{H}_k^{-1} \tilde{\nabla} f(w_k + \gamma_k \Delta w_k) + \beta_k \mathbf{H}_k^{-1} \mathbf{H}_{k-1} \Delta w_k \quad (677)$$

where  $H_k$  is a p.d. matrix involving the entire sequence of iterates  $(w_1, \dots, w_k)$ . For example, regular momentum would be  $\gamma_k=0$ , and Nesterov momentum would be  $\gamma_k=\beta_k$ . In practice, we basically always define  $H_k$  as the diagonal matrix:

$$H_k := \text{diag} \left[ \left( \sum_{i=1}^k \eta_i \mathbf{g}_i \odot \mathbf{g}_i \right)^{1/2} \right] \quad (678)$$

**The Potential Perils of Adaptivity (3).** Consider the binary least-squares classification problem, where we aim to minimize

$$R_s[w] := \frac{1}{2} \|Xw - y\|_2^2 \quad (679)$$

where  $X \in \mathbb{R}^{n \times d}$  and  $y \in \{-1, 1\}^n$ .

**Lemma 3.1**

If there exists a scalar  $c$  s.t.  $X \text{sign}(X^T y) = cy$ , then (assuming  $w_0 := 0$ ) AdaGrad, Adam, and RMSProp all converge to the unique solution  $w \propto \text{sign}(X^T y)$ .

---

<sup>154</sup>I'm defining  $\Delta w_{k+1} \triangleq w_{k+1} - w_k$ .

## A Theoretically Grounded Application of Dropout in Recurrent Neural Networks

Table of Contents Local

Written by Brandon McKinzie

Y. Gal and Z. Ghahramani, “A Theoretically Grounded Application of Dropout in Recurrent Neural Networks,” *University of Cambridge* (Oct 2016).

**Background** (3). In Bayesian regression, we want to infer the parameters  $\omega$  of some function  $y = f^\omega(x)$ . We define a prior,  $p(\omega)$ , and a likelihood,

$$p(y=d | x, \omega) = \text{Cat}_d(\text{softmax}(f^\omega(x))) \quad (680)$$

for a classification setting. Given a dataset  $X, Y$ , and some new point  $x^*$ , we can predict its output via

$$p(y^* | x^*, X, Y) = \int p(y^* | x^*, \omega) p(\omega | X, Y) d\omega \quad (681)$$

In a **Bayesian neural network**, we place the prior over the NNs weights (typically Gaussians). The posterior  $p(\omega | X, Y)$  is usually intractable, so we resort to **variational inference** to approximate it. We define our approximating distribution  $q(\omega)$  and aim to minimize the KLD:

$$\text{KL}(q(\omega) || p(\omega | X, Y)) \propto - \int q(\omega) \log p(Y | X, \omega) d\omega + \text{KL}(q(\omega) || p(\omega)) \quad (682)$$

$$= \sum_{i=1}^N \int q(\omega) \log p(y_i | f^\omega(x_i)) d\omega + \text{KL}(q(\omega) || p(\omega)) \quad (683)$$

**Variational Inference in RNNs** (4). The authors use MC integration to approximate the integral. They use only a single sample  $\hat{\omega} \sim q(\omega)$  for each of the  $N$  summations, resulting in an unbiased estimator. Plugging this in, we obtain our objective:

$$\mathcal{L} \approx - \sum_{i=1}^N \log p(y_i | f_y^{\hat{\omega}_i} (f_h^{\hat{\omega}_i}(x_{i,T}, h_{T-1}))) + \text{KL}(q(\omega) || p(\omega)) \quad (684)$$

The crucial observations here are:

- For each sequence  $x_i$ , we sample a new realization  $\hat{\omega}_i$ .
- For each of the  $T$  symbols in  $x_i$ , we use that *same* realization.

We define our approximating distribution to factorize over the weight matrices and their rows in  $\omega$ . For each weight matrix row  $w_k$ , we have

$$q(w_k) \triangleq p\mathcal{N}(w_k; \mathbf{0}, \sigma^2 I) + (1-p)\mathcal{N}(w_k; m_k, \sigma^2 I) \quad (685)$$

with  $m_k$  **variational parameter** (row vector).

## Improving Neural Language Models with a Continuous Cache

Table of Contents Local

Written by Brandon McKinzie

E. Grave, A. Joulin, and N. Usunier, “Improving Neural Language Models with a Continuous Cache,” *Facebook AI Research* (Dec 2016).

The cache stores pairs  $(h_t, x_{t+1})$  of the final hidden-state representation at time  $t$ , along with the word which was *generated*<sup>155</sup> based on this representation.

$$p_{vocab}(w | \mathbf{x}_{\langle 1 \dots t \rangle}) \propto \exp(h_t^T o_w) \quad (686)$$

$$p_{cache}(w | \mathbf{h}_{\langle 1 \dots t \rangle}, \mathbf{x}_{\langle 1 \dots t \rangle}) \propto \sum_{i=1}^{t-1} \mathbb{1}_{x_{i+1}=w} \exp(\theta h_i^T h_i) \quad (687)$$

$$= \sum_{\substack{(x,h) \in \text{cache} \\ s.t. x=w}} \exp(\theta h^T h) \quad (688)$$

$$p(w | \mathbf{h}_{\langle 1 \dots t \rangle}, \mathbf{x}_{\langle 1 \dots t \rangle}) = (1 - \lambda)p_{vocab}(w | h_t) + \lambda p_{cache}(w | \mathbf{h}_{\langle 1 \dots t \rangle}, \mathbf{x}_{\langle 1 \dots t \rangle}) \quad (689)$$

where  $\theta$  is a scalar parameter that controls the flatness of the cache distribution.

---

<sup>155</sup>They say this, but everything else in the paper strongly suggests they mean the next gold-standard input instead.

## Protection Against Reconstruction and Its Applications in Private Federated Learning

Table of Contents Local

Written by Brandon McKinzie

A. Bhowmick et al., “Protection Against Reconstruction and Its Applications in Private Federated Learning,” *Apple, Inc.* (Dec 2018).

**Introduction** (1). In many scenarios, it is possible to reconstruct model inputs  $x$  given just  $\nabla_{\theta} \ell(\theta; x, y)$ . **Differential privacy** is one approach for obscuring the gradients such that guarantees can be made regarding risk of compromising user data  $x$ . **Locally private** algorithms, however, are preferred to DP when the user wants to keep their data private even from the data collector. The authors want to find a way to perform SGD while providing both local privacy to individual data  $X_i$  and stronger guarantees on the global privacy of the output  $\hat{\theta}_n$  of their procedure.

Formally, say we have two users’ data  $x$  and  $x'$  (both in  $\mathcal{X}$ ) and some randomized mechanism  $M : \mathcal{X} \mapsto \mathcal{Z}$ . We say that  $M$  is  **$\varepsilon$ -local differentially private** if  $\forall x, x' \in \mathcal{X}$  and sets  $S \subset \mathcal{Z}$ :

$$\frac{\Pr[M(x) \in S]}{\Pr[M(x') \in S]} \leq e^{\varepsilon} \quad (690)$$

Clearly, the RHS will need to be pretty big for this to be achievable. The authors claim that allowing  $\varepsilon \gg 1$  “may [still] provide meaningful privacy protections.”

**Privacy Protections** (2). The focus here is on the **curious onlooker**: an individual (e.g. Apple PFL engineer) who can observe all updates to a model and communication from individual devices. Let  $X$  denote some user data. Let  $\Delta W$  denote the weights difference after some model update using the data  $X$ . Let  $Z$  be the result of the randomized mapping  $\Delta W \mapsto Z$ . Our setting can be described with the Markov chain  $X \rightarrow \Delta W \rightarrow Z$ . The onlooker observes  $Z$  and wants to estimate some function  $f(X)$  on the private data.

**Separated private mechanisms** (2.2). The authors propose, instead of a simple mapping  $\Delta W \rightarrow Z$ , to split it up into two parts:  $Z_1 = M_1(U)$  and  $Z_2 = M_2(R)$ , where

$$U = \frac{\Delta W}{\|\Delta W\|_2} \quad (691)$$

$$R = \|\Delta W\|_2 \quad (692)$$

**Separated Differential Privacy**

*A pair of mechanisms  $M_1, M_2$  mapping from  $\mathcal{U} \times \mathcal{R}$  to  $\mathcal{Z}_1 \times \mathcal{Z}_2$  is  $(\varepsilon_1, \varepsilon_2)$ -separated differentially private if  $M_1$  is  $\varepsilon_1$ -locally differentially private and  $M_2$  is  $\varepsilon_2$ -locally differentially private.*

**Privatizing Unit  $\ell_2$  Vectors with High Accuracy** (4.1). Given some vector  $u \in \mathbb{S}^{d-1}$ <sup>156</sup>, we want to generate an  $\varepsilon$ -differentially private vector  $Z$  such that

$$\mathbb{E}[Z | u] = u \quad \forall u \in \mathbb{S}^{d-1} \quad (693)$$

### Privatized Unit Vector: `PrivUnit2`

Sample random vector  $V$ :

$$V \sim \begin{cases} U(\{v \in \mathbb{S}^{d-1} | \langle v, u \rangle \geq \gamma\}) & \text{with probability } p \\ U(\{v \in \mathbb{S}^{d-1} | \langle v, u \rangle < \gamma\}) & \text{otherwise} \end{cases} \quad (694)$$

where  $\gamma \in [0, 1]$  and  $p \geq \frac{1}{2}$  together control *accuracy* and *privacy*.

Let  $\alpha = \frac{d-1}{2}$ ,  $\tau = \frac{1+\gamma}{2}$ , and

$$m = \frac{(1-\gamma^2)\alpha}{2^{d-2}(d-1)} \left[ \frac{p}{B(\alpha, \alpha) - B(\tau; \alpha, \alpha)} - \frac{1-p}{B(\tau; \alpha, \alpha)} \right] \quad (695)$$

where  $B(x, \alpha, \beta)$  is the incomplete beta function (see paper pg 17 for details).

Return  $Z = \frac{1}{m} \cdot V$

**Privatizing the Magnitude** (4.3). We also need to privatize the weight delta norms. We want to return values  $r \in [0, r_{max}]$  for some  $r_{max} < \infty$ .

---

<sup>156</sup>Here, this denotes an n-sphere:

$$\mathbb{S}^n \triangleq \{x \in \mathbb{R}^{n+1} : \|x\| = r\}$$

## Context Dependent RNN Language Model

Table of Contents Local

Written by Brandon McKinzie

T. Mikolov and G. Zweig, “Context Dependent Recurrent Neural Network Language Model,” *BRNO and Microsoft* (2012).

**Model Structure (2).** Given one-hot input vector  $\mathbf{x}_t$ , output a probability distribution  $\mathbf{y}_t$  for the next word. Incorporate a *feature vector*  $\mathbf{f}_t$  that will contain topic information.

$$\mathbf{y}_t = \text{Softmax}(\mathbf{V}\mathbf{h}_t + \mathbf{G}\mathbf{f}_t) \quad (696)$$

$$\mathbf{h}_t = \sigma(\mathbf{U}\mathbf{x}_t + \mathbf{W}\mathbf{h}_{t-1} + \mathbf{F}\mathbf{f}_t) \quad (697)$$

**LDA for Context Modeling (3).** “Documents” fed to LDA here will be individual sentences. The generative process assumed by LDA is compactly defined by the following sequence of operations<sup>157</sup>:

$$N \sim \text{Poisson}(\xi) \quad (698) \quad \begin{matrix} \text{N: number of words} \\ \Theta_i \equiv p(\text{topic}[i]) \end{matrix}$$

$$\Theta \sim \text{Dir}(\alpha) \quad (699) \quad z_n: \text{topic of word } n$$

$$z_n \sim \text{Multinomial}(\Theta) \quad (700)$$

$$w_n \sim \text{Pr}[w_n | z_n, \beta] \quad (701)$$

where  $\text{Pr}[w_n=a | z_n=b] = \beta_{b,a}$ , so we are really just sampling from row  $z_n$  of  $\beta$ , where  $\beta \in [0, 1]^{Z \times V}$  (where Z is number of topics). The result of LDA is a learned value for  $\alpha$ , and the topic distributions  $\beta$ .

$$\mathbf{f}_t = \frac{1}{Z} \prod_{i=0}^K \mathbf{t}_{x_{t-i}} \quad (702)$$

$$\mathbf{f}_t = \frac{1}{Z} \mathbf{f}_{t-1}^\gamma \mathbf{t}_{x_t}^{1-\gamma} \quad (703)$$

---

<sup>157</sup> $\alpha$  is a vector with number of elements equal to number of topics.

## Strategies for Training Large Vocabulary Neural Language Models

Table of Contents Local

Written by Brandon McKinzie

Chen et al., “Strategies for Training Large Vocabulary Neural Language Models,” *FAIR* (Dec 2018). arXiv:1512.04906

**Setup/Notation.** Note that in everything below, the authors are using a rather primitive feed-forward network as their language model. To predict  $w_t$  it just concatenates the embeddings of the previous  $n$  words and feeds it through a  $k$ -layer FF network. Then, layer  $k+1$  is the dense projection and softmax:

$$h^{k+1} = W^{k+1}h^k + b^{k+1} \in \mathbb{R}^V \quad (704)$$

$$y = \frac{1}{Z} \exp \left\{ h^{k+1} \right\} \quad (705)$$

Using cross-entropy loss, the derivative of  $\log p(w_t=i)$  wrt the  $j$ th element of the logits is:

$$\frac{\partial \log y_i}{\partial h_j^{k+1}} = \frac{\partial}{\partial h_j^{k+1}} \left[ h_i^{k+1} - \log Z \right] \quad (706)$$

$$= \delta_{ij} - y_j \quad (707)$$

When computing gradients of the cross-entropy loss,  $y_i$  here is the ground truth. Therefore, to increase the probability of the correct token, we need to increase the logits element for that index, and decrease the elements for the others. Note how this implies we must compute the final activations for *all words in the vocabulary*.

**Hierarchical Softmax** (2.2). Group words into one of two clusters  $\{c_1, c_2\}$ , based on unigram frequency<sup>158</sup>. Then model  $p(w_t | x) = p(c_t | x)p(w_t | c_t)$  where  $c_t$  is the class that word  $w_t$  was assigned to.

---

<sup>158</sup>For example, you could just put the top 50% in  $c_1$  and the rest in  $c_2$ .

**NCE** (2.5). Define  $P_{noise}(w)$  by the unigram frequency distribution. For each real token  $w_t$  in the training set, sample  $K$  noise tokens  $\{n_k\}_{k=1}^K$ . NCE aims to minimize

$$L_{NCE}(\{\mathbf{w}_1, \dots, \mathbf{w}_N\}) = \sum_{i=1}^N \sum_{t=1}^{len(\mathbf{w}^{(i)})} \left[ \log h(w_t^{(i)}) + \sum_{k=1}^K \log(1 - h(n_k^{(i)})) \right] \quad (708)$$

$$h(w_t) = \frac{P_{model}(w)}{P_{model}(w) + kP_{noise}(w)} \quad (709)$$

$$\approx \frac{\tilde{P}_{model}(w)}{\tilde{P}_{model}(w) + kP_{noise}(w)} \quad (710)$$

where the final approximation is what makes NCE less computationally expensive in practice than standard softmax. This would seem to imply that NCE should approach standard softmax (in terms of correctness) as  $k$  increases.

### Takeaways.

- Hierarchical softmax is the fastest.
- NCE performs well on large-volume large-vocab datasets.
- Similar NCE values can result in very different validation perplexities.
- Sampled softmax shows good results if the number of negative samples is at 30% of the vocab size or larger.
- Sampled softmax has a lower ppl reduction per step than others.

## Product quantization for nearest neighbor search

Table of Contents Local

Written by Brandon McKinzie

Jégou, “Product quantization for nearest neighbor search,” (2011)

**Vector Quantization.** Denote the *index set*  $\mathcal{I} = [0..k - 1]$  and the set of reproduction values (a.k.a. **centroids**)  $c_i$  as  $\mathcal{C} = \{\mathbf{c}_i \in \mathbb{R}^D : i \in \mathcal{I}\}$ . We refer to  $\mathcal{C}$  as the **codebook** of size  $k$ . A **quantizer** is a function  $q : \mathbf{x} \mapsto \mathbf{q}(\mathbf{x})$ , where  $\mathbf{x} \in \mathbb{R}^D$  and  $\mathbf{q}(\mathbf{x}) \in \mathcal{C}$ . We typically evaluate the quality of a quantizer with mean squared error of the reconstruction:

$$MSE(q) = \mathbb{E}_{\mathbf{x} \sim p(\mathbf{x})} \left[ \|\mathbf{x} - \mathbf{q}(\mathbf{x})\|_2^2 \right] \quad (711)$$

In order for an optimizer to be optimal, it must satisfy the **Lloyd optimality conditions**:

$$(1) \quad \mathbf{q}(\mathbf{x}) = \arg \min_{c_i \in \mathcal{C}} \|\mathbf{x} - \mathbf{c}_i\|_2 \quad (712)$$

$$(2) \quad \mathbf{c}_i = \mathbb{E}_{\mathbf{x}} [\mathbf{x} \mid i] \triangleq \int_{\mathcal{V}_i} \mathbf{x} p(\mathbf{x}) d\mathbf{x} \quad (713)$$

a.k.a. literally just K-means.

**Product Quantization.** Input vector  $\mathbf{x} \in \mathcal{R}^D$  is split into  $m$  distinct subvectors  $\mathbf{u}_j \in \mathbb{R}^{D/m}$ , where  $j \in [1..m]$ . Note that  $D$  must be an integer multiple of  $m$  (i.e.  $D = am$  for some  $a \in \mathbb{Z}$ ).

$$\mathbf{x} = \underbrace{\mathbf{x}_{\langle 1 \dots D^* \rangle}}_{\mathbf{u}_1(\mathbf{x})}, \dots, \underbrace{\mathbf{x}_{\langle D-D^*+1 \dots D \rangle}}_{\mathbf{u}_m(\mathbf{x})} \rightarrow \mathbf{q}_1(\mathbf{u}_1(\mathbf{x})), \dots, \mathbf{q}_m(\mathbf{u}_m(\mathbf{x})) \quad (714)$$

Note that each subquantizer  $q_j$  has its own index set  $\mathcal{I}_j$  and codebook  $\mathcal{C}_j$ . Therefore, the final reproduction value of a product quantizer is identified by an element of the product set  $\mathcal{I} = \mathcal{I}_1 \times \dots \times \mathcal{I}_m$ . The associated final codebook is  $\mathcal{C} = \mathcal{C}_1 \times \dots \times \mathcal{C}_m$ .

## Large Memory Layers with Product Keys

Table of Contents

Local

Written by Brandon McKinzie

Lample et al., “Large Memory Layers with Product Keys,” *FAIR* (July 2019). arXiv:1907.05242v1

**Memory Design** (3.1). The high-level structure (sequence of ops) is as follows:

1. **Query network** computes some query vector  $\mathbf{q}$ .
2. Compare  $\mathbf{q}$  with each **product key** via some scoring function.
3. Select the  $k$  highest scoring product keys.
4. Compute output  $\mathbf{m}(\mathbf{x})$  as weighted sum over the values associated with each of the top  $k$  keys from the previous step.

The **query network** is usually just a dense layer<sup>159</sup>. Since they want it to output query vectors with good coverage over the key space, they put a batch normalization layer before the query network<sup>160</sup>.

The *standard* way of doing key assignment/weighting is as follows:

$$[\text{KNN}] \quad \mathcal{I} \triangleq \text{TopK}(\mathbf{q}(\mathbf{x})^T \mathbf{k}_i) \quad 1 \leq i \leq \mathcal{K} \quad (715)$$

$$[\text{Normalize}] \quad \mathbf{w} = \text{Softmax}(\mathcal{I}) \quad (716)$$

$$[\text{Aggregate}] \quad \mathbf{m}(\mathbf{x}) = \sum_{i \in \mathcal{I}} w_i \mathbf{v}_i \quad (717)$$

where equation 715 is inefficient for large memory (key-value) stores. The authors propose instead a structured set of keys that they call **product keys**. Spoiler alert: it’s just product quantization with  $m=2$  subvectors. Instead of using the flat key set  $\mathcal{K} \triangleq \{\mathbf{k}_1, \dots, \mathbf{k}_{|\mathcal{K}|}\}$  with each  $\mathbf{k}_i \in \mathbb{R}^{d_q}$  from earlier, we redefine it as

$$\mathcal{K} \triangleq \{(\mathbf{c}, \mathbf{c}') \mid c \in \mathcal{C}, c' \in \mathcal{C}'\} \quad (718)$$

where both  $\mathcal{C}$  and  $\mathcal{C}'$  are sets of *sub-keys*  $\mathbf{k}_i \in \mathbb{R}^{d_q/2}$ .

Then...

1. Just run each of subvectors  $q_1$  and  $q_2$  through the standard TopK. You’ll have  $k$  sub-keys for both, defined by their index into their respective codebook.
2. Let  $\mathcal{K} := \{(\mathbf{c}_i, \mathbf{c}'_j) \mid i \in \mathcal{I}_{\mathcal{C}}, j \in \mathcal{I}_{\mathcal{C}'}\}$ . This new reduced-size key set  $\mathcal{K}$  has only  $k \times k$  entries.

<sup>159</sup>They choose  $d_q = 512$  as the output dimensionality of their query network.

<sup>160</sup>Recall that batch norm just normalizes the batch inputs to have 0 mean and unit standard deviation, followed by a scaling and bias factor.

3. Run the standard algorithm using the new reduced key set  $\mathcal{K}$ .

**TODO:** finish this note

## Show, Ask, Attend, and Answer

Table of Contents Local

Written by Brandon McKinzie

V. Kazemi and A. Elqursh, “Show, Ask, Attend, and Answer: A Strong Baseline For Visual Question Answering”  
*Google Research* (April 2017). arXiv:1704.03162v2

**TL;DR:** Good for a high-level overview of the VQA task, but extremely vague with so many details omitted it renders the paper fairly useless.

**Method (3).** Given a training set of image-question-answer triplets  $(I, q, a)$ , learn to estimate the most likely answer  $\hat{a}$  out of the set of most frequent answers<sup>161</sup> in the training data:

$$\hat{a} = \arg \max_a \Pr [a | I, q] \quad (719)$$

The method utilizes the following architectural components:

- **Image Embedding** (3.1). Extracts features  $\phi = \text{CNN}(I)$ .
- **Question Embedding** (3.2). Encode question  $q$  as the final state of LSTM:  $s = \text{LSTM}(\text{Embed}(q))$ .
- **Stacked Attention** (3.3)<sup>162</sup> Seems like they feed  $\text{Concat}[s, \phi]$  through two layers of convolution to produce an output  $F_c$  for  $c \in [1..C]$  (meaning they do  $C$  such convolutions separately and in parallel, like multi-head attention). This represents the scores for the attention function. The attention output, as usual, is computed as

$$\mathbf{x}_c = \sum_{\ell} \alpha_{c,\ell} \phi_{\ell} \quad (720)$$

$$\alpha_{c,\ell} \propto \exp F_c(s, \phi_{\ell}) \quad (721)$$

where  $\ell$  appears to be over all [flattened] spatial indices of  $\phi$ .

- **Classifier** (3.4). Concat the **image glimpses**  $\mathbf{x}_c$  with the LSTM output  $s$  and feed through a couple FC layers to eventually obtain softmax probabilities over each possible answer  $a_i$ ,  $i \in [1..M]$ .

---

<sup>161</sup>Same approach as how we define vocabulary in language modeling tasks.

<sup>162</sup>Authors do a laughably poor job at describing this part in any detail, so I’m taking the liberty of filling in the blanks. Blows my mind that papers this sloppy can even be published.

**Dataset.** Although, again, the authors are horribly vague/sloppy here, it *seems* like the data they use actually provides  $K$  “correct” answers for each image-question pair. The model loss is therefore an average NLL over the  $K$  true classes.

## Did the Model Understand the Question?

Table of Contents   Local

*Written by Brandon McKinzie*

Mudrakarta et al., “Did the Model Understand the Question?” *Univ. Chicago & Google Brain* (May 2018). arXiv:1805.05492v1

**TL;DR.** Basically all QA-related networks are dumb and don’t learn what we think they learn.

- Networks tend to make predictions based a tiny subset of the input words. Due to this, altering the non-important words in ways that may drastically change the meaning of the question can have virtually no impact on the network’s prediction.
- Networks assign high-importance to words like “there”, “what”, “how”, etc. These are actually low-information words that the network should not heavily rely on.
- Networks rely on the image far more than the question.

**Integrated Gradients (IG)** (3). Purpose: “isolate question words that a DL system uses to produce an answer.”

$$F(\mathbf{x} = \{x_1, \dots, x_n\}) \in [0, 1] \quad (722)$$

$$A_F(\mathbf{x}, \mathbf{x}') = \{a_1, \dots, a_n\} \in \mathbb{R}^n \quad (723)$$

where  $\mathbf{x}'$  is some baseline input we use to compute the relative attribution of input  $x$ . The authors set  $\mathbf{x}'$  as the “empty question” (sequence of padding values)<sup>163</sup>.

*Given an input  $x$  and baseline  $x'$ , the **integrated gradient** along the  $i$ th dimension is as follows.*

$$IG_i(x, x') \triangleq (x_i - x'_i) \times \int_{\alpha=0}^1 \frac{\partial F(x' + \alpha \times (x - x'))}{\partial x_i} d\alpha \quad (724)$$

Interpretation: seems like IG gives us a better idea of the *total* “attribution” of each input dimension  $x_i$  relative to baseline  $x'_i$  along the line connecting  $x_i$  and  $x'_i$ , instead of just the immediate derivative around  $x_i$ . Although, the fact that infinitesimal contributions could cancel each other out seems problematic (positive and negative gradients along the interpolation).

---

<sup>163</sup>The use the same context though (e.g. the associated image for VQA). Only the question is changed.

## XLNet

Table of Contents Local

Written by Brandon McKinzie

Yang et al., “XLNet: Generalized Autoregressive Pretraining for Language Understanding” *CMU & Google Brain* (May 2018).

**TL;DR:** Instead of minimizing the NLL using  $p(w_1, \dots, w_T)$ , minimize over NLL’s using every possible order of the given word sequence.

**Background.** Recall that BERT does denoising auto-encoding. Given text sequence  $\mathbf{x} = \{x_1, \dots, x_T\}$ , BERT constructs a corrupted version  $\hat{\mathbf{x}}$  by randomly masking out some tokens. Let  $\bar{\mathbf{x}}$  denote the tokens that were masked. The BERT training objective is then

$$[\text{BERT}] \quad \max_{\theta} \log p_{\theta}(\bar{\mathbf{x}} | \hat{\mathbf{x}}) \approx \sum_{\bar{x} \in \bar{\mathbf{x}}} \log p_{\theta}(\bar{x} | \hat{\mathbf{x}}) \quad (725)$$

$$p(\bar{x} | \hat{\mathbf{x}}) = \text{Softmax} \left( H_{\theta}(\hat{\mathbf{x}})^T e(\bar{x}) \right) \quad (726)$$

**Objective & Architecture.** Their proposed **permutation language modeling objective** is:

$$\max_{\theta} \mathbb{E}_{\mathbf{z} \sim \mathcal{Z}_T} \left[ \sum_{t=1}^T \log p_{\theta}(x_{z_t} | \mathbf{x}_{\langle z_1 \dots z_{t-1} \rangle}) \right] \quad (727)$$

where  $\mathcal{Z}_T$  is the set of all possible permutations of the length- $T$  index sequence  $[1..T]$ . To implement this, the authors had to re-parameterize the next-token distribution to be **target position aware**:

$$p_{\theta}(X_{z_t} = x | \mathbf{x}_{\langle z_1 \dots z_{t-1} \rangle}) = \text{Softmax} \left( g_{\theta} \left( \mathbf{x}_{\langle z_1 \dots z_{t-1} \rangle}, \textcolor{red}{z_t} \right)^T e(x) \right) \quad (728)$$

They accomplish this via **two-stream self-attention**, a technique that utilizes two sets of hidden representations (instead of one):

- **Content representation:**  $h_{z_t} \triangleq h_{\theta}(\mathbf{x}_{\langle z_1 \dots z_t \rangle})$ .
- **Query representation:**  $g_{z_t} \triangleq g_{\theta}(\mathbf{x}_{\langle z_1 \dots z_{t-1} \rangle}, z_t)$ .

The query stream is initialized with some vector  $g_i^{(0)}=w$ , and the content stream is initialized with word embedding  $h_i^{(0)}=e(x_i)$ . For the subsequent attention layers  $1 \leq m \leq M$ , they are computed respectively as follows:

$$g_{z_t}^{(m)} \leftarrow \text{Attention}(Q=g_{z_t}^{(m-1)}, K=V=\mathbf{h}_{z_{<t}}^{(m-1)}) \quad (729)$$

$$h_{z_t}^{(m)} \leftarrow \text{Attention}(Q=h_{z_t}^{(m-1)}, K=V=\mathbf{h}_{z_{\leq t}}^{(m-1)}) \quad (730)$$

In practice, in order to speed up optimization, the authors do **partial prediction**: only train to predict over  $\mathbf{x}_{z_{>c}}$  targets rather than all of them.

**Incorporating Ideas from Transformer-XL.** Often times, sequences are too long to feed all at once. The authors adopt relative positional encoding and segment-level recurrence from Transformer-XL. To compute the attention update with memory on a given segment, we use the content representations from the *previous* segment,  $\tilde{\mathbf{h}}$ , along with the current segment,  $\mathbf{h}_{z_{\leq t}}$  as follows:

$$h_{z_t}^{(m)} \leftarrow \text{Attention}\left(Q=h_{z_t}^{(m-1)}, K=V=\left[\tilde{\mathbf{h}}^{(m-1)}; \mathbf{h}_{z_{\leq t}}^{(m-1)}\right]\right) \quad (731)$$

## Transformer-XL

Table of Contents Local

Written by Brandon McKinzie

Dai et al., “Transformer-XL: Attentive Language Models Beyond a Fixed-Length Context” *CMU & Google Brain* (Jan 2019).

**Segment-Level Recurrence with State Reuse.** Denote two consecutive segments of length  $L$  as  $\mathbf{s}_\tau = [x_{\tau,1}, \dots, x_{\tau,L}]$  and  $\mathbf{s}_{\tau+1} = [x_{\tau+1,1}, \dots, x_{\tau+1,L}]$ . Denote output of layer  $n$  given input segment  $\mathbf{s}_\tau$  as  $\mathbf{h}_\tau^n \in \mathbb{R}^{L \times d}$ , where  $d$  is the hidden dimension. To obtain the output of layer  $n$  given the next segment,  $\mathbf{s}_{\tau+1}$ , do:

$$\mathbf{h}_{\tau+1}^n = \text{TransformerLayer}(\mathbf{q}_{\tau+1}^n, \mathbf{k}_{\tau+1}^n, \mathbf{v}_{\tau+1}^n) \quad (732)$$

$$= \text{TransformerLayer}(\mathbf{h}_{\tau+1}^{n-1} \mathbf{W}_q^T, \tilde{\mathbf{h}}_{\tau+1}^{n-1} \mathbf{W}_k^T, \tilde{\mathbf{h}}_{\tau+1}^{n-1} \mathbf{W}_v^T) \quad (733)$$

$$\tilde{\mathbf{h}}_{\tau+1}^{n-1} = [\text{SG}(\mathbf{h}_\tau^{n-1}); \mathbf{h}_{\tau+1}^{n-1}] \quad (734)$$

where the concat in 734 is along the length (time) dimension. In other words,  $\mathbf{Q}$  remains the same, but  $\mathbf{K}$  and  $\mathbf{V}$  get the previous segment prepended. Ultimately this only changes the inner dot products in the attention mechanism to attend over both segments. The  $L$  output attention vectors are therefore each weighted sums over the previous  $2L$  timesteps instead of just  $L$ .

**Relative Positional Encodings.** Instead of absolute positional encodings (as regular transformers do), only encode the *relative* positional information in the hidden states. Ignoring the scale factor of  $1/\sqrt{d_k}$ , we can write the score for query vector  $q_i = W_q(e_{x_i} + u_i)$  and key vector  $k_j = W_k(e_{x_j} + u_j)$ , for input embeddings  $e$  and positional encodings  $u$  as follows. Below it we show the authors proposed re-parameterized relative encoding version.

$$A_{i,j}^{abs} = e_{x_i}^T W_q^T W_k e_{x_j} + e_{x_i}^T W_q^T W_k u_j + u_i^T W_q^T W_k e_{x_j} + u_i^T W_q^T W_k u_j \quad (735)$$

$$A_{i,j}^{res} = \underbrace{e_{x_i}^T W_q^T W_k e_{x_j}}_{\text{cont-based addr}} + \underbrace{e_{x_i}^T W_q^T W_k, R r_{i-j}}_{\text{cont-dep pos bias}} + \underbrace{u_i^T W_k, E e_{x_j}}_{\text{global cont bias}} + \underbrace{u_i^T W_k, R r_{i-j}}_{\text{global pos bias}} \quad (736)$$

where content is abbreviated as “cont” and positional is abbrev as “pos”. I’ve shown all differences introduced by the second version in red font. It appears that  $r_{i-j}$  is literally just  $u_{i-j}$  but I guess using new letters is cool. Note that they separate  $W_k$  into **content-based**  $W_{k,E}$  and **location-based**  $W_{k,R}$ .

## Efficient Softmax Approximation for GPUs

Table of Contents   Local

*Written by Brandon McKinzie*

Grave et al., “Efficient Softmax Approximation for GPUs” *FAIR* (June 2017).

### Adaptive Softmax: Two-Level

Partition the vocabulary  $\mathcal{V}$  into two clusters  $\mathcal{V}_h$  and  $\mathcal{V}_t$ , where

- $\mathcal{V}_h$  denotes the *head*, consisting of the most frequent words.
- $\mathcal{V}_t$  denotes the *tail*, associated with a **large number** of rare words.
- $|\mathcal{V}_h| << |\mathcal{V}_t|$  and  $P(\mathcal{V}_h) >> P(\mathcal{V}_t)$ .

To compute the probability of some word  $w$  given context  $h$ , do:

$$\Pr[w | h] = \begin{cases} P_{\mathcal{V}_h}(w | h) & \text{if } w \in \mathcal{V}_h \\ P_{\mathcal{V}_t}(w | h)P_{\mathcal{V}_h}(\text{tail} | h) & \text{otherwise} \end{cases} \quad (737)$$

where both  $P_{\mathcal{V}_h}$  and  $P_{\mathcal{V}_t}$  are modeled with a softmax over the words in their respective clusters ( $P_{\mathcal{V}_h}$  also includes the special “tail” token).

More generally, we can extend the above algorithm to  $N$  clusters (instead of 2). We can also adapt the *capacity* of each cluster (varying their embedding size). The authors recommend, for each successive tail cluster, reducing the output size by a factor of 4. Of course, this then has to be followed by projecting back up to the number of words associated with the given cluster.

**TODO:** detail out how cross entropy loss is computed under this setup.

## Adaptive Input Representations for Neural Language Modeling

Table of Contents Local

*Written by Brandon McKinzie*

A. Baevski and M. Auli, “Adaptive Input Representations for Neural Language Modeling” *FAIR* (Feb 2019).

**TL;DR:** Literally just adaptive softmax but for the input embeddings. Official implementation can be found [here](#).

**Adaptive Input Representations** (3). Same as Grave et al., they partition the vocabulary  $\mathcal{V}$  into

$$\mathcal{V} = \mathcal{V}_1 \cup \mathcal{V}_2 \cup \dots \cup \mathcal{V}_n \quad (738)$$

where  $\mathcal{V}_1$  is the head and the rest are the tails (ordered by decreasing frequency). They reduce the capacity of each cluster by a factor of  $k=4$  (also same as Grave et al.). Finally, they add linear projections for each cluster’s embeddings in order to ensure they all result in  $d$ -dimensional output embeddings (even  $\mathcal{V}_1$ ).

## Neural Module Networks

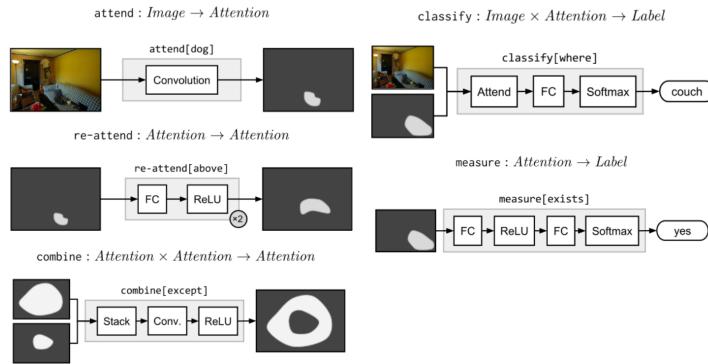
Table of Contents Local

Written by Brandon McKinzie

Andreas et al., “Deep Compositional Question Answering with Neural Module Networks” *UC Berkeley* (Nov 2015).

**NMNs for Visual QA** (4). Model and task overview:

- **Data:** 3-tuples  $(w, x, y)$  containing the question, image, and answer, respectively.
- **Model:** fully specified by a collection of **modules**  $\{m\}$ . Each module  $m$  has parameters  $\theta_m$  and a **network layout predictor**  $P(w)$  that maps from strings to networks. The high-level procedure is, for each  $(w, x, y)$ , do:
  1. Instantiate a network based on  $P(w)$ .
  2. Pass the image  $x$  (and possibly  $w$  again) as inputs to the network.
  3. Obtain network outputs encoding  $p(y | w, x; \theta)$ .
- **Modules:**



**From strings to networks** (4.2).

1. Parse question  $w$  with the Stanford Parser to obtain universal dependency representation.
2. Filter dependencies to those connected to the wh-word in the question. Some examples:
  - *what is standing in the field*  $\mapsto$  `what(stand)`
  - *what color is the truck*  $\mapsto$  `color(truck)`
  - *is there a circle next to a square*  $\mapsto$  `is(circle, next-to(square))`
3. Assign identities of modules (already have full network structure).
  - Leaves become **attend** modules.
  - Internal nodes become **re-attend** or **combine** modules.
  - Root nodes become **measure** (y/n questions) or **classify** (everything else) modules.

**Answering natural language questions** (4.3). They combine the results from the module network with an LSTM, which is fed the question as input and outputs a predictive distribution over the set of answers<sup>164</sup>. The final prediction is a geometric average of the LSTM output probabilities and the root module output probabilities.

---

<sup>164</sup>This is the same distribution that the root module is trying to predict

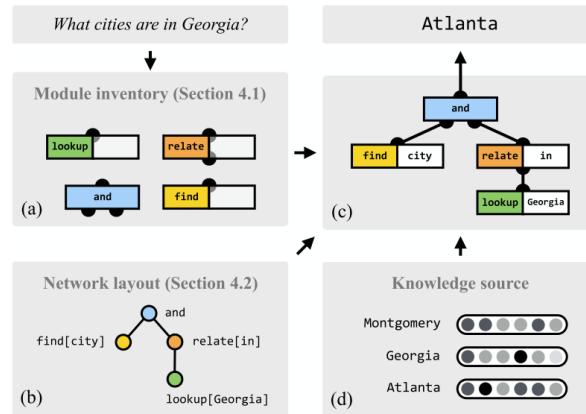
## Learning to Compose Neural Networks for QA

Table of Contents Local

Written by Brandon McKinzie

Andreas et al., “Learning to Compose Neural Networks for Question Answering” UC Berkeley (June 2016).

**TL;DR.** Improve initial NMN work (previous note) by (1) learning network predictor ( $P(w)$  in previous paper) instead of manually specifying it, and (2) extending visual primitives from previous work to reason over structured world representations.



**Model** (4). Training data consists of (world, question, answer) triplets  $(w, x, y)$ . The model is built around two distributions:

- **layout model**  $p(z | x; \theta_\ell)$  which predicts a layout  $z$  for sentence  $x$ .
- **execution model**  $p_z(y | w; \theta_e)$  which applies the network specified by  $z$  to the world representation  $w$ .

**Evaluating Modules** (4.1). The execution model is defined as

$$p_z(y | w) = (\llbracket z \rrbracket_w)_y \quad (739)$$

where  $\llbracket z \rrbracket_w$  denotes the output of network with layout  $z$  on input world  $w$ . The defining equations for all modules is as follows ( $\sigma \equiv \text{ReLU}$ ,  $sm \equiv \text{softmax}$ ):

$$\llbracket \text{lookup}[i] \rrbracket = e_{f(i)} \quad (740)$$

$$\llbracket \text{find}[i] \rrbracket = sm(a \odot \sigma(Bv^i \oplus CW \oplus d)) \quad (741)$$

$$\llbracket \text{relate}[i](h) \rrbracket = sm(a \odot \sigma(Bv^i \oplus CW \oplus D\bar{w}(h) \oplus e)) \quad (742)$$

$$\llbracket \text{and}(h^1, h^2, \dots) \rrbracket = h^1 \odot h^2 \odot \dots \quad (743)$$

$$\llbracket \text{describe}[i](h) \rrbracket = sm(A\sigma(B\bar{w}(h)) + v^i) \quad (744)$$

$$\llbracket \text{exists}(h) \rrbracket = sm\left(\left(\max_k h_k\right)a + b\right) \quad (745)$$

$$\bar{w}(h) \triangleq \sum_k h_k w^{(k)}$$

To train, maximize

$$\sum_{(w,y,z)} \log p_z(y | w; \theta_e) \quad (746)$$

**Assembling Networks** (4.2). **TODO:** finish note

## End-to-End Module Networks for VQA

Table of Contents Local

Written by Brandon McKinzie

R. Hu, J. Andreas, et al., “Learning to Reason: End-to-End Module Networks for Visual Question Answering”  
*UC Berkeley, FAIR, BU* (Sep 2017).

**End-to-End Module Networks** (3). High-level sequence of operations, given some input question and image:

1. Layout policy predicts a coarse functional expression that describes the structure of the computation.
2. Some subset of function applications within the expression receive parameter vectors predicted from the text.
3. Network is assembled with the modules according to layout expression to output an answer.

**Attentional Neural Modules** (3.1). A neural module  $m$  is a parameterized function  $y = f_m(a_1, a_2, \dots; x_{vis}, x_{txt}, \theta_m)$ , where the  $a_i$  are image attention maps and the output  $y$  is either an image attention map or a probability distribution over answers. The full set of modules used by the authors, along with their inputs/outputs, is tabulated below.

Module name	Att-inputs	Features	Output	Implementation details
find	(none)	$x_{vis}, x_{txt}$	att	$a_{out} = \text{conv}_2(\text{conv}_1(x_{vis}) \odot Wx_{txt})$
relocate	$a$	$x_{vis}, x_{txt}$	att	$a_{out} = \text{conv}_2(\text{conv}_1(x_{vis}) \odot W_1\text{sum}(a \odot x_{vis}) \odot W_2x_{txt})$
and	$a_1, a_2$	(none)	att	$a_{out} = \min(a_1, a_2)$
or	$a_1, a_2$	(none)	att	$a_{out} = \max(a_1, a_2)$
filter	$a$	$x_{vis}, x_{txt}$	att	$a_{out} = \text{and}(a, \text{find}[x_{vis}, x_{txt}]())$ , i.e. reusing find and and
[exist, count]	$a$	(none)	ans	$y = W^T \text{vec}(a)$
describe	$a$	$x_{vis}, x_{txt}$	ans	$y = W_1^T (W_2\text{sum}(a \odot x_{vis}) \odot W_3x_{txt})$
[eq.count, more, less]	$a_1, a_2$	(none)	ans	$y = W_1^T \text{vec}(a_1) + W_2^T \text{vec}(a_2)$
compare	$a_1, a_2$	$x_{vis}, x_{txt}$	ans	$y = W_1^T (W_2\text{sum}(a_1 \odot x_{vis}) \odot W_3\text{sum}(a_2 \odot x_{vis}) \odot W_4x_{txt})$

Note that, whereas the original NMN paper (see previous note) instantiated module types based on words (e.g. `describe[shape]` vs `describe[where]`) and gave different instantiations different parameters, this paper has a single module for each module type (no “instances” anymore). To distinguish between cases where e.g. `describe` should describe a shape vs describing a location, the module incorporates a text feature  $x_{txt}^{(m)}$  computed separately/identically for each module  $m$ :

$$x_{txt}^{(m)} = \sum_{i=1}^T \alpha_i^{(m)} w_i \quad (747)$$

**Layout Policy with Seq2Seq RNN** (3.2). **TODO** finish note

## Fast Multi-language LSTM-based Online Handwriting Recognition

Table of Contents Local

Written by Brandon McKinzie

Carbune et al., “Fast Multi-language LSTM-based Online Handwriting Recognition” *Google AI Perception* (Feb 2019).

**Introduction** (1). Task: given input strokes, i.e. sequences of points  $(x, y, t)$ , output it in the form of text.

**Model Architecture** (2). The high-level sequences of operations is:

1. Input time series  $(v_1, \dots, v_T)$  encoding user input. The authors experiment with two representations:
  - (a) *Raw touch points*: sequence of 5-dimensional points  $(x_i, y_i, t_i, p_i, n_i)$ , where  $t_i$  is seconds since first touch point in current observation,  $p_i$  is binary-valued equal to 0 if pen-up, else 1 if pen-down, and  $n_i$  is binary on start-of-new-stroke (1 if True).
  - (b) *Bézier curves*: TL;DR is that they model x, y, t each as a cubic polynomial over a new variable  $s \in [0, 1]$ . Ultimately this means solving for some coefficients  $\Omega$  of a linear system of equations:  $V^T Z = V^T V \Omega$ .
2. Several BiLSTM layers for contextual character embedding.
3. Softmax layer providing character probabilities at each time step.
4. CTC decoding with beam search. They also incorporate **feature functions** into the output logits to help with decoding. They use the following 3 feature functions:
  - (a) Character language models. A 7-gram LM over Unicode codepoints using Stupid back-off.
  - (b) Word language models. 3-grams pruned to between 1.25M and 1.5M entries.
  - (c) Character classes. Scoring heuristic which boosts the score of characters from the LM’s alphabet.

**Training** (3). Training happens in two stages, each on a different dataset:

1. Training neural network model with CTC loss on large dataset.
2. Decoder tuning using **Bayesian optimization** through **Gaussian Processes** in Vizier<sup>165</sup>.

---

<sup>165</sup>Vizier is a program made by Google for black-box tuning

## Multi-Language Online Handwriting Recognition

Table of Contents Local

Written by Brandon McKinzie

Keyser et al., “Multi-Language Online Handwriting Recognition” *Google* (June 2017).

**System Architecture** (3). Segment-and-decode approach consisting of the following steps:

- Preprocessing (4).
  1. Resampling.
  2. Slope correction.
- **Segmentation**<sup>166</sup> and search lattice creation (5).
  1. Segmentation goal: obtain high *recall* of all actual character boundaries. Accomplished via a heuristic which creates a set of potential cut points and then a neural net which assigns a score to each.
  2. Segmentation lattice: a graph  $(V, E)$  of ink segments. Each segment is identified by a unique integer index.
    - Nodes (in  $V$ ) define the path of ink segments up to that point (e.g.  $\{1, 0, 2\}$  (i.e. a character hypothesis)
    - Edges (in  $E$ ) from a given node  $v$  indicate the ink segments which are grouped in a character hypothesis. For example, if  $v=\{i, j\}$  has some edge  $k$ , then that edge will have node  $\{i, j, k\}$  on the other end, and  $\{i, j, k\}$  is a valid character hypothesis.

It appears that each node (assign from the empty start node) is passed to the next stage as a character hypothesis to be scored/classified.
- Generation & scoring of **character hypotheses**<sup>167</sup> (5.3). Goal: determine the characters most likely to have been written.
  1. Feature extraction: they make a fixed-length dense feature vector containing *point-wise* and *character-global* features.
  2. Classification: single hidden layer NN with tanh activation followed by softmax.
  3. Create a labeled lattice which will later be decoded to find the final recognition result.
- Best path search in the resulting lattice using additional knowledge sources (6).

---

<sup>166</sup> **Segmentation/cut point**: a point at which another character may start. **Segment**: the (partial) strokes between 2 consecutive segmentation points.

<sup>167</sup> **Character hypothesis**: a set of one or more segments (not necessarily consecutive).

**Language Models** (6.1). They utilize two types of language models:

- Stupid-backoff entropy-pruned 9-gram character LM. This is their “main” LM. Depending on the language, they use about 10M to 100M n-grams.
- Word-based probabilistic finite automaton. Creating using 100K most frequent words of a language.

**Search** (6.2). Goal: obtain a recognition result by finding the best path from the source node (no ink recognized) to the target node (all ink recognized). Algorithm: ink-aligned beam search that starts at the start node and proceeds through the lattice in topological order.

## Modular Generative Adversarial Networks

Table of Contents Local

Written by Brandon McKinzie

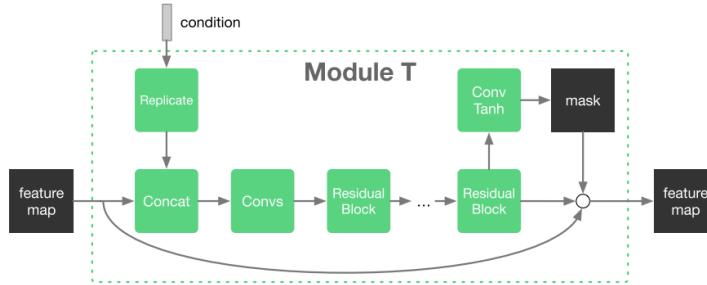
Zhao et al., “Modular Generative Adversarial Networks” *UBC, Tencent AI Lab* (April 2018).

**TL;DR.** Task(s): multi-domain image **generation** and image-to-image **translation**.

**Network Construction** (3.2). Let  $x$  and  $y$  denote the input and target image, respectively, wherever applicable. Let  $\mathbf{A} = \{A_1, A_2, \dots, A_n\}$  denote an **attribute set**. Four types of modules are used:

1. Initial module is task-dependent (below). Output is feature map in  $\mathbb{R}^{C \times H \times W}$ .
  - [translation] **encoder**  $\mathbf{E}$ :  $x \mapsto \mathbf{E}(x)$
  - [generation] **generator**  $\mathbf{G}$ :  $(z, a_0) \mapsto \mathbf{G}(z, a_0)$  where  $z$  is random noise and  $a_0$  is a condition vector representing auxiliary information.
2. **transformer(s)**  $\mathbf{T}_i$ :  $E(x) \mapsto \mathbf{T}_i(E(x), a_i)$ . Modifies repr of attrib  $a_i$  in the FM.
3. **reconstructor**  $\mathbf{R}$ :  $(\mathbf{T}_i, \mathbf{T}_j, \dots) \mapsto y$ . Reconstructs image from an intermediate FM.
4. **discriminator**  $\mathbf{D}_i$ :  $\mathbf{R} \mapsto \{0, 1\} \times \text{Val}(a_i)$ . Predicts probability that  $R$  came from  $p_{true}$ , and the [transformed] value of  $a_i$ .

The authors emphasize that the transformer module is their core module. It’s architecture is illustrated below.



**Loss Function (3.4).**

$$\mathcal{L}_D(\mathbf{D}) = -\sum_{i=1}^n \mathcal{L}_{adv_i} + \lambda_{cls} \sum_{i=1}^n \mathcal{L}_{cls_i}^r \quad (748)$$

$$\mathcal{L}_G(\mathbf{E}, \mathbf{T}, \mathbf{R}) = \sum_{i=1}^n \mathcal{L}_{adv_i} + \lambda_{cls} \sum_{i=1}^n \mathcal{L}_{cls_i}^f + \lambda_{cyc} \left( \mathcal{L}_{cyc}^{\mathbf{ER}} + \sum_{i=1}^n \mathcal{L}_{cyc}^{\mathbf{T}_i} \right) \quad (749)$$

$$\mathcal{L}_{adv_i}(\mathbf{E}, \mathbf{T}_i, \mathbf{R}, \mathbf{D}_i) = \mathbb{E}_{y \sim p_{data}(y)} [\log \mathbf{D}_i(y)] + \mathbb{E}_{x \sim p_{data}(x)} [\log (1 - \mathbf{D}_i(\mathbf{R}(\mathbf{T}_i(\mathbf{E}(x)))))] \quad (750)$$

$$\mathcal{L}_{cls_i}^r = -\mathbb{E}_{x, c_i} [\log \mathbf{D}_{cls_i}(c_i | x)] \quad (751)$$

$$\mathcal{L}_{cls_i}^f = -\mathbb{E}_{x, c_i} [\log \mathbf{D}_{cls_i}(c_i | \mathbf{R}(\mathbf{E}(\mathbf{T}_i(x))))] \quad (752)$$

$$\mathcal{L}_{cyc}^{\mathbf{ER}} = \mathbb{E}_x [||\mathbf{R}(\mathbf{E}(x)) - x||_1] \quad (753)$$

$$\mathcal{L}_{cyc}^{\mathbf{T}_i} = \mathbb{E}_x [||\mathbf{T}_i(\mathbf{E}(x)) - \mathbf{E}(\mathbf{R}(\mathbf{T}_i(\mathbf{E}(x))))||_1] \quad (754)$$

where  $n$  is the total number of controllable attributes.

## Transfer Learning from Speaker Verification to TTS

Table of Contents Local

*Written by Brandon McKinzie*

Jia et al., “Transfer Learning from Speaker Verification to Multispeaker Text-To-Speech Synthesis” *Google* (Jan 2019).

**TL;DR:** TTS that’s able to generate speech in the voice of different speakers, including those unseen during training.

**Multispeaker Speech Synthesis Model** (2). System is composed of three independently trained NNs:

1. **Speaker Encoder**. Computes a fixed-dimensional vector from a speech signal.
2. **Synthesizer**. Predicts a **mel spectrogram** from a sequence of **grapheme** or **phoneme** inputs, conditioned on the speaker vector. Extension of **Tacotron 2** to support multiple speakers.
3. **Vocoder**. Autoregressive WaveNet, which converts the spectrogram into time domain waveforms.

# NLP WITH DEEP LEARNING

## CONTENTS

5.1	Word Vector Representations (Lec 2) . . . . .	252
5.2	GloVe (Lec 3) . . . . .	255

## Word Vector Representations (Lec 2)

Table of Contents Local

Written by Brandon McKinzie

**Meaning of a word.** Common answer is to use a *taxonomy* like WordNet that has hypernyms (is-a) relationships. Problems with this discrete representation: misses nuances, e.g. the words in a set of synonyms can actually have quite different meanings/connotations. Furthermore, viewing words as atomic symbols is a bit like using one-hot vectors of words in a vocabulary space (inefficient).

**Distributed representations.** Want a way to encode word vectors such that two similar words have a similar structure/representation. The **distributional similarity-based**<sup>168</sup> approach represents words by means of its *neighbors* in the sentences in which it appears. You end up with a dense “vector for each word type, chosen so that it is good at predicting other words appearing in its context.”

**Skip-gram prediction.** Given a word  $w_t$  at position  $t$  in a sentence, learn to predict [probability of] some number of surrounding words, given  $w_t$ . Standard minimization with negative log-likelihood:

$$J(\theta) = -\frac{1}{T} \sum_{t=1}^T \sum_{-m \leq j \leq m} \log \Pr(w_{t+j} | w_t) \quad (755)$$

$$\Pr(o | c) = \frac{e^{\mathbf{u}_o^T \mathbf{v}_c}}{\sum_{w=1}^{\text{vocab size}} e^{\mathbf{u}_w^T \mathbf{v}_c}} \quad (756)$$

I cannot believe this actually works.

where

- The params  $\theta$  are the vector representation of the words (they are the *only* learnable parameters here).
- $m$  is our radius/window size.
- $o$  and  $c$  are indices into our vocabulary (somewhat inconsistent notation).
- Yes, they are using different vector representations for  $\mathbf{u}$  (context words) and  $\mathbf{v}$  (center words). I’m assuming one reason this is done is because it makes the model architecture simpler/easier to build.

Some subtleties:

---

<sup>168</sup>Note that this is distinct from the way “distributed” is meant in “distributed representation.” In contrast, distributional similarity-based representations refers to the notion that you can describe the meaning of words by understanding the context in which they appear.

- Looks like e.g.  $Pr(w_{t+j} | w_t)$  doesn't really care what the value of  $j$  is, it is just modeling the probability that it is *somewhere* in the context window. The  $w_t$  are one-hot vectors into the vocabulary. Standard tricks for simplifying the cross-entropy loss apply.
- Equation 756 should be interpreted as the probability that the  $o^{th}$  word in our vocabulary occurs in the context window of the  $c^{th}$  word in our vocabulary.
- The model architecture is *identical* to an autoencoder. However, the (big) difference is the training procedure and interpretation of the model parameters going “in” versus the parameters going “out”.

**Sentence Embeddings.** It turns out that simply taking a weighted average of word vectors and doing some PCA/SVD is a competitive way of getting unsupervised word embeddings. Apparently it *beats* supervised learning with LSTMs (?!). The authors claim the theoretical explanation for this method lies in a latent variable generative model for sentences (of course). Approach:

Discussion based on paper by Arora et al., (2017).

1. Compute the weighted average of the word vectors in the sentence:

$$\frac{1}{N} \sum_i^N \frac{a}{a + p(\mathbf{w}_i)} \mathbf{w}_i \quad (757)$$

The authors call their weighted average the **Smooth Inverse Frequency (SIF)**.

where  $\mathbf{w}_i$  is the word vector for the  $i$ th word in the sentence,  $a$  is a parameter, and  $p(\mathbf{w}_i)$  is the (estimated) word frequency [over the entire corpus].

2. Remove the projections of the average vectors on their first principal component (“common component removal”) (y tho?).

1

---

**Algorithm 1** Sentence Embedding

---

**Input:** Word embeddings  $\{v_w : w \in \mathcal{V}\}$ , a set of sentences  $\mathcal{S}$ , parameter  $a$  and estimated probabilities  $\{p(w) : w \in \mathcal{V}\}$  of the words.  
**Output:** Sentence embeddings  $\{v_s : s \in \mathcal{S}\}$

```

1: for all sentence  $s$  in  $\mathcal{S}$  do
2:    $v_s \leftarrow \frac{1}{|s|} \sum_{w \in s} \frac{a}{a + p(w)} v_w$ 
3: end for
4: Compute the first principal component  $u$  of  $\{v_s : s \in \mathcal{S}\}$ 
5: for all sentence  $s$  in  $\mathcal{S}$  do
6:    $v_s \leftarrow v_s - uu^\top v_s$ 
7: end for

```

---

## **Further Reading.**

- Learning representations by back-propagating errors (Rumelhard et al., 1986)
- A Neural Probabilistic Language Model (Bengio et al., 2003)
- NLP (almost) from Scratch (Collobert & Watson, 2008)
- Word2Vec (Miklov et al. 2013)

## GloVe (Lec 3)

Table of Contents Local

Written by Brandon McKinzie

**Skip-gram and negative sampling.** Main idea:

- Split the loss function from last lecture into two (additive) terms corresponding to the numerator and denominator respectively (you've done this a trillion times).
- The second term is an expectation over all the words in your vocab space. That is huge, so instead we only use a subsample of size  $k$  (the negative samples).
- Interpretation: First term is maximizing  $\Pr(o | c)$ , the probability that the true outside word (given by index o) occurs given context (index) c. Second term is minimizing the probability of random words (the negative samples) occurring around the center (context) word given by c.

To sample the negative samples, draw from  $P(w) = U(w)^{3/4}/Z$ , where  $U$  is the **unigram distribution**.

**GloVe** (Global Vectors). Given some co-occurrence matrix we computed with previous methods, we can use the following GloVe loss function over all pairs of co-occurring words in our matrix:

$$J(\theta) = \sum_{i,j=1}^W f(P_{ij})(u_i^T v_j - \log P_{ij})^2 \quad (758)$$

where  $P_{ij}$  is computed simply from the counts of words i and j co-occurring (empirical probability) and  $f$  is some squashing function that really isn't discussed in this lecture (**TODO**).

### Evaluating word vectors.

- **Word Vector Analogies:** Basically, determining if we can do standard analogy fill-in-the-blank problems: “man [a] is to woman [b] as king [c] is to <blank>” (if you answered “queen”, you'd make a good AI). We can determine this using a standard cosine distance measure:

$$d = \arg \max_i \frac{(x_b - x_a + x_c)^T x_i}{\|x_b - x_a + x_c\|} \quad (759)$$

Woah that is pretty neat. The solution is  $x_i = \text{queen}$ .  $x_b - x_a$  is the vector pointing from man to woman, which encodes the type of similarity we are looking for with the other pair. Therefore, we take the vector to “king” and *add* the aforementioned difference vector – the resultant vector should point to “queen”. Neat!

**Derivation.** Based on the descriptions in the original paper<sup>169</sup> We want to develop a model for learning word vectors.

1. The authors argue that “the appropriate starting point for word vector learning should be with ratios of co-occurrence probabilities rather than the probabilities themselves.” The most general such model takes the form,

$$F(w_i, w_j, \tilde{w}_k) = \frac{\Pr[\tilde{w}_k | w_i]}{\Pr[\tilde{w}_k | w_j]} \equiv \frac{P_{ik}}{P_{jk}} \quad \text{where all } w \in \mathbb{R}^d \quad (760)$$

and the tilde in  $\tilde{w}_k$  denotes that  $\tilde{w}_k$  is a *context* word vector, which are given a distinct space from the word vectors  $w_i$  and  $w_j$  being compared. We compute all  $P_{ik}$  via frequency counts over the corpus.

2. Now that we’ve specified the inputs and ratio of interest, we can start specifying some desirable constraints on the function  $F$  that we’re trying to find. The authors speculate that, since vector spaces are linear structures, we should have  $F$  encode the information of the ratio in the vector space via vector differences:

$$F(w_i - w_j, \tilde{w}_k) = \frac{P_{ik}}{P_{jk}} \quad (761)$$

which basically says “our representation of the word vectors should be s.t. the **relative** probability of some word  $\tilde{w}_k$  occurring in the context of a word  $w_i$  compared to it occurring in the context of a different word  $w_j$  can be captured by  $w_i - w_j$  and  $\tilde{w}_k$  alone.”

3. Next we notice that  $F$  maps arguments in  $\mathbb{R}^d$  to a scalar in  $\mathbb{R}$ . The most straightforward way of doing so while maintaining the linear structure we are trying to capture is via a dot product:

$$F((w_i - w_j)^T \tilde{w}_k) = \frac{P_{ik}}{P_{jk}} \quad (762)$$

Note that now  $F : \mathbb{R} \mapsto \mathbb{R}$ .

4. We want our model to be invariant under the exchanges  $w \leftrightarrow \tilde{w}$  and  $X \leftrightarrow X^T$ . We can restore this symmetry by first requiring that  $F$  be a **homomorphism**<sup>170</sup> between the groups  $(\mathbb{R}, +)$  and  $(\mathbb{R}_{>0}, \times)$  (in our case, negation and division would be better symbols, but it’s equivalent). This requires that the following relation hold<sup>171</sup>

$$F(w_i^T \tilde{w}_k - w_j^T \tilde{w}_k) = \frac{F(w_i^T \tilde{w}_k)}{F(w_j^T \tilde{w}_k)} \quad (763)$$

where I’ve grouped terms on the LHS to emphasize how this is the definition of homomorphism. The solution for this equation is that  $F(\cdot) \equiv \exp(\cdot)$ . Combining this realization with equations 762 and 763 yields,

$$F(w_i^T \tilde{w}_k) = e^{w_i^T \tilde{w}_k} = P_{ik} \quad (764)$$

$$\Rightarrow w_i^T \tilde{w}_k = \log(P_{ik}) = \log(X_{ik}) - \log(X_i) \quad (765)$$

---

<sup>169</sup>Pennington et al., “GloVe: Global Vectors for Word Representation.”

<sup>170</sup>In more detail,  $F : (\mathbb{R}, +) \mapsto (\mathbb{R}_{>0}, \times)$ , which reads “the function  $F$  maps elements in  $\mathbb{R}$  and any summation of elements in  $\mathbb{R}$  to elements in  $\mathbb{R}$  that are greater than zero or any product of positive elements in  $\mathbb{R}$ .”

<sup>171</sup>Note that we do not need to know anything about the RHS of the equations above to state this relation. We write it by definition of homomorphism.

where  $X_{ik}$  is the number of times word  $k$  appears in the context of word  $i$ , and  $X_i = \sum_k X_{ik}$  is the number of times any word appears in the context of  $i$ .

5. Restore symmetry under the exchanges  $w \leftrightarrow \tilde{w}$  and  $X \leftrightarrow X^T$ . We absorb  $X_i$  into a bias  $b_i$  for  $w_i$  since it is independent of  $k$ .

$$w_i^T \tilde{w}_k + b_i + \tilde{b}_k = \log(X_{ik}) \quad (766)$$

6. A main drawback to this model is that it weighs all co-occurrences equally, even those that happen rarely or never. The authors propose a new weighted least squares regression model, introducing a weighting function  $f(X_{ij})$  into the cost function of our model:

$$J = \sum_{i,j=1}^V f(X_{ij})(w_i^T \tilde{w}_j + b_i + \tilde{b}_j - \log X_{ij})^2 \quad (767)$$

# SPEECH AND LANGUAGE PROCESSING

## CONTENTS

6.1	Introduction (Ch. 1 2nd Ed.) . . . . .	259
6.2	Morphology (Ch. 3 2nd Ed.) . . . . .	260
6.3	N-Grams (Ch. 6 2nd Ed.) . . . . .	261
6.4	Naive Bayes and Sentiment (Ch. 6 3rd Ed.) . . . . .	263
6.5	Hidden Markov Models (Ch. 9 3rd Ed.) . . . . .	265
6.6	POS Tagging (Ch. 10 3rd Ed.) . . . . .	268
6.7	Formal Grammars (Ch. 11 3rd Ed.) . . . . .	271
6.8	Vector Semantics (Ch. 15) . . . . .	272
6.9	Semantics with Dense Vectors (Ch. 16) . . . . .	275
6.10	Information Extraction (Ch. 21 3rd Ed) . . . . .	278

## Introduction (Ch. 1 2nd Ed.)

Table of Contents Local

*Written by Brandon McKinzie*

**Overview.** Going to rapidly jot down what seems most important from this chapter.

- **Morphology:** captures information about the shape and behavior of words in context (Ch. 2/3).
- **Syntax:** knowledge needed to order and group words together.
- **Lexical semantics:** knowledge of the meanings of individual words.
- **Compositional semantics:** knowledge of how these components (words) combine to form larger meanings.
- **Pragmatics:** the appropriate use of polite and indirect language.
- The knowledge of language needed to engage in complex language behavior can be separated into the following 6 distinct categories:
  1. Phonetics and Phonology – The study of linguistic sounds.
  2. Morphology – The study of the meaningful components of words.
  3. Syntax – The study of the structural relationships between words.
  4. Semantics – The study of meaning.
  5. Pragmatics – The study of how language is used to accomplish goals.
  6. Discourse – The study of linguistic units larger than a single utterance.
- Methods for **resolving ambiguity:** pos-tagging, word sense disambiguation, probabilistic parsing, and speech act interpretation.
- **Models and Algorithms.** Among the most important are **state space search** and **dynamic programming** algorithms.

## Morphology (Ch. 3 2nd Ed.)

Table of Contents Local

*Written by Brandon McKinzie*

**English Morphology.** Morphology is the study of the way words are built up from smaller meaning-bearing units, **morphemes**. A morpheme is often defined as the minimal meaning-bearing unit in a language<sup>172</sup>. The two classes of morphemes are **stems** (the “main” morpheme of the word) and **affixes** (the “additional” meanings of various kinds).

Affixes are further divided into prefixes (precede stem), suffixes (follow stem), circumfixes (do both), and infixes (inside the stem).

Two classes of ways to form words from morphemes: **inflection** and **derivation**. Inflection is the combination of a word stem with a grammatical morpheme, usually resulting in a word of the same class as the original stem, and usually filling some syntactic function like agreement. Derivation is the combination of a word stem with a grammatical morpheme, usually resulting in a word of a different class, often with a meaning hard to predict exactly.

---

<sup>172</sup>Examples: “fox” is its own morpheme, while “cats” consists of the morpheme “cat” and the morpheme “-s”.

## N-Grams (Ch. 6 2nd Ed.)

Table of Contents Local

Written by Brandon McKinzie

**Counting Words.** Most  $N$ -gram based systems deal with the *wordform*, meaning they treat words like “cats” and “cat” as distinct. However, we may want to treat the two as instances of a single abstract word, or **lemma**: a set of lexical forms having the same stem, the same major part of speech, and the same word-sense.

**Simple (Unsmoothed) N-Grams.** An  $N$ -gram is a  $N$ -th order Markov model (because it looks  $N-1$  steps in the past). Notation: the authors use the convention that  $w_1^n \triangleq w_1, w_2, \dots, w_n$  to denote a sequence of  $n$  words. Given this, we can write the general equation for the  $N$ -gram approximation for the probability of the  $n$ th word ( $n > N$ ) in a sentence:

$$\Pr[w_n | w_1^{n-1}] \approx \Pr[w_n | w_{n-N+1}^{n-1}] \quad (768)$$

for  $N \geq 2$ . We can compute these probabilities by simply counting:

$$\Pr[w_n | w_1^{n-1}] = \frac{C(w_{n-N+1}^{n-1} w_n)}{C(w_{n-N+1}^{n-1})} \quad (769)$$

where  $C(\cdot)$  is the number of times the sequence, denoted as  $\cdot$ , occurred in the corpus.

**Entropy.** Denote the random variable of interest as  $\mathbf{x}$  with probability function  $p(\mathbf{x})$ . The entropy of this random variable is:

$$H(\mathbf{x}) = - \sum_x p(\mathbf{x} = x) \log_2 p(\mathbf{x} = x) \quad (770)$$

which should be thought of as a lower bound on the number of bits it would take to encode a certain decision or piece of information in the optimal coding scheme. The value  $2^H$  is the **perplexity**, which can be interpreted as the weighted average number of choices a random variable has to make.

**Cross Entropy for Comparing Models.** Useful when we don't know the actual probability distribution  $p$  that generated some data. Assume we have some model  $m$  that's an approximation of  $p$ . The cross-entropy of  $m$  on  $p$  is defined by:

$$H(p, m) = \lim_{n \rightarrow \infty} \frac{1}{n} \sum_{W \in L} p(w_1, \dots, w_n) \log m(w_1, \dots, w_n) \quad (771)$$

That is we draw sequences according to the probability distribution  $p$ , but sum the log of their probability according to  $m$ .

## Naive Bayes and Sentiment (Ch. 6 3rd Ed.)

Table of Contents Local

Written by Brandon McKinzie

[3rd Ed.] [Quick Review]

**Overview.** This chapter is concerned with **text categorization**, the task of classifying an entire text by assigning it a label drawn from some set of labels. *Generative classifiers* like naive Bayes build a model of each class. Given an observation, they return the class most likely to have generated the observation. *Discriminative classifiers* like logistic regression instead learn what features from the input are most useful to discriminate between the different possible classes. Notation: we will assume we have a training set of  $N$  documents, each hand-labeled with some class:  $\{(d_1, c_1), \dots, (d_N, c_N)\}$ .

Discriminative systems are often more accurate and hence more commonly used.

**Naive Bayes.** A multinomial<sup>173</sup> classifier with a naive assumption about how the features interact. We model a text document as a **bag of words**, meaning we store (1) the words that occurred and (2) their frequencies. It's a probabilistic classifier, meaning it estimates the label/class of a document  $d$  as

$$\hat{c} = \arg \max_{c \in C} \Pr [c | d] \quad (773)$$

$$= \arg \max_{c \in C} \frac{\Pr [d | c] \Pr [c]}{\Pr [d]} \quad (774)$$

$$= \arg \max_{c \in C} \Pr [d | c] \Pr [c] \quad (775)$$

Computing the class-conditional distribution (the likelihood)  $\Pr [d | c]$  over all possible  $d \in D$  is typically intractable; we must introduce some simplifying assumptions and use an approximation of it. Our assumptions in this section will be:

- **Bag-of-Words:** Assume that word position within a document is irrelevant. (Counts still matter)

---

<sup>173</sup>

Rapid review: **multinomial distribution.** Let  $\mathbf{x} = (x_1, \dots, x_k)$  denote the result of an experiment with  $n$  independent trials ( $n = \sum_i^k x_i$ ) and  $k$  possible outcomes for any given trial. i.e.  $x_i$  is the number of trials that had outcome  $i$  ( $1 \leq i \leq k$ ). The pmf of this multinomial distribution, over all possible  $\mathbf{x}$  constrained by  $n = \sum_i^k x_i$ , is:

$$\Pr [\mathbf{x} = (x_1, \dots, x_k); n] = \frac{n!}{x_1! \cdots x_k!} p_1^{x_1} \times \cdots \times p_k^{x_k} \quad (772)$$

where  $p_i$  is the probability of outcome  $i$  for any single trial.

- **Naive Bayes Assumption:** First, recall that  $d$  is typically modeled as a (random) vector consisting of features  $f_1, \dots, f_n$ , each of which has an associated probability distribution  $\Pr [f_i | c]$ . The NB assumption is that the features are mutually independent given the class  $c$ :

$$\Pr [f_1, \dots, f_n | c] = \Pr [f_1 | c] \cdots \Pr [f_n | c] \quad (776)$$

$$c_{NB} = \arg \max_{c \in C} \Pr [c] \prod_{f \in F} \Pr [f | c] \quad (777)$$

where 777 is the final equation for the class chosen by the naive Bayes classifier.

In text classification we typically use the word at position  $i$  in the document as  $f_i$ , and move to log space to avoid underflow/increase speed:

$$c_{NB} = \arg \max_{c \in C} \log \Pr [c] + \sum_i^{\text{len}(d)} \log \Pr [w_i | c] \quad (778)$$

Classifiers that use a linear combination of the inputs to make a classification decision (e.g. NB, logistic regression) are called linear classifiers.

**Training the NB Classifier.** No real "training" in my opinion, just simple counting from the data:

$$\hat{P}[c] = \frac{N_c}{N_{\text{docs}}} \quad (779)$$

$$\hat{P}[w_i | c] = \frac{\text{count}(w_i, c) + 1}{(\sum_{w \in V} \text{count}(w, c)) + |V|} \quad (780)$$

The Laplace smoothing is added to avoid the occurrence of zero-probabilities in equation 777.

### Optimizing for Sentiment Analysis.

- It often improves performance [for sentiment] to clip word counts in each document to 1 (“binary NB”).
- deal with *negations* in some way.
- Use sentiment lexicons, lists of words that are pre-annotated with positive or negative sentiment.

## Hidden Markov Models (Ch. 9 3rd Ed.)

Table of Contents Local

Written by Brandon McKinzie

**Overview.** Here we will first go over the math behind HMMs: the **Viterbi**, **Forward**, and **Baum-Welch** (EM) algorithms for unsupervised or semi-supervised learning. Recall that a HMM is defined by specifying the set of  $N$  states  $Q$ , transition matrix  $A$ , sequence of  $T$  observations  $O$ , sequence of observation likelihoods  $B$ , and the initial/final states. They can be characterized by three fundamental problems:

1. **Likelihood.** Given an HMM  $\lambda = (A, B)$  and observation sequence, compute the likelihood (prob. of the observations given the model). (**Forward**)
2. **Decoding.** Given an HMM  $\lambda = (A, B)$  and observation sequence, discover the best hidden state sequence. (**Viterbi**)
3. **Learning.** Given an observation sequence and the set of states in the HMM, learn the HMM parameters  $A$  and  $B$ . (**Baum-Welch/Forward-Backward/EM**)

**The Forward Algorithm.** For likelihood computation. We want to compute the probability of some sequence of observations  $O$ , without knowing the sequence of hidden states (that emitted the observations)  $Q$ . In general, this can be expressed by summing over all possible hidden state sequences:

$$\Pr [O] = \sum_Q \Pr [Q] \Pr [O | Q] \quad (781)$$

However, for  $N$  hidden states and  $T$  observations, this summation involves  $N^T$  terms, which becomes intractable rather quickly. Instead, we can use the  $\mathcal{O}(N^2T)$  **Forward Algorithm**. The forward algorithm can be defined by initialization, recursion definition, and termination, shown respectively as follows:

$$\alpha_1(j) = a_{0j} b_j(o_1) \quad 1 \leq j \leq N \quad (782)$$

$$\alpha_t(j) = \sum_{i=1}^N \alpha_{t-1}(i) a_{ij} b_j(o_t) \quad 1 \leq j \leq N, 1 \leq t \leq T \quad (783)$$

$$\Pr [O | \lambda] = \alpha_T(q_F) = \sum_{i=1}^N \alpha_T(i) a_{iF} \quad (784)$$

**Viterbi Algorithm.** For decoding. Want the most likely hidden state sequence given observations. Let  $v_t(j)$  represent the probability that we are in state  $j$  after  $t$  observations and passing through the most probable state sequence  $q_0, q_1, \dots, q_{t-1}$ . Similar to the forward algorithm, we show the defining equations for the Viterbi algorithm below:

$$v_1(j) = a_{0j} b_j(o_1) \quad 1 \leq j \leq N \quad (785)$$

$$v_t(j) = \max_{i=1}^N v_{t-1}(i) a_{ij} b_j(o_t) \quad (786)$$

$$P^* = v_T(q_F) = \max_{i=1}^N v_T(i) a_{iF} \quad (787)$$

N.B.: At each step, the best path up to that point can be found by taking the *argmax* instead of max.

**Baum-Welch Algorithm.** AKA forward-backward algorithm, a special case of the EM algorithm. Given an observation sequence  $O$  and the set of best possible states in the HMM, learn the HMM parameters  $A$  and  $B$ . First, we must define some new notation. The **backward probability**  $\beta$  is defined as:

$$\beta_t(i) \triangleq \Pr [o_{t+1}, \dots, o_T \mid q_t = i, \lambda] \quad (788) \quad \text{Remember } \lambda \equiv (A, B)$$

As usual, we can compute its values inductively as follows:

$$\beta_T(i) = a_{iF} \quad 1 \leq i \leq N \quad (789)$$

$$\beta_t(i) = \sum_{j=1}^N a_{ij} b_j(o_{t+1}) \beta_{t+1}(j) \quad 1 \leq i \leq N, 1 \leq t \leq T \quad (790)$$

$$\Pr [O \mid \lambda] = \alpha_T(q_F) = \beta_1(q_0) \quad (791)$$

$$= \sum_{j=1}^N a_{0j} b_j(o_1) \beta_1(j) \quad (792)$$

We can use the forward and backward probabilities  $\alpha$  and  $\beta$  to compute the transition probabilities  $a_{ij}$  and observation probabilities  $b_i(o_t)$  from an observation sequence. The derivation steps are as follows:

1. **Estimating  $\hat{a}_{ij}$ .** Begin by defining quantities that will prove useful:

$$\xi_t(i, j) \triangleq \Pr [q_t = i, q_{t+1} = j | O, \lambda] \quad (793)$$

$$\tilde{\xi}_t(i, j) \triangleq \Pr [q_t = i, q_{t+1} = j, O | \lambda] \quad (794)$$

$$= \alpha_t(i) a_{ij} b_j(t+1) \beta_{t+1}(j) \quad (795)$$

Remember, knowing the observation sequence does NOT give us the sequence of hidden states.

where you should be able to derive eq. 795 in your head using just logic. If you cannot, review before continuing. We can then derive  $\xi_t(i, j)$  using basic definitions of conditional probability, combined with eq. 791. Finally, we estimate  $\hat{a}_{ij}$  as the expected number of transitions  $q_i \rightarrow q_j$  divided by the expected number of transitions from  $q_i$  total:

$$\hat{a}_{ij} = \frac{\sum_{t=1}^{T-1} \xi_t(i, j)}{\sum_{t=1}^{T-1} \sum_{k=1}^N \xi_t(i, k)} \quad (796)$$

2. **Estimating  $\hat{b}_j(v_k)$ .** We define our estimate as the expected number of times we are in  $q_j$  and emit observation  $v_k$ , divided by the expected number of times we are in state  $j$ . Similar to our approach for  $\hat{a}_{ij}$  we define helper quantities for these values at a given timestep, then sum over them (all  $t$ ) to obtain our estimate.

$$\gamma_t(j) \triangleq \Pr [q_t = j | O, \lambda] \quad (797)$$

$$= \frac{\Pr [q_t = j, O | \lambda]}{\Pr [O | \lambda]} \quad (798)$$

$$= \frac{\alpha_t(j) \beta_t(j)}{\Pr [O | \lambda]} \quad (799)$$

Thus, we obtain  $\hat{b}_j(v_k)$  by summing over all timesteps where  $o_t = v_k$ , denoted as the set  $T_{v_k}$ , divided by the summation over all  $t$  regardless of  $o_t$ :

$$\hat{b}_j(v_k) = \frac{\sum_{t \in T_{v_k}} \gamma_t(j)}{\sum_{t=1}^T \gamma_t(j)} \quad (800)$$

At last we can finally define the **Forward-Backward Algorithm** as follows:

1. Initialize  $A$  and  $B$ .
2. **E-step.** Compute  $\gamma_t(j)$  ( $\forall t, j$ ), and compute  $\xi_t(i, j)$  ( $\forall t, i, j$ ).
3. **M-step.** Update all  $\hat{a}_{ij}$  and  $\hat{b}_j(v_k)$ .
4. Upon convergence, return  $A$  and  $B$ .

## POS Tagging (Ch. 10 3rd Ed.)

Table of Contents Local

*Written by Brandon McKinzie*

**English Parts-of-Speech.** POS are traditionally defined based on syntactic and morphological function, grouping words that have similar neighboring words or take similar affixes.

	Part of Speech	Definition	Properties
<b>Open classes:</b>	noun	people, places, things	occur with determiners, take possessives,
	verb	actions, processes	3rd-person-sg, progressive, past participle
	adjective	properties, qualities	
	adverb	modify verbs, adverbs, verb phrases	directional, locative, degree, manner, temporal
Common nouns can be divided into <i>count</i> (e.g. goat/goats) and <i>mass</i> (e.g. snow) nouns.			

**Closed classes.** POS with relatively fixed membership. Some of the most important in English are:

**prepositions:** on, under, over, near, by, at, from, to, with  
**determiners:** a, an, the  
**pronouns:** she, who, I, others  
**conjunctions:** and, but, or, as, if, when  
**auxiliary verbs:** can, may, should, are  
**particles:** up, down, on, off, in, out, at, by  
**numerals:** one, two, three, first, second, third

Some subtleties: the **particle** resembles a preposition or an adverb and is used in combination with a verb. An example case where “over” is a particle: “she turned the paper over.” When a verb and a particle behave as a single syntactic and/or semantic unit, we call the combination a **phrasal verb**. Phrasal verbs cause widespread problems with NLP because they often behave as a semantic unit with a noncompositional meaning – one that is not predictable from the distinct meanings of the verb and the particle. Thus, “turn down” means something like “reject”, “rule out” means “eliminate”, “find out” is “discover”, and “go on” is “continue”.

**HMM POS Tagging.** Since we typically train on labeled data, we need only use the Viterbi algorithm for decoding<sup>174</sup>. In the POS case, we wish to find the sequence of  $n$  tags,  $\hat{t}_1^n$ , given the observation sequence of  $n$  words  $w_1^n$ .

$$\hat{t}_1^n = \arg \max_{t_1^n} \Pr [t_1^n | w_1^n] = \arg \max_{t_1^n} \Pr [w_1^n | t_1^n] \Pr [t_1^n] \quad (801)$$

where we've dropped the denominator after using Bayes' rule (since argmax is the same). HMM taggers made two further simplifying assumptions:

$$\Pr [w_1^n | t_1^n] \approx \prod_{i=1}^n \Pr [w_i | t_i] \quad (802)$$

$$\Pr [t_1^n] \approx \prod_{i=1}^n \Pr [t_i | t_{i-1}] \quad (803)$$

We can thus plug-in these values into eq. 801 to obtain the equation for  $\hat{t}_1^n$ .

$$\hat{t}_1^n = \arg \max_{t_1^n} \prod_{i=1}^n \Pr [w_i | t_i] \Pr [t_i | t_{i-1}] \quad (804)$$

$$= \arg \max_{t_1^n} \prod_{i=1}^n b_i(w_i) a_{i-1,i} \quad (805)$$

where I've written a “translated” version on the second line using the familiar syntax from the previous chapter. In practice, we can obtain quick estimates for the two probabilities on the RHS by taking counts/averages over our tagged training data. We then run through the Viterbi algorithm to find all the argmaxes over states for the most likely hidden state sequence.

**Maximum Entropy Markov Models** (MEMMs). A sequence model adaptation of the MaxEnt (multinomial logistic regression) classifier<sup>175</sup>. Since HMMs are generative models, they decompose  $\Pr [T | W]$  into  $\Pr [W | T] \Pr [T]$  when computing the best tag sequence  $\hat{T}$ . Since MEMMs are discriminative, they compute/model  $\Pr [T | W]$  directly:

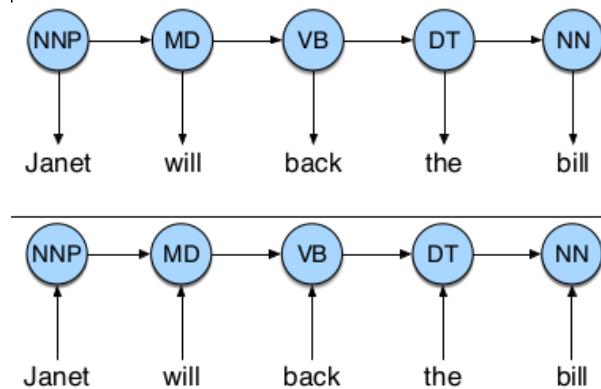
$$\hat{T} = \arg \max_T \Pr [T | W] = \arg \max_T \prod_i \Pr [t_i | w_i, t_{i-1}] \quad (806)$$

Visually, we can think of the difference between HMMs and MEMMs via the direction of arrows, as illustrated below.

---

<sup>174</sup>Recall that decoding is the problem of finding the best hidden state sequence, given  $\lambda = (A, B)$  and observation sequence  $O$ .

<sup>175</sup>Because it is based on logistic regression, the MEMM is a **discriminative sequence model**. By contrast, the HMM is a **generative sequence model**.



The top shows the HMM representation, while the bottom is MEMM.

*The reason to use a discriminative sequence model is that discriminative models make it easier to incorporate a much wider variety of features.*

**Bidirectionality.** The one problem with the MEMM and HMM models as presented is that they are exclusively run left-to-right. MEMMs<sup>176</sup> have a weakness known as the **label bias problem**. Consider the tagged fragment: “will/NN to/TO fight/VB<sup>177</sup>.” Even though the word “will” is followed by “to”, which strongly suggests “will” is a NN, a MEMM will incorrectly label “will” as MD (modal verb). The culprit lies in the fact that  $\Pr[\text{TO} | \text{to}, t_{\text{will}}]$  is essentially 1 regardless of  $t_{\text{will}}$ ; i.e. the fact that “to” must have the tag TO has **explained away** the presence of TO and so the model doesn’t learn the importance of the previous NN tag for predicting TO.

One way to implement bidirectionality (and thus allowing e.g. the link between TO being available when tagging the NN) is to use a **Conditional Random Field** (CRF) model. However, CRFs are much more computationally expensive than MEMMs and don’t work better for tagging.

---

<sup>176</sup> And other non-generative finite-state models based on next-state classifiers

<sup>177</sup> Note on the tag meanings: TO literally means “to”. MD means “modal” and refers to modal verbs such as *will*, *shall*, etc.

## Formal Grammars (Ch. 11 3rd Ed.)

Table of Contents Local

Written by Brandon McKinzie

**Constituency and CFGs.** Discovering the inventory of constituents present in the language. Groups of words like *noun phrases* or *prepositional phrases* can be thought of as single units which can only be placed within certain parts of a sentence.

The most widely used formal system for modeling constituent structure in English is the **Context-Free Grammar**<sup>178</sup>. A CFG consists of a set of **productions** (rules), e.g.

$$\text{NP} \longrightarrow \text{Det Nominal} \tag{807}$$

$$\text{NP} \longrightarrow \text{ProperNoun} \tag{808}$$

$$\text{Nominal} \longrightarrow \text{Noun} \mid \text{Nominal Nominal} \tag{809}$$

where the arrow is to be read “is composed of” or “consists of.”

The sequence of rule expansions going from left to right is called a **derivation** of the string of words, commonly represented by a **parse tree**. The formal language defined by a CFG is the set of strings that are derivable from the designated **start symbol**.

---

<sup>178</sup> Also called Phrase-Structure Grammars. Equiv formalism as Backus-Naur Form (BNF)

## Vector Semantics (Ch. 15)

Table of Contents Local

Written by Brandon McKinzie

**Words and Vectors.** Vector models are generally based on a **co-occurrence**, an example of which is a **term-document matrix**: Each row is identified by a word, and each column a document. A given cell value is the number of times the assoc. word occurred in the assoc. document. Can also view each column as a document vector.

Information Retrieval:  
task of finding document  
 $d$ , from  $D$  docs total,  
that best matches a  
query  $q$ .

For individual word vectors, however, it is most common to instead use a **term-term** matrix<sup>179</sup>, in which columns are also identified by individual words. Now, cell values are the number of times the row (target) word and the column (context) word co-occur in some context in some training corpus. The context is most commonly a window around the row/target word, meaning a cell gives the number of times the column word occurs in a window of  $\pm N$  words from the row word.

- **Q:** What about the co-occurrence of a word with itself (row  $i$ , col  $i$ )?
  - **A:** It is included, yes. Source: “The size of the window ... is generally between 1 and 8 words on each side of the target word (*for a total context of 3-17 words*).”
- **Q:** Why is the size of each vector generally  $|V|$  (vocab size)? Shouldn’t this vary substantially with window and corpus size?
  - **A:** idk

(**TODO:** revisit end of page 5 in my pdf of this)

**Pointwise Mutual Information (PMI).** Motivation: raw frequencies in a co-occurrence matrix aren’t that great, and words like “the” (which aren’t useful and occur everywhere) can really skew things. *The best weighting or measure of association between words should tell us how much more often than chance the two words co-occur.* PMI is such a measure.

$$[\text{Mutual Information}] \quad I(X, Y) = \sum_x \sum_y P(X = x, Y = y) \text{PMI}(x, y) \quad (810)$$

$$[\text{PMI}] \quad \text{PMI}(x, y) = \ln \frac{P(x, y)}{P(x)P(y)} \quad (811)$$

which can be applied for our specific use case as  $\text{PMI}(w, c) = \ln \frac{P(w, c)}{P(w)P(c)}$ . The interpretation is simple: the denominator tells the joint probability of the given target word  $w$  occurring with context word  $c$  if they were independent of each other, while the numerator tells us how often we observed the two words together (assuming we compute probability by using the MLE).

---

<sup>179</sup> Also called the word-word or term-context matrix

Therefore, the ratio gives us how much more the target and feature co-occur than we expect by chance<sup>180</sup>. Most people use **Positive PMI**, which is just  $\max(0, \text{PMI})$ . We can compute a **PPMI matrix** (to replace our co-occurrence matrix), where  $\text{PPMI}_{ij}$  gives the PPMI value of word  $w_i$  with context  $c_j$ . The authors show a few formulas which is really distracting since all we actually need is the counts  $f_{ij} = \text{counts}(w_i, c_j)$ , and from there we can use basic probability and Bayes rule to get the PPMI formula.

- **Q:** Explain why the following is true: very rare words tend to have very high PMI values.
  - **A:** hi
- **Q:** What is the range of  $\alpha$  used in  $\text{PPMI}_\alpha$ ? What is the intuition behind doing this?
  - **A:** For reference:

$$\text{PPMI}_\alpha(w, c) = \max \left( \ln \frac{P(w, c)}{P(w)P_\alpha(c)}, 0 \right) \quad (812)$$

$$P_\alpha(c) = \frac{\text{count}(c)^\alpha}{\sum_{c'} \text{count}(c')^\alpha} \quad (813)$$

Although there are better methods than PPMI for weighted co-occurrence matrices, most notably **TF-IDF**, things like tf-idf are not generally used for measuring *word similarity*. For that, PPMI and significance-testing metrics like t-test and likelihood-ratio are more common. The **t-test** statistic, like PMI, measures how much more frequent the association is than chance.

$$t = \frac{\bar{x} - \mu}{\sqrt{s^2/N}} \quad (814)$$

$$\text{t-test}(a, b) = \frac{P(a, b) - P(a)P(b)}{\sqrt{P(a)P(b)}} \quad (815)$$

where  $\bar{x}$  is the observed mean, while  $\mu$  is the expected mean [under our null-hypothesis of independence].

**Measuring Similarity.** By far the most common similarity metric is the **cosine** of the angle between the vectors:

$$\text{cosine}(\mathbf{v}, \mathbf{w}) = \frac{\mathbf{v} \cdot \mathbf{w}}{|\mathbf{v}| |\mathbf{w}|} \quad (816)$$

Note that, since we've been defining vector elements as frequencies/PPMI values, they won't have negative elements, and thus our cosine similarities will be between 0 and 1 (not -1 and 1).

---

<sup>180</sup>Computing PMI this way can be problematic for word pairs with small probability, especially if we have a small corpus. Recognize that PMI should never really be negative, but in practice this happens for such cases

Alternatives to cosine:

- **Jaccard measure:** Described as "weighted number of overlapping features, normalized", but looks like a silly hack in my opinion:

$$\text{sim}_{Jac}(\mathbf{v}, \mathbf{w}) = \frac{\sum_{i=1}^N \min(\mathbf{v}_i, \mathbf{w}_i)}{\sum_{i=1}^N \max(\mathbf{v}_i, \mathbf{w}_i)} \quad (817)$$

- **Dice measure:** Another hack. This displeases me.

$$\frac{2 \times \sum_{i=1}^N \min(\mathbf{v}_i, \mathbf{w}_i)}{\sum_{i=1}^N (\mathbf{v}_i + \mathbf{w}_i)} \quad (818)$$

- **Jensen-Shannon Divergence:** An alternative to the KL-divergence<sup>181</sup>, which represents the divergence of each distribution from the mean of the two:

$$\text{sim}_{JS}(\mathbf{v} || \mathbf{w}) = D\left(\mathbf{v} \middle\| \frac{\mathbf{v} + \mathbf{w}}{2}\right) + D\left(\mathbf{w} \middle\| \frac{\mathbf{v} + \mathbf{w}}{2}\right) \quad (820)$$

---

<sup>181</sup> Idea: if two vectors,  $\mathbf{v}$  and  $\mathbf{w}$ , each express a probability distribution (their values sum to one), then they are similar to the extent that these probability distributions are similar. The basis of comparing two probability distributions  $P$  and  $Q$  is the **Kullback-Leibler** divergence or relative entropy, defined as:

$$D(P || Q) = \sum_x P(x) \log \frac{P(x)}{Q(x)} \quad (819)$$

## Semantics with Dense Vectors (Ch. 16)

Table of Contents Local

Written by Brandon McKinzie

**Overview.** This chapter introduces three methods for generating short, dense vectors: (1) dimensionality reduction like SVD, (2) neural networks like skip-gram or CBOW, and (3) Brown clustering.

**Dense Vectors via SVD.** Method for finding more important dimensions of a dataset, “important” defined as dimensions wherein the data most varies. First applied (for language) for generating embeddings from term-document matrices in a model called **Latent Semantic Analysis** (LSA). LSA is just SVD on a  $|V| \times c$  term-document matrix  $\mathbf{X}$ , factorized into  $\mathbf{W}\Sigma\mathbf{C}^T$ . By using only the top  $k < m$  dimensions of these three matrices, the product becomes a least-squares approx. to the original  $\mathbf{X}$ . It also gives us the reduced  $|V| \times k$  matrix  $\mathbf{W}_k$ , where each row (word) is a  $k$ -dimensional vector (embedding). Voilà, we have our dense vectors!

$$\mathbf{W} \in \mathbb{R}^{|V| \times m}$$

$$\Sigma \in \mathbb{R}^{m \times m}$$

$$\mathbf{C}^T \in \mathbb{R}^{m \times c}$$

Note that LSA implementations typically use a particular weighting of each cell in the term-document matrix called the **local** and **global** weights.

$$[\text{local}] \quad \log f(i, j) + 1 \quad (821)$$

$$[\text{global}] \quad 1 + \frac{\sum_j p(i, j) \log p(i, j)}{\log D} \quad (822)$$

$f(i, j)$  is the raw frequency of word  $i$  in context  $j$ .  $D$  is number of docs.

For the case of a word-word matrix, it is common to use PPMI weighting.

**Skip-Gram and CBOW.** Neural models learn an embedding by starting with a random vector and then iteratively shifting a word’s embeddings to be more like the embeddings of neighboring words, and less like the embeddings of words that don’t occur nearby<sup>182</sup> Word2vec, for example, learns embeddings by training to predict neighboring words<sup>183</sup>.

- **Skip-Gram:** Learns two embeddings for each word  $w$ : the **word embedding**  $v$  (within matrix  $\mathbf{W}$ ) and **context embedding**  $c$  (within matrix  $\mathbf{C}$ ). Visually:

$$\mathbf{W} = \begin{pmatrix} \mathbf{v}_0^T \\ \mathbf{v}_1^T \\ \vdots \\ \mathbf{v}_{|V|}^T \end{pmatrix} \quad \mathbf{C} = \begin{pmatrix} \mathbf{c}_0 & \mathbf{c}_1 & \cdots & \mathbf{c}_{|V|} \end{pmatrix} \quad (823)$$

<sup>182</sup>Why? Why is this a sensible assumption? I see no reason a priori why it ought to be true.

<sup>183</sup>Note that the prediction task is not the goal – it just happens to result in good word embeddings. Hacky.

For a context window of  $L = 2$ , and at a given word  $\mathbf{v}^{(t)}$  inside the corpus<sup>184</sup>, our goal is to predict the context [words] denoted as  $[\mathbf{c}^{(t-2)}, \mathbf{c}^{(t-1)}, \mathbf{c}^{(t+1)}, \mathbf{c}^{(t+2)}]$ .

- Example: Consider one of the context words, say  $\mathbf{c}^{(t+1)} \triangleq \mathbf{c}_k$ , where we also assume it's the  $k$ th word in our vocab. Also assume that our target word  $\mathbf{v}^{(t)} \triangleq \mathbf{v}_j$  is the  $j$ th word in our vocab.
- Our task is to compute  $\Pr[\mathbf{c}_k | \mathbf{v}_j]$ . We do this with a softmax:

$$\Pr[\mathbf{c}_k | \mathbf{v}_j] = \frac{e^{\mathbf{c}_k^T \mathbf{v}_j}}{\sum_{i \in |V|} e^{\mathbf{c}_i^T \mathbf{v}_j}} \quad (824)$$

- **CBOW**: Continuous bag of words. Basically the mirror-image of skip-gram. Goal is to predict current word  $\mathbf{v}^{(t)}$  from the context window of  $2L$  words  $[\mathbf{c}^{(t-2)}, \mathbf{c}^{(t-1)}, \mathbf{c}^{(t+1)}, \mathbf{c}^{(t+2)}]$ .

As usual, the denominator of the softmax is computationally expensive, and usually we approximate it with **negative sampling**.

*In the training phase, the algorithm walks through the corpus, at each target word choosing the surrounding context words as positive examples, and for each positive example also choosing  $k$  noise samples or negative samples: non-neighbor words. The goal will be to move the embeddings toward the neighbor words and away from the noise words.*

Example: Suppose we come along the following window (in “[]”) ( $L=2$ ) in our corpus:

lemon, a [tablespoon of apricot preserves or] jam

Ultimately, we want dot products,  $\mathbf{c}_i \cdot \text{vector}(\text{"apricot"})$ , to be high for all four of the context words  $\mathbf{c}_i$ . We do negative sampling by sampling  $k$  random noise words according to their [unigram] frequency. So here, for e.g.  $k = 2$ , this would amount to 8 noise words, 2 for each context word. We want the dot products between “apricot” and these noise words to be low. For a given single context-word pair  $(w, c)$ , our training objective is to maximize:

$$\log \sigma(c \cdot w) + \sum_{i=1}^k \mathbb{E}_{w_i \sim p(w)} [\log \sigma(-w_i \cdot w)] \quad (825)$$

In practice, common to use  $p^{3/4}(w)$  instead of  $p(w)$

Again, the above is for a single context-target word-pair and, accordingly, the summation is only over  $k = 2$  (for our example). Don't try to split the expectation into a summation or anything – just view it as an expected value. To iteratively shift parameters, we use an optimizer like SGD.

---

<sup>184</sup>Note that, technically, the position  $t$  of  $\mathbf{v}^{(t)}$  is irrelevant for our computation; we are predicting those words based on which word  $\mathbf{v}^{(t)}$  is in the vocabulary, not its position in the corpus.

The actual model architecture is a typical neural net, progressing as follows:

$$\text{"apricot"} \rightarrow \mathbf{w}^{\text{one-hot}} = [0 \ 0 \ \dots \ 1 \ \dots \ 0] \quad (826)$$

$$\rightarrow \mathbf{h} = \mathbf{W}^T \mathbf{w}^{\text{one-hot}} \quad (827)$$

$$\rightarrow \mathbf{o} = \mathbf{C}^T \mathbf{h} = [\mathbf{c}_0^T \mathbf{h}, \mathbf{c}_1^T \mathbf{h}, \dots, \mathbf{c}_{|V|}^T \mathbf{h}]^T \quad (828)$$

$$\rightarrow \mathbf{y} = \text{softmax}(\mathbf{o}) = [\Pr[\mathbf{c}_0|\mathbf{h}], \Pr[\mathbf{c}_2|\mathbf{h}], \dots, \Pr[\mathbf{c}_{|V|}|\mathbf{h}]]^T \quad (829)$$

$$(830)$$

**Brown Clustering.** An agglomerative clustering algorithm for deriving vector representations of words by clustering words based on their associations with the preceding or following words. Makes use of the **class-based language model** (CBLM), wherein each word  $w$  belongs to some class  $c \in C$  via the probability  $P(w | c)$ . CBLMs define

$$P(w_i | w_{i-1}) = P(c_i | c_{i-1})P(w_i | c_i) \quad (831)$$

$$P(\text{corpus} | C) = \prod_{i=1}^n P(w_i | w_{i-1}) \quad (832)$$

A naive and extremely inefficient version of Brown clustering, a hierarchical clustering algorithm, is as follows:

1. Each word is initially assigned to its own cluster.
2. For each cluster pair  $(c_i, c_{j \neq i})$ , compute the value of eq 832 that would result from merging  $c_i$  and  $c_j$  into a single cluster. The pair whose merger results in the smallest decrease in eq 832 is merged.
3. Clustering proceeds until all words are in one big cluster.

This process builds a binary tree from the bottom-up, and the binary string corresp. to a word's traversal from leaf-to-root is its representation.

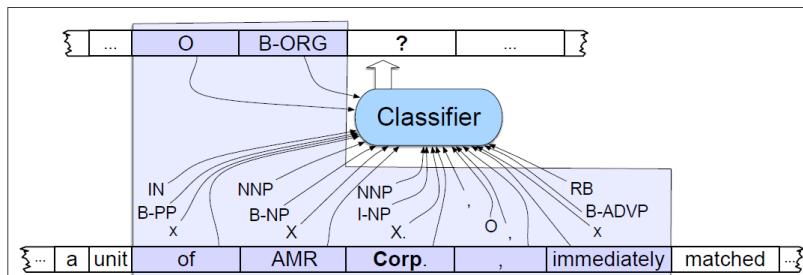
## Information Extraction (Ch. 21 3rd Ed)

Table of Contents Local

Written by Brandon McKinzie

**Overview.** The first step in most IE tasks is **named entity recognition** (NER). Next we can do **relation extraction**: finding and classifying semantic relations among the entities, e.g. “spouse-of.” **Event extraction** is finding the events in which the entities participate, and **event coreference** for figuring out which event mentions actually refer to the same event. It’s also common to extract dates/times (temporal expression) and perform **temporal expression normalization** to map them onto specific calendar dates. Finally, we can do **template filling**: finding recurring/stereotypical situations in documents and filling the template slots with appropriate material.

**Named Entity Recognition.** Standard algorithm is word-by-word sequence labeling task by a MEMM or CRF, trained to label tokens with tags indicating the presence of particular kinds of NEs. It is common to label with **BIO tagging**, for beginning, inside, and outside of entities. If we have  $n$  unique entity types, then we’d have  $2n + 1$  BIO tags<sup>185</sup>. A helpful illustration is shown below:



**Figure 21.7** Named entity recognition as sequence labeling. The features available to the classifier during training and classification are those in the boxed area.

Here we see a classifier determining the label for *Corp.* with a context window of size 2 and various features shown in the boxed region. For evaluation of NER, we typically use the familiar **recall**, **precision**, and **F1 measure**.

<sup>185</sup>  $2n$  for  $B- <NE>$  and  $I- <NE>$ , with  $+1$  for the blanket  $O$  tag (not any of our NEs)

**Relation Extraction.** The four main algorithm classes used are (1) hand-written patterns, (2) supervised ML, (3) semi-supervised, and (4) unsupervised. Terminology:

- **Infobox:** structured tables associated with certain articles/topics/etc. For example, the Wikipedia infobox for Stanford includes structured facts like state = 'California'.
- **Resource Description Framework (RDF):** a metalanguage of RDF triples, tuples of (entity, relation, entity), called a subject-predicate-object expression. For example: (Golden Gate Park, location, San Francisco).
- **hypernym:** the “is-a” relation.
- **hyponym:** the “kind-of” relation. *Gelidium is a kind of red algae.*

Overview of the four algorithm classes:

1. **Patterns.** Consider a sentence that has the following form:

$NP_0 \text{ such as } NP_1\{\}, NP_2, \dots, (\text{and/or}) NP_i\}, i \geq 1$

also known as a *lexico-syntactic pattern*, which implies  $\forall NP_i, i \geq 1, \text{hyponym}(NP_i, NP_0)$ <sup>186</sup>. Patterns typically have high precision but low-recall.

2. **Supervised.** The general approach for finding relations in a given sequence of words is the following:

- (a) Find all pairs of named entities in the sequence (typically a single sentence).
- (b) For each pair, use a trained binary classifier to predict whether or not the entities in the pair are indeed related.
- (c) If related, use a classifier trained to predict the relation given the entity-pair.

As with NER, **the most important step** in this process is to identify useful surface features that will be useful for relation classification, including word features, NER features, syntactic paths (chunk seqs, constituent paths, dependency-tree paths), and more.

3. **Semi-supervised** via bootstrapping. Suppose we have a few high-precision **seed patterns** (or seed tuples)<sup>187</sup>. **Bootstrapping** proceeds by taking the entities in the seed pair, and then finding sentences (on the web, or whatever dataset we are using) that contain both entities. From all such sentences, we extract and generalize the context around the entities to learn new patterns.

```

function BOOTSTRAP(Relation R) returns new relation tuples
    tuples  $\leftarrow$  Gather a set of seed tuples that have relation R
    iterate
        sentences  $\leftarrow$  find sentences that contain entities in seeds
        patterns  $\leftarrow$  generalize the context between and around entities in sentences
        newpairs  $\leftarrow$  use patterns to grep for more tuples
        newpairs  $\leftarrow$  newpairs with high confidence
        tuples  $\leftarrow$  tuples + newpairs
    return tuples

```

**Figure 21.14** Bootstrapping from seed entity pairs to learn relations.

<sup>186</sup>Here, hyponym(A, B) means “A is a kind-of (hyponym) of B.”

<sup>187</sup>seed tuples are tuples of the general form (M1, M2) where M1 and M2 are each specific named entities we know have the relation of interest R.

4. **Unsupervised.** The Re Verb system extracts a relation from a sentence  $s$  in 4 steps:

1. Run a part-of-speech tagger and entity chunker over  $s$
2. For each verb in  $s$ , find the longest sequence of words  $w$  that start with a verb and satisfy syntactic and lexical constraints, merging adjacent matches.
3. For each phrase  $w$ , find the nearest noun phrase  $x$  to the left which is not a relative pronoun, wh-word or existential “there”. Find the nearest noun phrase  $y$  to the right.
4. Assign confidence  $c$  to the relation  $r = (x, w, y)$  using a confidence classifier and return it.

**Event Extraction.** An event mention is any expression denoting an event or state that can be assigned to a particular point, or interval, in time. Note that this is quite different than the colloquial usage of the word “event,” you should think of the two as distinct. Here, most event mentions correspond to verbs, and most verbs introduce events. Event extraction is typically modeled via ML, detecting events via sequence models with BIO tagging, and assigning event classes/attributes with multi-class classifiers.

**Template Filling.** The task is creation of one template for each event in the input documents, with the slots filled with text from the document. For example, an event could be “Fare-Raise Attempt” with corresponding template (slots to be filled) “(<Lead Airline>, <Amount>, <Effective Date>, <Follower>)”. This is generally modeled by training two separate supervised systems:

1. **Template recognition.** Trained to determine if template T is present in sentence S. Here, “present” means there is a sequence within the sentence that could be used to fill a slot within template T.
2. **Role-filler extraction.** Trained to detect each role (slot-name), e.g. “Lead Airline”.

# PROBABILISTIC GRAPHICAL MODELS

## CONTENTS

7.1	Foundations (Ch. 2) . . . . .	282
7.1.1	Appendix . . . . .	284
7.1.2	L-BFGS . . . . .	287
7.1.3	Exercises . . . . .	290
7.2	The Bayesian Network Representation (Ch. 3) . . . . .	292
7.3	Undirected Graphical Models (Ch. 4) . . . . .	296
7.3.1	Exercises . . . . .	303
7.4	Local Probabilistic Models (Ch. 5) . . . . .	305
7.5	Template-Based Representations (Ch. 6) . . . . .	306
7.6	Gaussian Network Models (Ch. 7) . . . . .	309
7.7	Variable Elimination (Ch. 9) . . . . .	310
7.8	Clique Trees (Ch. 10) . . . . .	314
7.9	Inference as Optimization (Ch. 11) . . . . .	319
7.10	Parameter Estimation (Ch. 17) . . . . .	320
7.11	Partially Observed Data (Ch. 19) . . . . .	322

## Foundations (Ch. 2)

Table of Contents Local

Written by Brandon McKinzie

(Brief summary of the book's notation, for ease of reference.)

**Graphs.** Authors denote directed graphs as  $\mathcal{G}$  and undirected graphs as  $\mathcal{H}$ .

- **Induced Subgraph.** Let  $\mathcal{K} = (\mathcal{X}, \mathcal{E})$  and  $\mathbf{X} \subset \mathcal{X}$ . Define the *induced subgraph*  $\mathcal{K}[\mathbf{X}]$  to be the graph  $(\mathbf{X}, \mathcal{E}')$  where  $\mathcal{E}'$  are all edges  $X \Rightarrow Y$  such that  $X, Y \in \mathbf{X}$ .
- **Complete Subgraph.** A subgraph over  $\mathbf{X}$  is *complete* if every two nodes in  $\mathbf{X}$  are connected by some edge. The set  $\mathbf{X}$  is often called a *clique*; we say that a clique  $\mathbf{X}$  is *maximal* if for any superset of nodes  $\mathbf{Y} \supset \mathbf{X}$ ,  $\mathbf{Y}$  is not a clique.
- **Upward Closure.** We say that a subset of nodes  $\mathbf{X} \subset \mathcal{X}$  is *upwardly closed* in  $\mathcal{K}$  if  $\forall X \in \mathbf{X}$ , we have that  $\text{Boundary}_X \subset \mathbf{X}$ <sup>188</sup>. We define the *upward closure* of  $\mathbf{X}$  to be the minimally upwardly closed subset  $\mathbf{Y}$  that contains  $\mathbf{X}$ . We define the *upwardly closed subgraph* of  $\mathbf{X}$ , denoted  $\mathcal{K}^+[\mathbf{X}]$ , to be the induced subgraph over  $\mathbf{Y}$ ,  $\mathcal{K}[\mathbf{Y}]$ .

**Paths and Trails.** Definitions for longer-range connections in graphs. We use the notation  $X_i \Rightarrow X_{j+1}$  to denote that  $X_i$  and  $X_j$  are connected via some edge, whether directed (in any direction) or undirected.

- **Trail/Path.** We say that  $X_1, \dots, X_k$  form a *trail* in the graph  $\mathcal{K} = (\mathcal{X}, \mathcal{E})$  if  $\forall i = 1, \dots, k-1$ , we have that  $X_i \Rightarrow X_{i+1}$ . A *path* makes an additional restriction: either  $X_i \rightarrow X_{i+1}$  or  $X_i \leftarrow X_{i+1}$ .
- **Connected Graph.** A graph is *connected* if  $\forall X_i, X_j$  there is a trail between  $X_i$  and  $X_j$ .
- **Cycle.** A *cycle* in  $\mathcal{K}$  is a directed path  $X_1, \dots, X_k$  where  $X_1 = X_k$ .
- **Loop.** A *loop* in  $\mathcal{K}$  is a trail where  $X_1 = X_k$ . A graph is *singly connected* if it contains no loops. A node in a singly connected graph is called a *leaf* if it has exactly one adjacent node.
- **Polytree/Forest.** A singly connected graph is also called a *polytree*. A singly connected undirected graph is called a *forest*; if a forest is also connected, it is called a *tree*.
  - A directed graph is a forest if each node has at most one parent. A directed forest is a tree if it is also connected.
- **Chordal Graph.** Let  $X_1 - X_2 - \dots - X_k - X_1$  be a loop in the graph. A *chord* in the loop is an edge connecting  $X_i$  and  $X_j$  for two nonconsecutive nodes  $X_i, X_j$ . An undirected graph  $\mathcal{H}$  is said to be *chordal* if any loop  $X_1 - X_2 - \dots - X_k - X_1$  for  $k \geq 4$  has a chord.

---

<sup>188</sup> $\text{Boundary}_X \triangleq \text{Pa}_X \cup \text{Nb}_X$ . For DAGs, this is simply  $X$ 's parents, and for undirected graphs  $X$ 's neighbors.

**Probability.** Some notational reminders for this book. Let  $\Omega$  denote a space of possible outcomes, and let  $S$  denote a set of measurable **events**  $\alpha$ , each of which are a subset of  $\Omega$ .

A probability distribution  $P$  over  $(\Omega, S)$  is a mapping from events in  $S$  to real values that satisfy:

- $P(\alpha) \geq 0$  for all  $\alpha \in S$ .
- $P(\Omega) = 1$ .
- If  $\alpha, \beta \in S$  and  $\alpha \cap \beta = \emptyset$ , then  $P(\alpha \cup \beta) = P(\alpha) + P(\beta)$ .

Some useful independence properties:

$$\text{Symmetry : } (X \perp Y \mid Z) \implies (Y \perp X \mid Z) \quad (833)$$

$$\text{Decomposition : } (X \perp (Y, W) \mid Z) \implies (X \perp Y \mid Z) \quad (834)$$

$$\text{Weak Union : } (X \perp (Y, W) \mid Z) \implies (X \perp Y \mid Z, W) \quad (835)$$

$$\text{Contraction : } (X \perp W \mid Z, Y) \& (X \perp Y \mid Z) \implies (X \perp Y, W \mid Z) \quad (836)$$

### My Proofs: Independence Properties

I'll be using the definition that  $(X \perp Y \mid Z) \Leftrightarrow P(X, Y \mid Z) = P(X \mid Z)P(Y \mid Z)$ . Given this definition the proof for the symmetry property is trivial. In what follows, I'll assume the LHS of the given implication is true, and then show that the RHS must hold as well.

**Decomposition:**

$$P(X, Y \mid Z) = \sum_w P(X, Y, w \mid Z) = \sum_w P(X \mid Z)P(Y, w \mid Z) = P(X \mid Z)P(Y \mid Z) \quad \checkmark$$

**Weak Union:**

$$P(X, Y \mid Z, W) = \frac{P(X, Y, W \mid Z)}{P(W \mid Z)} \quad (837)$$

$$= \frac{P(X \mid Z)P(Y, W \mid Z)}{P(W \mid Z)} \quad (838)$$

$$= \frac{P(X \mid Z)P(W \mid Z)P(Y \mid Z, W)}{P(W \mid Z)} \quad (839)$$

$$= P(X \mid Z, W)P(Y \mid Z, W) \quad \checkmark \quad (840)$$

**Contraction:**

$$P(X, Y, W \mid Z) = P(Y \mid Z)P(X, W \mid Z, Y) \quad (841)$$

$$= P(Y \mid Z)P(X \mid Z, Y)P(W \mid Z, Y) \quad (842)$$

$$= P(X \mid Z)[P(Y \mid Z)P(W \mid Z, Y)] \quad (843)$$

$$= P(X \mid Z)P(Y, W \mid Z) \quad \checkmark \quad (844)$$

We now define what “positive distribution” means, and a useful property of such distributions.

*A distribution  $P$  is said to be **positive** if for all events  $\alpha \in S$ , such that  $\alpha \neq \emptyset$ , we have that  $P(\alpha) > 0$ .*

For positive distributions, and for mutually disjoint sets  $\mathbf{X}, \mathbf{Y}, \mathbf{Z}, \mathbf{W}$ , the **intersection** property also holds:

$$\text{Intersection : } (\mathbf{X} \perp \mathbf{Y} | \mathbf{Z}, \mathbf{W}) \& (\mathbf{X} \perp \mathbf{W} | \mathbf{Z}, \mathbf{Y}) \implies (\mathbf{X} \perp \mathbf{Y}, \mathbf{W} | \mathbf{Z}) \quad (845)$$


---

### 7.1.1 APPENDIX

---

Figured this would be a good place to put some of the definitions in the Appendix, too.

#### Information Theory (A.1).

- **Entropy.** Let  $P(X)$  be a distribution over a random variable  $X$ . The *entropy* of  $X$  is defined as.

$$H_P(X) = \mathbb{E}_P \left[ \lg \frac{1}{P(X)} \right] = \sum_x P(X = x) \lg \frac{1}{P(X = x)} \quad (846)$$

$$0 \leq H_P(X) \leq \lg |Val(X)| \quad (847)$$

We treat  $0 \log \frac{1}{0} = 0$

$H_p(X)$  is a lower bound for the expected number of bits required to encode instances sampled from  $P(X)$ . Another interpretation is that the entropy is a measure of our uncertainty about the value of  $X$ .

- **Conditional Entropy.** The *conditional entropy* of  $X$  given  $Y$  is

$$H_P(X | Y) = H_P(X, Y) - H_P(Y) = \mathbb{E}_P \left[ \lg \frac{1}{P(X | Y)} \right] \quad (848)$$

$$H_P(X | Y) \leq H_P(X) \quad (849)$$

which captures the additional cost (in bits) of encoding  $X$  when we’re already encoding  $Y$ .

- **Mutual Information.** The *mutual information* between  $X$  and  $Y$  is

$$I_P(X; Y) = H_P(X) - H_P(X | Y) = \mathbb{E}_P \left[ \lg \frac{P(X | Y)}{P(X)} \right] \quad (850)$$

which captures how many bits we save (on average) in the encoding of  $X$  if we know the value  $Y$ .

- **Distance Metric.** A *distance metric* is any distance measure  $d$  evaluating the distance between two distributions that satisfies all of the following properties:

- **Positivity:**  $d(P, Q) \geq 0$  and  $d(P, Q) = 0$  if and only if  $P = Q$ .

- **Symmetry:**  $d(P, Q) = d(Q, P)$ .
- **Triangle inequality:** For any three distributions  $P, Q, R$ , we have that

$$d(P, R) \leq d(P, Q) + d(Q, R) \quad (851)$$

- **Kullback-Liebler Divergence.** Let  $P$  and  $Q$  be two distributions over random variables  $X_1, \dots, X_n$ . The relative entropy, or KL-divergence, of  $P$  and  $Q$  is

$$D(P\|Q) = \mathbb{E}_P \left[ \lg \frac{P(X_1, \dots, X_n)}{Q(X_1, \dots, X_n)} \right] \quad (852)$$

Note that this only satisfies the positivity property, and is thus not a true distance metric.

### Algorithms and Algorithmic Complexity (A.3).

- **Decision Problems.** A decision problem  $\Pi$  is a task that accepts an input (instance)  $\omega$  and decides whether it satisfies a certain condition or not. The **SAT** problem accepts a formula in propositional logic and decides whether there is an assignment to the variables in the formula such that it evaluates to true. **3-SAT** restricts this to accepting only formulas in conjunctive normal form (CNF), and further restricted s.t. each clause contains at most 3 literals.
- **P.** A decision problem is in the class  $\mathcal{P}$  if there exists a deterministic algorithm that takes an instance  $\omega$  and determines whether  $\omega \in \mathcal{L}_\Pi$  (the set of instances for which a correct algorithm must return true), in polynomial time in the size of the input  $\omega$ .
- **NP.** A *non-deterministic algorithm* takes the general form: (1) nondeterministically guess some assignment  $\gamma$  to the variables of  $\omega$ , (2) deterministically verify whether  $\gamma$  satisfies the condition of the problem. The algorithm will repeat these steps until it produces a  $\gamma$  that satisfies the problem. A decision problem  $\Pi$  is in the class  $\mathcal{NP}$  if there exists a nondeterministic algorithm that accepts  $\omega$  if and only if  $\omega \in \mathcal{L}_\Pi$ , and if the verification stage can be executed in polynomial time in the length of  $\omega$ .
- **NP-hard.**  $\Pi$  is  $\mathcal{NP}$ -hard if for every DP  $\Pi' \in \mathcal{NP}$ , there is a polynomial-time transformation of inputs such that an input for  $\Pi'$  belongs to  $\mathcal{L}_{\Pi'}$  if and only if the transformed instance belongs to  $\mathcal{L}_\Pi$ . Note that  $\mathcal{NP}$ -hard is a superset of  $\mathcal{NP}$ .
- **NP-complete.** A problem  $\Pi$  is said to be  $\mathcal{NP}$ -complete if it is both  $\mathcal{NP}$ -hard and in  $\mathcal{NP}$ .

The SAT problem is  $\mathcal{NP}$ -hard.

**Combinatorial Optimization and Search** (A.4). Below, I'll outline some common search algorithms. These are designed to address the following task:

*Given initial candidate solution  $\sigma_{cur}$ , a score function  $score$ , and a set of **search operators**  $\mathcal{O}$ , search for the optimal solution  $\sigma_{best}$  that maximizes the value of  $score(\sigma_{best})$*

### Greedy local search (Algorithm A.5)

Repeat the following until *didUpdate* evaluates to *false* at the end of an iteration.

1. Initialize  $\sigma_{best} := \sigma_{cur}$ .
2. Set *didUpdate* := *false*.
3. For each operator  $o \in \mathcal{O}$ , do:
  - (a) Let  $\sigma_o := o(\sigma_{best})$ .
  - (b) If  $\sigma_o$  is legal solution, and  $score(\sigma_o) > score(\sigma_{best})$ , reassign  $\sigma_{best} := \sigma_o$ , and set *didUpdate* := *true*.
4. If *didUpdate* == *true*, go back to step 2. Otherwise terminate and return  $\sigma_{best}$ .

### Beam search (Algorithm A.7)

We are given a **beam width**  $K$ . Initialize our *beam*, the set of at most  $K$  solutions we are currently tracking, to  $\{\sigma_{cur}\}$ . Repeat the following until termination<sup>a</sup>:

1. Initial the set of successor states  $H := \emptyset$ .
2. For each solution  $\sigma \in Beam$ , and each operator  $o \in \mathcal{O}$ , insert a candidate successor state  $o(\sigma)$  into  $H$ .
3. Set  $Beam := K\text{BestScore}(H)$ <sup>b</sup>.

Once termination is reached, return the best solution  $\sigma_{best}$  in *Beam*.

<sup>a</sup>Termination condition could be e.g. an upper bound on number of iterations or on the improvement achieved in the last iteration.

<sup>b</sup>Notice that this implies an underlying assumption of beam search: all successor states  $\sigma \in H$  have scores greater than any of the states in the current beam. We always assume improvement.

## Continuous Optimization (A.5).

- **Line Search.** Method for adaptively choosing the step size (learning rate)  $\eta$  at each training step. Assuming we are doing gradient ascent. We'd usually set the parameters  $\theta$  at step  $t+1$  to  $\theta^{(t)} + \eta \nabla f(\theta^{(t)})$ . Line search modifies this by instead defining the “line”  $g(\eta)$  below, and searching for the optimal value of  $\eta$  along that line.

$$g(\eta) = \vec{\theta}^{(t)} + \eta \nabla f(\theta^{(t)}) \quad (853)$$

$$\eta^{(t)} = \arg \max_{\eta} [g(\eta)] \quad (854)$$

$$\theta^{(t+1)} \leftarrow \theta^{(t)} + \eta^{(t)} \nabla f(\theta^{(t)}) \quad (855)$$

At risk of stating the obvious,  $g(\eta)$  is referred to as a “line” because it’s a function of the form  $\mathbf{m}\mathbf{x} + \mathbf{b}$ .

### 7.1.2 L-BFGS

Some notes from the textbook “Numerical Optimization” (chapters 8 and 9).

**The BFGS Method** (8.1). Begin the derivation by forming the following quadratic model of the objective function  $f$  at the current iterate<sup>189</sup>  $\theta_t$ :

$$m_t(p) = f_t + \nabla f_t^T p + \frac{1}{2} p^T B_t p \quad (856)$$

For  $m_t(p)$ ,  $p$  denotes the deviation at step  $t$  from the current parameters  $\theta_t$ .

where  $B_t$  is an  $n \times n$  symmetric p.d. matrix that will be revised/updated every iteration (it is *not* the Hessian!). The minimizer  $p_t$  of this function can be written explicitly

$$p_t = -B_t^{-1} \nabla f_t \quad (857) \quad \frac{\partial}{\partial p} \frac{1}{2} p^T B_t p = B_t p$$

is used as the search direction, and the new iterate is

$$\theta_{t+1} \leftarrow \theta_t + \alpha_t p_t \quad (858)$$

where the step length  $\alpha_t$  is chosen to satisfy the **Wolfe conditions**<sup>190</sup>. This is basically Newton’s method with line search, except that we’re using the *approximate Hessian*  $B_t$  instead of the true Hessian. It would be nice if we could somehow avoid recomputing  $B_t$  at each step. One proposed method involves imposing conditions on  $B_{t+1}$  based on the previous step(s). Require that  $\nabla m_{t+1}$  equal  $\nabla f$  at the latest two iterates  $\theta_t$  and  $\theta_{t+1}$ . Formally, the two conditions can be written as

$$\nabla m_{t+1}(-\alpha_t p_t) = \nabla f_{t+1} - \alpha_t B_{t+1} p_t = \nabla f_t \quad (861)$$

$$\nabla m_{t+1}(0) = \nabla f_{t+1} \quad (862)$$

We can rearrange the first condition to obtain the **secant equation**:

$$H_{t+1} y_t = s_t \quad \text{where} \quad (863)$$

$$H_{t+1} \triangleq B_{t+1}^{-1} \quad (864)$$

$$s_t \triangleq \theta_{t+1} - \theta_t \quad (865)$$

$$y_t \triangleq \nabla f_{t+1} - \nabla f_t \quad (866)$$

which is true only if  $s_t$  and  $y_t$  satisfy the **curvature condition**,  $s_t^T y_t > 0$ . The curvature condition is guaranteed to hold if we impose the Wolfe conditions on the line search. As is, this

<sup>189</sup>Recall that an “iterate” is just some variable that gets iteratively computed/updated. Fancy people with fancy words.

<sup>190</sup>The Wolfe conditions are the following sufficient decrease and curvature conditions for line search:

$$f(\theta_t + \alpha_t p_t) \leq f(\theta_t) + c_1 \alpha_t \nabla f_t^T p_t \quad (859)$$

$$\nabla f(\theta_t + \alpha_t p_t)^T p_t \geq c_2 \nabla f_t^T p_t \quad (860)$$

for some constant  $c_1 \in (0, 1)$  and  $c_2 \in (c_1, 1)$ .

still has infinitely many solutions for  $H_{t+1}$ . To determine it uniquely, we impose the additional condition that  $H_{t+1}$  is the closest of all possible solutions to the current  $H_t$ :

$$\min_H \|H - H_t\|_W \quad \text{s.t.} \quad H = H^T, \quad Hy_t = s_t \quad (867)$$

where  $\|\cdot\|_W$  is the *weighted Frobenius norm*<sup>191</sup>, and  $W$  can be any matrix satisfying  $Ws_t = y_t$ . For concreteness, assume that  $W = \tilde{G}_t$ , where

$$\tilde{G}_t = \int_0^1 \nabla^2 f(\theta_t + \tau \alpha_t p_t) d\tau \quad (870)$$

The unique solution to  $H_{t+1}$  is given by

$$(\text{BFGS}) \quad H_{t+1} = (I - \rho_t s_t y_t^T) H_t (I - \rho_t y_t s_t^T) + \rho_t s_t s_t^T \quad (8.16)$$

where  $\rho_t = 1/(y_t^T s_t)$ . The BFGS is summarized in algorithm 8.1 below.

**Algorithm 8.1** (BFGS Method). Given starting point  $\theta_0$ , convergence tolerance  $\epsilon > 0$ , and inverse Hessian approximation  $H_0$ . Initialize  $t = 0$ . While  $\|\nabla f_t\| > \epsilon$  do:

1. Compute search direction  $p_t = -H_t \nabla f_t$ .
2. Set  $\theta_{t+1} = \theta_t + \alpha_t p_t$ , where  $\alpha_t$  is computed via line search to satisfy the Wolfe conditions.
3. Define  $s_t = \theta_{t+1} - \theta_t$  and  $y_t = \nabla f_{t+1} - \nabla f_t$ .
4. Compute  $H_t$  by means of equation 8.16.
5. Increment  $t += 1$  and go back to step 1.

---

<sup>191</sup>1

$$\|H\|_W \triangleq \|W^{1/2} H W^{1/2}\|_F \quad (868)$$

$$\|C\|_F^2 \triangleq \sum_{i,j} c_{ij}^2 \quad (869)$$

**L-BFGS.** Modifies BFGS to store a modified version of  $H_t$  implicitly, by storing some number  $m$  of vector pairs  $\{s_i, y_i\}$ , corresponding to the  $m$  most recent time steps. We use a recursive procedure to compute  $H_t \nabla f_t$  given the set of vectors.

**Algorithm 9.1** (L-BFGS two-loop recursion) Subroutine of L-BFGS for computing  $H_t \nabla f_t$ . We're given the current value of  $\nabla f_t$ , and we initialize  $q$  to this value.

1. For  $i$  in the range  $[t - 1, t - m]$ , compute

$$\alpha_i \leftarrow \rho_i s_i^T q \quad (871)$$

$$q \leftarrow q - \alpha_i y_i \quad (872)$$

2. Set  $r \rightarrow H_t^0 q$ .
3. For  $i$  in the range  $[t - m, t - 1]$ , compute

$$\beta \leftarrow \rho_i y_i^T r \quad (873)$$

$$r \leftarrow r + s_i(\alpha_i - \beta) \quad (874)$$

4. Return result  $H_t \nabla f_t = r$ .

**Algorithm 9.2** (L-BFGS). Given starting point  $\theta_0$ , integer  $m > 0$ , and initial  $t = 0$ . Repeat the following until convergence.

1. Choose  $H_t^0$ . A popular choice is  $H_t^0 := \gamma_t I$ , where

$$\gamma_t \triangleq \frac{s_{t-1}^T y_{t-1}}{y_{t-1}^T y_{t-1}}$$

2. Compute  $p_t \leftarrow -H_t \nabla f_t$  from Algorithm 9.1.
3. Compute  $\theta_{t+1} \leftarrow \theta_t + \alpha_t p_t$ , where  $\alpha_t$  is chosen to satisfy the Wolfe conditions.
4. if  $t > m$ , discard  $\{s_{t-m}, y_{t-m}\}$ . Compute and save  $s_t$  and  $y_t$ .
5. Increment  $t += 1$  and go back to step 1.

---

### 7.1.3 EXERCISES

---

Going through all the problems with a star for review.

#### Exercise 2.4

Let  $\alpha \in S$  be an event s.t.  $P(\alpha) > 0$ . Show that  $P(\cdot | \alpha)$  satisfies the properties of a valid probability distribution.

- Show  $P(\beta | \alpha) \geq 0$  for all  $\beta \in S$ . By definition,

$$P(\beta | \alpha) = \frac{1}{P(\alpha)} P(\alpha \cap \beta) \quad (875)$$

and since the full joint  $P \geq 0$  and since  $P(\alpha) > 0$ , we have the desired result.

- Show  $P(\Omega_\alpha) = 1$ . Again, using just the definitions,

$$P(\Omega_\alpha) = \sum_{\beta \in S} P(\beta | \alpha) \quad (876)$$

$$= \frac{1}{P(\alpha)} \sum_{\beta \in S} P(\alpha \cap \beta) \quad (877)$$

$$= \frac{1}{P(\alpha)} \left( \sum_{\beta \in \alpha} P(\beta) + \sum_{\gamma \notin \alpha} P(\emptyset) \right) \quad (878)$$

$$= \frac{1}{P(\alpha)} (P(\alpha) + 0) \quad (879)$$

$$= 1 \quad (880)$$

- Show, for any  $\beta, \gamma \in S$ , where  $\beta \cap \gamma = \emptyset$ , that  $P(\beta \cup \gamma | \alpha) = P(\beta | \alpha) + P(\gamma | \alpha)$ .

$$P(\beta \cup \gamma | \alpha) = \frac{1}{P(\alpha)} P((\beta \cup \gamma) \cap \alpha) \quad (881)$$

$$= \frac{1}{P(\alpha)} P((\beta \cap \alpha) \cup (\gamma \cap \alpha)) \quad (882)$$

$$= \frac{1}{P(\alpha)} (P(\beta \cap \alpha) + P(\gamma \cap \alpha)) \quad (883)$$

$$= P(\beta | \alpha) + P(\gamma | \alpha) \quad (884)$$

### Exercise 2.16: Jensen's Inequality

Let  $f$  be a concave function and  $P$  a distribution over a random variable  $X$ . Then

$$\mathbb{E}_P [f(X)] \leq f(\mathbb{E}_P [X]) \quad (885)$$

Use this inequality to prove the following 3 properties.

- $H_P(X) \leq \log |Val(X)|$ . Let  $f(u) := \lg(u)$  be our concave function.

$$H_P(X) \triangleq \mathbb{E}_P \left[ \lg \frac{1}{P(X)} \right] = \mathbb{E}_P [f(u)] \quad (886)$$

$$\leq f(\mathbb{E}_P [u]) = f \left( \sum_x u(x) P(x) \right) = f(|Val(X)|) = \lg |Val(X)| \quad (887)$$

(888)

- $H_P(X) \geq 0$ .

$$-H_P(X) = -\mathbb{E}_P \left[ \lg \frac{1}{P(X)} \right] \quad (889)$$

$$= \mathbb{E}_P [\lg P(X)] \leq 0 \quad (890)$$

since  $0 \leq P(X = x) \leq 1 \forall x$  (any term in the expectation where  $P(X = x) = 0$  is equal to 0, by definition).

- $D(P||Q) \geq 0$ . Use the same idea as in the first proof, but let  $u(x) = Q(x)/P(x)$ .

$$D(P||Q) \triangleq \mathbb{E}_P [\lg (P(X)/Q(X))] = -\mathbb{E}_P [f(u)] \quad (891)$$

$$-\mathbb{E}_P [f(u)] \geq f \left( \sum_x u(x) P(x) \right) = f(1) = 0 \quad (892)$$

## The Bayesian Network Representation (Ch. 3)

Table of Contents Local

Written by Brandon McKinzie

Koller and Friedman (2009). The Bayesian Network Representation.

*Probabilistic Graphical Models: Principles and Techniques.*

**Goal:** represent a joint distribution  $P$  over some set of variables  $\mathcal{X} = \{X_1, \dots, X_n\}$ . Consider the case where each  $X_i$  is binary-valued. A single joint distribution requires access to the probability for each of the  $2^n$  possible assignments for  $\mathcal{X}$ . The set of all such possible joint distributions,

$$\{(p_1, \dots, p_{2^n}) \in \mathbb{R}^{2^n} : \sum_{i=1}^{2^n} p_i = 1\} \quad (893)$$

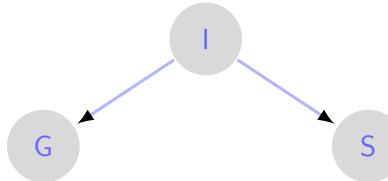
Understanding the exponential blowup.

is a  $2^n - 1$  dimensional subspace of  $\mathbb{R}^{2^n}$ . Note that each  $p_i$  represents the probability for a unique instantiation of  $\mathcal{X}$ . Furthermore, in the general case, knowing  $p_i$  tells you nearly nothing about  $p_{j \neq i}$  – i.e. you require an instantiation of  $2^n - 1$  independent parameters to specify a given joint distribution.

But it would be foolish to parameterize any joint distribution in this way, since we can often take advantage of independencies. Consider the case where each  $X_i$  gives the outcome (H or T) of coin  $i$  being tossed. Then our distribution satisfies  $(\mathbf{X} \perp \mathbf{Y})$  for any disjoint subsets of the variables  $\mathbf{X}$  and  $\mathbf{Y}$ . Let  $\theta_i$  denote the probability that coin  $i$  lands heads. The key observation is that you only need each of the  $n \theta_i$  to specify a unique joint distribution over  $\mathcal{X}$ , reducing the  $2^n - 1$  dimensional subspace to an  $n$  dimensional manifold in  $\mathbb{R}^{2^n}$ <sup>192</sup>.

Taking advantage of independencies.

**The Naive Bayes Model.** Say we want to determine the intelligence of an applicant based on their grade G in some course, and their score S on the SAT. A naive bayes model can be illustrated as below.



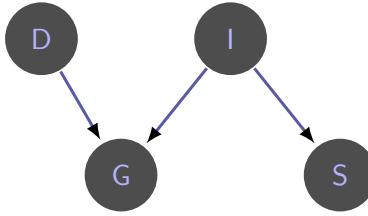
This encodes our general assumption that<sup>193</sup>  $P \models (S \perp G \mid I)$ .

<sup>192</sup>**TODO:** Formally define this manifold using the same notation as in 893. Edit: Not sure how to actually write it out, but intuitively it's because each of the  $2^n$   $p_i$  values, when going from the general case to the case of i.i.d, become *functions* of the  $n \theta_i$  values. Whereas before they were independent free parameters.

<sup>193</sup>We say that an event  $\alpha$  is independent of event  $\beta$  in  $P$  with the notation  $P \models (\alpha \perp \beta)$ . (Def 2.2, pg 23 of book)

In general, a naive bayes model assumes that instances fall into one of a number of mutually exclusive and exhaustive *classes*, defined as the set of values that the top variable in the graph can take on<sup>194</sup>. The model also includes some number of *features*  $X_1, \dots, X_k$ , whose values are typically observed. The **naive Bayes assumption** is that the features are conditionally independent given the instance's class.

**Bayesian Networks.** A Bayesian network  $\mathcal{B}$  is defined by a network structure together with its set of CPDs. **Causal reasoning** (or prediction) refers to computing the downstream effects of various factors (such as intelligence). **Evidential reasoning** (or explanation) is the reverse case, where we reason from effects to causes. Finally, **intercausal reasoning** (or explaining away) is when different causes of the same effect can interact. For our student example, we could be trying to determine  $\Pr [I | G]$ , the intelligence of an a student given his/her grade in a class. In addition to intelligence being a cause for the grade, we could have another causal variable  $D$  for the difficulty of the class:



An example of intercausal reasoning would be observing  $D$ , so that we now want  $\Pr [I | G, d]$ ; the diffulty of the course can help *explain away* good/bad grades, thus changing our value for the probability of intelligence based on the grade alone.

A **Bayesian network structure**  $\mathcal{G}$  is a DAG whose nodes represent RVs  $X_1, \dots, X_n$ . Let  $Pa_{X_i}^{\mathcal{G}}$  denote the parents of  $X_i$  in  $\mathcal{G}$ , and  $\text{NonDesc}_{X_i}$  denotes the variables that are NOT descendants of  $X_i$ . Then  $\mathcal{G}$  encodes the following set of **local independencies**,

$$\mathcal{I}_{\ell}(\mathcal{G}) \triangleq \{\forall X_i : (X_i \perp \text{NonDesc}_{X_i} | Pa_{X_i}^{\mathcal{G}})\} \quad (894)$$

---

<sup>194</sup>For the intelligence example, the classes are high intelligence and low intelligence.

**Graphs and Distributions.** Here we see that a distribution  $P$  satisfies the local independencies associated with a graph  $\mathcal{G}$  iff  $P$  is representable as a set of CPDs associated with  $\mathcal{G}$ .

- **Local independencies.** Let  $P$  be a distribution over  $\mathcal{X}$ . We define  $\mathcal{I}(P)$  to be the set of **independence assertions** of the form  $(\mathbf{X} \perp \mathbf{Y} | \mathbf{Z})$  that hold in  $P$ . The statement “ $P$  satisfies the local independencies associated with  $\mathcal{G}$ ” can thus be succinctly written:

$$\mathcal{I}_\ell(\mathcal{G}) \subseteq \mathcal{I}(P) \quad (895)$$

and we'd say that  $\mathcal{G}$  is an **I-map** (independency map) for  $P$ .

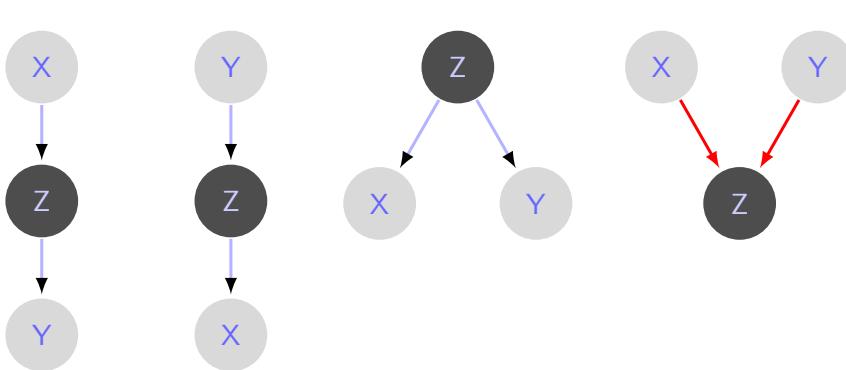
- **I-maps.** More generally, let  $\mathcal{K}$  be *any* graph object associated with a set of independencies  $\mathcal{I}(\mathcal{K})$ . We say that  $\mathcal{K}$  is an **I-map** for a set of independencies  $\mathcal{I}$  if  $\mathcal{I}(\mathcal{K}) \subseteq \mathcal{I}$ . Note that the complete graph (every two nodes connected) is an I-map for any distribution.
- **Factorization.** Let  $\mathcal{G}$  be a BN graph over  $X_1, \dots, X_n$ . We say that a distribution  $P$  over the same space **factorizes** according to  $\mathcal{G}$  if  $P$  can be expressed as

$$P(X_1, \dots, X_n) = \prod_{i=1}^n P(X_i | Pa_{X_i}^{\mathcal{G}}) \quad (896)$$

For BNs, the following equivalence holds:

$$\mathcal{G} \text{ is an I-map for } P \iff P \text{ factorizes according to } \mathcal{G}$$

**D-separation.** We want to understand when we can *guarantee* that an independence  $(\mathbf{X} \perp \mathbf{Y} | \mathbf{Z})$  holds in a distribution associated with a BN structure  $\mathcal{G}$ . Consider a 3-node network consisting of  $X$ ,  $Y$ , and  $Z$ , where  $X$  and  $Y$  are not directly connected. There are four such cases, which I've drawn below.



The fourth trail is called a *v-structure*.

From left-to-right: indirect causal effect, indirect evidential effect, common cause, common effect. The first 3 satisfy  $(X \perp Y | Z)$ , but the 4th does not. Another way of saying this is that the first 3 trails are **active**<sup>195</sup> IFF  $Z$  is not observed, while the 4th trail is active IFF  $Z$  (or a descendent of  $Z$ ) is observed.

---

<sup>195</sup>When influence can flow from  $X$  to  $Y$  via  $Z$ , we say that the trail  $X \rightleftharpoons{} Y \rightleftharpoons{} Z$  is *active*.

**General case:**

Let  $\mathcal{G}$  be a BN structure, and  $X \rightleftharpoons{\dots} X_n$  a trail in  $\mathcal{G}$ . Let  $\mathbf{Z}$  be a subset of observed variables. The trail is **active** given  $\mathbf{Z}$  if

- Any v-structure  $X_{i-1} \rightarrow X_i \leftarrow X_{i+1}$  has  $X_i$  or one of its descendants in  $\mathbf{Z}$ .
- No other node along the trail is in  $\mathbf{Z}$ .

**D-separation:** (Directed separation)

Let  $\mathbf{X}, \mathbf{Y}, \mathbf{Z}$  be three sets of nodes in  $\mathcal{G}$ . We say that  $\mathbf{X}$  and  $\mathbf{Y}$  are **d-separated** in  $\mathcal{G}$ , denoted  $d\text{-sep}_{\mathcal{G}}(\mathbf{X}; \mathbf{Y} \mid \mathbf{Z})$ , if there is no active trail between any node  $X \in \mathbf{X}$  and  $Y \in \mathbf{Y}$  given  $\mathbf{Z}$ . We use  $\mathcal{I}(\mathcal{G})$  to denote the set of independencies that correspond to d-separation,

$$\mathcal{I}(\mathcal{G}) = \{(\mathbf{X} \perp \mathbf{Y} \mid \mathbf{Z}) : d\text{-sep}_{\mathcal{G}}(\mathbf{X}; \mathbf{Y} \mid \mathbf{Z})\} \quad (897)$$

also called the set of **global Markov independencies**.

**Soundness and Completeness** of d-separation as a method for determining independence.

- **Soundness** (Thm 3.3). If a distribution  $P$  factorizes according to  $\mathcal{G}$ , then  $\mathcal{I}(\mathcal{G}) \subseteq \mathcal{I}(P)$ .
- **Completeness**. For any distribution  $P$  that factorizes over  $\mathcal{G}$ , we have that  $P$  is **faithful**<sup>196</sup> to  $\mathcal{G}$ :

$$X \perp Y \mid \mathbf{Z} \in \mathcal{I}(P) \implies d\text{-sep}_{\mathcal{G}}(X; Y \mid \mathbf{Z}) \quad (898)$$

To see the detailed algorithm for finding nodes reachable from  $X$  given  $\mathbf{Z}$  via active trails, see Algorithm 3.1 on pg. 75 of the book.

---

<sup>196</sup>  $P$  is faithful to  $\mathcal{G}$  if, whenever  $X \perp Y \mid \mathbf{Z} \in \mathcal{I}(P)$ , then  $d\text{-sep}_{\mathcal{G}}(X; Y \mid \mathbf{Z})$ .

## Undirected Graphical Models (Ch. 4)

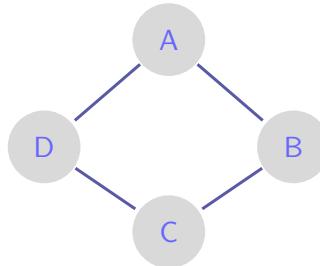
Table of Contents Local

Written by Brandon McKinzie

Koller and Friedman (2009). Undirected Graphical Models.

*Probabilistic Graphical Models: Principles and Techniques.*

**The Misconception Example.** Consider a scenario where we have four students who get together in pairs to work on homework for a class. Only the following pairs meet: (Alice, Bob), (Bob, Charles), (Charles, Debbie), (Debbie, Alice). The professor misspoke in class, giving rise to a possible misconception among the students. We have four binary random variables,  $\{A, B, C, D\}$ , representing whether the student has the misconception (1) or not (0)<sup>197</sup>. Intuitively, we want to model a distribution that satisfies  $(A \perp C | \{B, D\})$ , and  $(B \perp D | \{A, C\})$ , but no other independencies<sup>198</sup>. Note that the interactions between variables seem symmetric here – students influence each other (out of the ones they have a pair with).



The nodes in the graph of a **Markov network** represent the variables, and the edges correspond to a notion of direct probabilistic interaction between the neighboring variables – an interaction that is not mediated by any other variable in the network. So, how should we parameterize our network? We want to capture the **affinities** between the related variables (e.g. Alice and Bob are more likely to agree than disagree).

Let  $\mathbf{D}$  be a set of random variables. We define a **factor**  $\phi$  to be a function from  $\text{Val}(\mathbf{D})$  to  $\mathbb{R}$ . A factor is nonnegative if all its entries are nonnegative.  $\mathbf{D}$  is called the **scope** of the factor, denoted  $\text{Scope}[\phi]$ .

We restrict our attention to nonnegative factors.

The factors need not be normalized. Therefore, to interpret probabilities over factors, we must

<sup>197</sup> A student might not have the misconception if e.g. they went home and figured out the problem via reading the textbook instead.

<sup>198</sup> These independences cannot be naturally captured in a Bayesian (i.e. directed) network.

normalize it with what we'll call the **partition function**,  $Z$ :

$$\Pr[a, b, c, d] = \frac{1}{Z} \phi_1(a, b) \cdot \phi_2(b, c) \cdot \phi_3(c, d) \cdot \phi_4(d, a) \quad (899)$$

$$Z = \sum_{a,b,c,d} \phi_1(a, b) \cdot \phi_2(b, c) \cdot \phi_3(c, d) \cdot \phi_4(d, a) \quad (900)$$

**Parameterization.** Associating the graph structure with a set of parameters. We parameterize undirected graphs by associating a set of factors with it. First, we introduce the definition of **factor product**:

Let  $\mathbf{X}$ ,  $\mathbf{Y}$ , and  $\mathbf{Z}$  be three disjoint sets of variables, and let  $\phi_1(\mathbf{X}, \mathbf{Y})$  and  $\phi_2(\mathbf{Y}, \mathbf{Z})$  be two factors. We define the factor product  $\phi_1 \times \phi_2$  to be a factor  $\psi : \text{Val}(\mathbf{X}, \mathbf{Y}, \mathbf{Z}) \mapsto \mathbb{R}$  as follows:

$$\psi(\mathbf{X}, \mathbf{Y}, \mathbf{Z}) = \phi_1(\mathbf{X}, \mathbf{Y}) \cdot \phi_2(\mathbf{Y}, \mathbf{Z}) \quad (901)$$

We use this to define an undirected parameterization of a distribution:

A distribution  $P_\Phi$  is a **Gibbs distribution** parameterized by a set of factors  $\Phi = \{\phi_1(\mathbf{D}_1), \dots, \phi_K(\mathbf{D}_K)\}$  if it is defined as follows:

$$P_\Phi = \frac{1}{Z} \tilde{P}_\Phi(X_1, \dots, X_n) \quad (902)$$

$$\tilde{P}_\Phi(X_1, \dots, X_n) = \phi_1(\mathbf{D}_1) \times \phi_2(\mathbf{D}_2) \times \dots \times \phi_K(\mathbf{D}_K) \quad (903)$$

where the authors (Koller and Friedman) have made a point to emphasize: **A factor is only one contribution to the overall joint distribution. The distribution as a whole has to take into consideration the contributions from all of the factors involved.** Now, we relate the parameterization of a Gibbs distribution to a graph structure.

A distribution  $P_\Phi$  with  $\Phi = \{\phi_1(\mathbf{D}_1), \dots, \phi_K(\mathbf{D}_K)\}$  **factorizes** over a Markov network  $\mathcal{H}$  if each  $\mathbf{D}_k (k = 1, \dots, K)$  is a complete subgraph<sup>199</sup> of  $\mathcal{H}$ .

The factors that parameterize a Markov network are often called **clique potentials**. Although it can be used without loss of generality, the parameterization using maximal clique potentials generally obscures structure that is present in the original set of factors. Below are some useful definitions that we will use often.

---

<sup>199</sup> A subgraph is complete if every two nodes in the subgraph are connected by some edge. The set of nodes in such a subgraph is often called a **clique**. A clique  $\mathbf{X}$  is maximal if for any superset of nodes  $\mathbf{Y} \supset \mathbf{X}$ ,  $\mathbf{Y}$  is not a clique.

### Factor reduction:

Let  $\phi(\mathbf{Y})$  be a factor, and  $\mathbf{U} = \mathbf{u}$  an assignment for  $\mathbf{U} \subseteq \mathbf{Y}$ . Define the **reduction** of the factor  $\phi$  to the context  $\mathbf{U} = \mathbf{u}$ , denoted  $\phi[\mathbf{u}]$ , to be a factor over the scope  $\mathbf{Y}' = \mathbf{Y} - \mathbf{U}$ , such that

$$\phi[\mathbf{u}](\mathbf{y}') = \phi(\mathbf{y}', \mathbf{u}) \quad (904)$$

For  $\mathbf{U} \not\subseteq \mathbf{Y}$ , define  $\phi[\mathbf{u}]$  only for the assignments in  $\mathbf{u}$  to the variables in  $\mathbf{U}' = \mathbf{U} \cap \mathbf{Y}$ .

### Reduced Gibbs distribution:

Let  $P_{\Phi}(\mathbf{X})$  be a Gibbs distribution parameterized by  $\Phi = \{\phi_1, \dots, \phi_K\}$  and let  $\mathbf{u}$  be a context. The **reduced Gibbs distribution**  $P_{\Phi}[\mathbf{u}]$  is the Gibbs distribution defined by the set of factors  $\Phi[\mathbf{u}] = \{\phi_1[\mathbf{u}], \dots, \phi_K[\mathbf{u}]\}$ . More formally:

$$P_{\Phi}[\mathbf{u}] = P_{\Phi}(\mathbf{W} \mid \mathbf{u}) \quad \text{where} \quad \mathbf{W} = \mathbf{X} - \mathbf{U} \quad (905)$$

### Reduced Markov Network:

Let  $\mathcal{H}$  be a Markov network over  $\mathbf{X}$  and  $\mathbf{U} = \mathbf{u}$  a context. The **reduced Markov network**  $\mathcal{H}[\mathbf{u}]$  is a Markov network over the nodes  $\mathbf{W} = \mathbf{X} - \mathbf{U}$ , where we have an edge  $X - Y$  if there's an edge  $X - Y$  in  $\mathcal{H}$ .

Note that if a Gibbs distribution  $P_{\Phi}(\mathbf{X})$  factorizes over  $\mathcal{H}$ , then  $P_{\Phi}[\mathbf{u}]$  factorizes over  $\mathcal{H}[\mathbf{u}]$ .

**Markov Network Independencies.** A formal presentation of the undirected graph as a representation of independence assertions.

- **Active Path.** Let  $\mathcal{H}$  be a Markov network structure, and let  $X_1 - \dots - X_k$  be a path in  $\mathcal{H}$ . Let  $\mathbf{Z} \subseteq \mathcal{X}$  be a set of *observed variables*. The path is *active* given  $\mathbf{Z}$  if none of the  $X_i$ 's is in  $\mathbf{Z}$ .
- **Separation.** A set of nodes  $\mathbf{Z}$  separates  $\mathbf{X}$  and  $\mathbf{Y}$  in  $\mathcal{H}$ , denoted  $\text{sep}_{\mathcal{H}}(\mathbf{X}; \mathbf{Y} \mid \mathbf{Z})$ , if there is no active path between any node  $X \in \mathbf{X}$  and  $Y \in \mathbf{Y}$  given  $\mathbf{Z}$ .
- **Global Independencies.** The *global independencies* associated with  $\mathcal{H}$  are defined as

$$\mathcal{I}(\mathcal{H}) = \{(\mathbf{X} \perp \mathbf{Y} \mid \mathbf{Z}) : \text{sep}_{\mathcal{H}}(\mathbf{X}; \mathbf{Y} \mid \mathbf{Z})\} \quad (906)$$

This is the *separation criterion*. Note that the definition of separation is monotonic in  $\mathbf{Z}$ : if it holds for  $\mathbf{Z}$ , then it holds for any  $\mathbf{Z}' \supset \mathbf{Z}$  as well<sup>200</sup>.

---

<sup>200</sup>Which means that Markov networks are fundamentally incapable of representing nonmonotonic independence relations!

- **Soundness** of the separation criterion for detecting independence properties in distributions over  $\mathcal{H}$ . In other words, we want to prove that

$$(P \text{ factorizes over } \mathcal{H}) \implies [\text{sep}_{\mathcal{H}}(\mathbf{X}; \mathbf{Y} | \mathbf{Z}) \implies P \models (\mathbf{X} \perp \mathbf{Y} | \mathbf{Z})]$$

where the portion in brackets can equivalently be said as “ $\mathcal{H}$  is an I-map for  $P$ ”.

### Proof

- Consider the case where  $\mathbf{X} \cup \mathbf{Y} \cup \mathbf{Z} = \mathcal{X}$ .
- Then, any clique in  $\mathcal{X}$  is fully contained in either  $\mathbf{X} \cup \mathbf{Z}$  or  $\mathbf{Y} \cup \mathbf{Z}$ . In other words,

$$P(\mathbf{X}, \mathbf{Y}, \mathbf{Z}) = \frac{1}{Z} f(\mathbf{X}, \mathbf{Z}) g(\mathbf{Y}, \mathbf{Z}) \quad (907)$$

$$\text{which implies } P(\mathbf{X}, \mathbf{Y} | \mathbf{Z}) = \frac{f(\mathbf{X}, \mathbf{Z}) g(\mathbf{Y}, \mathbf{Z})}{\sum_{x,y} f(x, \mathbf{Z}) g(y, \mathbf{Z})} \quad (908)$$

- We can prove that  $P \models (\mathbf{X} \perp \mathbf{Y} | \mathbf{Z})$  by showing that  $P(\mathbf{X}, \mathbf{Y} | \mathbf{Z}) = P(\mathbf{X} | \mathbf{Z}) P(\mathbf{Y} | \mathbf{Z})$ .

$$P(\mathbf{X}, \mathbf{Y} | \mathbf{Z}) = \frac{f(\mathbf{X}, \mathbf{Z}) g(\mathbf{Y}, \mathbf{Z})}{P(\mathbf{Z})} \quad (909)$$

$$= \frac{f(\mathbf{X}, \mathbf{Z}) g(\mathbf{Y}, \mathbf{Z})}{P(\mathbf{Z})} \frac{P(\mathbf{Z})}{P(\mathbf{Z})} \quad (910)$$

$$= \frac{f(\mathbf{X}, \mathbf{Z}) g(\mathbf{Y}, \mathbf{Z})}{P(\mathbf{Z})} \frac{\sum_x f(x, \mathbf{Z}) \sum_y g(y, \mathbf{Z})}{P(\mathbf{Z})} \quad (911)$$

$$= \frac{f(\mathbf{X}, \mathbf{Z})}{P(\mathbf{Z})} \frac{\sum_y g(y, \mathbf{Z})}{P(\mathbf{Z})} \frac{g(\mathbf{Y}, \mathbf{Z}) \sum_x f(x, \mathbf{Z})}{P(\mathbf{Z})} \quad (912)$$

$$= P(\mathbf{X} | \mathbf{Z}) P(\mathbf{Y} | \mathbf{Z}) \quad (913)$$

- For the general case where  $\mathbf{X} \cup \mathbf{Y} \cup \mathbf{Z} = \mathcal{X}$ , let  $\mathbf{U} = \mathcal{X} - (\mathbf{X} \cup \mathbf{Y} \cup \mathbf{Z})$ . Since we know that  $\text{sep}_{\mathcal{H}}(\mathbf{X}; \mathbf{Y} | \mathbf{Z})$ , we can partition  $\mathbf{U}$  into two disjoint sets  $\mathbf{U}_1$  and  $\mathbf{U}_2$  such that  $\text{sep}_{\mathcal{H}}(\mathbf{X} \cup \mathbf{U}_1; \mathbf{Y} \cup \mathbf{U}_2 | \mathbf{Z})$ . Combining the previous result with the decomposition property<sup>a</sup> give us the desired result that  $P \models (\mathbf{X} \perp \mathbf{Y} | \mathbf{Z})$ .

---

<sup>a</sup>The decomposition property:

$$(\mathbf{X} \perp (\mathbf{Y}, \mathbf{W}) | \mathbf{Z}) \implies (\mathbf{X} \perp \mathbf{Y} | \mathbf{Z})$$

### Pairwise Independencies:

Let  $\mathcal{H}$  be a Markov network. We define the **pairwise independencies** assoc. with  $\mathcal{H}$ :

$$\mathcal{I}_p(\mathcal{H}) = \{(X \perp Y | \mathcal{X} - \{X, Y\}) : X - Y \notin \mathcal{H}\} \quad (914)$$

which just says “ $X$  is indep. of  $Y$  given everything else if there’s no edge between  $X$  and  $Y$ .”

### Local Independencies:

For a given graph  $\mathcal{H}$ , define the **Markov blanket** of  $X$  in  $\mathcal{H}$ , denoted  $\text{MB}_{\mathcal{H}}(X)$ , to be the neighbors of  $X$  in  $\mathcal{H}$ . We define the **local independencies** associated with  $\mathcal{H}$ :

$$\mathcal{I}_{\ell}(\mathcal{H}) = \{(X \perp \mathcal{X} - \{X\} - \text{MB}_{\mathcal{H}}(X) \mid \text{MB}_{\mathcal{H}}(X)) : X \in \mathcal{X}\} \quad (915)$$

which just says “ $X$  is indep. of the rest of the nodes in the graph given its immediate neighbors.”

**Log-Linear Models.** Certain patterns involving particular values of variables for a given factor can often be more easily seen by converting factors into log-space. More precisely, we can rewrite a factor  $\phi(\mathbf{D})$  as

$$\phi(\mathbf{D}) = e^{-\epsilon(\mathbf{D})} \quad (916)$$

$$\epsilon(\mathbf{D}) \triangleq -\ln \phi(\mathbf{D}) \quad (917)$$

$$\Pr [X_1, \dots, X_n] \propto e^{-\sum \epsilon_i(\mathbf{D}_i)} \quad (918)$$

where  $\epsilon(\mathbf{D})$  is often called an **energy function**. Note how  $\epsilon$  can take on any value along the real line (i.e. removes our nonnegativity constraint)<sup>201</sup>. Also note that as the  $\epsilon$  summation approaches 0, the probability approaches one.

This motivates introducing the notion of a **feature**, which is just a factor without the nonnegativity requirement. A popular type of feature is the **indicator feature** that takes on value 1 for some values  $\mathbf{y} \in \text{Val}(\mathbf{D})$  and 0 otherwise. We can now provide a more general definition for our notion of log-linear models:

A distribution  $P$  is a **log-linear model** over a Markov network  $\mathcal{H}$  if it is associated with:

- A set of features  $\mathcal{F} = \{f_1(\mathbf{D}_1), \dots, f_k(\mathbf{D}_k)\}$  where each  $\mathbf{D}_i$  is a complete subgraph (i.e. a clique) in  $\mathcal{H}$ .
- A set of weights  $w_1, \dots, w_k$ .

such that

$$\Pr [X_1, \dots, X_n] = \frac{1}{Z} \exp \left[ - \sum_{i=1}^k w_i f_i(\mathbf{D}_i) \right] \quad (919)$$

The log-linear model provides a much more compact representation for many distributions, especially in situations where variables have large domains (such as text).

---

<sup>201</sup>We seem to be implicitly assuming that the original factors are all *positive* (not just non-negative).

### Box 4.C – Concept: Ising Models and Boltzmann Machines

The **Ising model**: Each atom is modeled as a binary RV  $X_i \in \{+1, -1\}$  denoting its spin. Each pair of neighboring atoms is associated with energy function  $\epsilon_{i,j}(x_i, x_j) = w_{i,j}x_i x_j$ . We also have individual energy functions  $u_i x_i$  for each atom. This defines our distribution:

$$P(\xi) = \frac{1}{Z} \exp \left( - \sum_{i < j} w_{i,j} x_i x_j - \sum_i u_i x_i \right) \quad (920)$$

The **Boltzmann distribution**: Now the variables are  $X_i \in \{0, 1\}$ . The distribution of each  $X_i$  given its neighbors is

$$P(x_i^1 | Nb(X_i)) = \text{sigmoid}(z) \quad (921)$$

$$z = -\left( \sum_j w_{i,j} x_j \right) - u_i \quad (922)$$

### Box 4.D – Concept: Metric MRFs

Consider the pairwise graph  $X_1, \dots, X_n$  in the context of sequence labeling. We want to assign each  $X_i$  a label. We also want adjacent nodes to prefer being similar to each other. We usually use the MAP objective, so our goal will be to minimize the total energy over the parameters (which are given by the individual energy functions  $\epsilon_i$ ).

$$E(x_1, \dots, x_n) = \sum_i \epsilon_i(x_i) + \sum_{(i,j) \in \mathcal{E}} \epsilon_{i,j}(x_i, x_j) \quad (923)$$

The simplest place to start for preferring neighboring labels to take on similar values is to define  $\epsilon_{i,j}$  to have low energy when  $x_i = x_j$  and some positive  $\lambda_{i,j}$  otherwise. We want to have finer granularity for our similarities between labels. To do this, we introduce the definition of a **metric**: a function  $\mu : \mathcal{V} \times \mathcal{V} \mapsto [0, \infty)$  that satisfies

$$[\text{reflexivity}] \quad \mu(v_k, v_l) = 0 \quad \text{IFF} \quad k = l \quad (924)$$

$$[\text{symmetry}] \quad \mu(v_k, v_l) = \mu(v_l, v_k) \quad (925)$$

$$[\text{triangle inequality}] \quad \mu(v_k, v_l) + \mu(v_l, v_m) \geq \mu(v_k, v_m) \quad (926)$$

and we can now let  $\epsilon_{i,j}(v_k, v_l) := \mu(v_k, v_l)$ .

**Canonical Parameterization** (4.4.2.1). Markov networks are generally overparameterized<sup>202</sup>. The **canonical parameterization**, which requires that  $P$  be positive, avoids this. First, some notation and requirements:

- $P$  must be positive.
- Let  $\xi^* = (x_1^*, \dots, x_n^*)$  denote some fixed assignment to the network variables  $\mathcal{X}$ .
- Define  $\mathbf{x}_Z \triangleq \mathbf{x}(Z)$  as the assignment of variables in some subset  $Z$ .<sup>203</sup>
- Define  $\xi_{-Z}^* \triangleq \xi^*(\mathcal{X} - Z)$  be our fixed assignment for the variables outside  $Z$ .
- Let  $\ell(\xi)$  denote  $\ln P(\xi)$ .

<sup>202</sup>Meaning: for any  $P_\Phi$ , there are often infinitely many ways to choose its set of parameter values for a given  $\mathcal{H}$ .

<sup>203</sup>And  $\mathbf{x}$  is some assignment to some subset of  $\mathcal{X}$  that also contains  $Z$ .

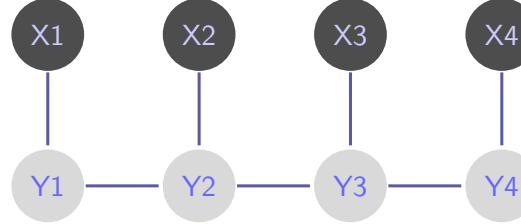
The **canonical energy function** for a clique  $\mathbf{D}$  is defined below, as well as the associated total  $P(\xi)$  over a full network assignment:

$$\epsilon_{\mathbf{D}}^*(\mathbf{d}) = \sum_{\mathbf{Z} \subseteq \mathbf{D}} \ell(\mathbf{d}_{\mathbf{Z}}, \xi_{-\mathbf{Z}}^*) \cdot (-1)^{|\mathbf{D} - \mathbf{Z}|} \quad (927)$$

$$P(\xi) = \exp \left[ \sum_i \epsilon_{\mathbf{D}_i}^*(\xi \langle \mathbf{D}_i \rangle) \right] \quad (928)$$

The sum is over all subsets of  $\mathbf{D}$ , including  $\mathbf{D}$  itself and  $\emptyset$ .

**Conditional Random Fields.** So far, we've only described Markov network representation as encoding a joint distribution over  $\mathcal{X}$ . The same undirected graph representation and parameterization can also be used to encode a *conditional distribution*  $\Pr[\mathbf{Y} | \mathbf{X}]$ , where  $\mathbf{Y}$  is a set of *target variables* and  $\mathbf{X}$  is a (disjoint) set of *observed variables*.



Note how there are no connection between any of the  $X$ s

Formally, a CRF is an undirected graph  $\mathcal{H}$  whose nodes correspond to  $\mathbf{Y} \cup \mathbf{X}$ . Since we want to avoid representing<sup>204</sup> a probabilistic model over  $\mathbf{X}$ , we disallow potentials that involve only variables in  $\mathbf{X}$ ; our set of factors is  $\phi_1(\mathbf{D}_1), \dots, \phi_m(\mathbf{D}_m)$ , such that each  $\mathbf{D}_i \not\subseteq \mathbf{X}$ . The network encodes a conditional distribution as follows:

$$P(\mathbf{Y} | \mathbf{X}) = \frac{1}{Z(\mathbf{X})} \tilde{P}(\mathbf{Y}, \mathbf{X}) \quad (929)$$

$$\tilde{P}(\mathbf{Y}, \mathbf{X}) = \prod_{i=1}^m \phi_i(\mathbf{D}_i) \quad (930)$$

$$Z(\mathbf{X}) = \sum_{\mathbf{Y}} \tilde{P}(\mathbf{Y}, \mathbf{X}) \quad (931)$$

where now the partition function is a function of the assignment  $\mathbf{x}$  to  $\mathbf{X}$ .

---

<sup>204</sup>Also note how we never have to deal with a summation over all possible  $\mathbf{X}$ , due to restricting ourselves to  $Z(\mathbf{X})$ .

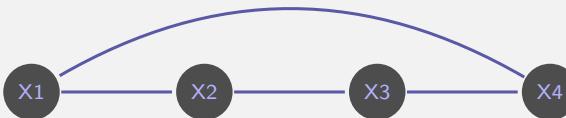
## Rapid Summary.

- **Gibbs distribution:** any probability distribution that can be written as a product of factors divided by some partition function  $Z$ .
- **Factorizes:** A Gibbs distribution factorizes over  $\mathcal{H}$  if each factor [in its product of factors] is a clique.

### 7.3.1 EXERCISES

#### Exercise 4.1

Let  $\mathcal{H}$  be the graph of binary variables below. Show that  $P$  does not factorize over  $\mathcal{H}$  (Hint: proof by contradiction).



- Example 4.4 recap:  $P$  satisfies the global independencies w.r.t  $\mathcal{H}$ . They showed this by manually checking the two global indeps of  $\mathcal{H}$ ,  $(X_1 \perp X_3 | X_2, X_4)$  and  $(X_2 \perp X_4 | X_1, X_3)$ , against the tabulated list of possible assignments for  $P$  (given in example). Nothing fancy.
- $P$  factorizes over  $\mathcal{H}$  if it can be written as a product of clique potentials.
- My proof:
  1. Assume that  $P$  does factorize over  $\mathcal{H}$ .
  2. Then  $P$  can be written as

$$P(X_1, X_2, X_3, X_4) = \frac{1}{Z} \phi_1(X_1, X_2) \phi_2(X_2, X_3) \phi_3(X_3, X_4) \phi_4(X_4, X_1) \quad (932)$$

Let the above assertion be denoted as  $C$ , and the statement that  $P$  factorizes according to  $\mathcal{H}$  be denoted simply as  $P_{\mathcal{H}}$ . Since  $P_{\mathcal{H}} \iff C$ , if we can prove that  $C$  does not hold, then we've found our contradiction, meaning  $P_{\mathcal{H}}$  also must not hold.

3. I know that the proof must take advantage of the fact that we know  $P$  is zero for certain assignments to  $\mathcal{X}$ . For example  $P(0100) = 0$ . Furthermore, by looking at the assignments where  $P$  is *not* zero, I can see that all possible combinations of  $(X_1, X_2)$  are present, which means  $\phi_1$  never evaluates to zero.
4. From the example, we know that  $P(1100) = 1/8 \neq 0$ . However, since

$$0 = \frac{P(0100)}{P(1100)} = \frac{\phi_1(0, 1)\phi_4(0, 0)}{\phi_1(1, 1)\phi_4(0, 1)} \quad (933)$$

$$= \frac{\phi_4(0, 0)}{\phi_4(0, 1)} \quad (\phi_1 > 0) \quad (934)$$

and we also know that both the numerator and denominator of eq. 934 are positive, and thus we have a contradiction.

### Exercise 4.4

Prove theorem 4.7 for the case where  $\mathcal{H}$  consists of a single clique. Theorem 4.7 is equation 928 in my notes. For a single clique  $D$ , the question reduces to: Show that, for any assignment  $d$  to  $D$ :

$$P(d) = \frac{\exp(-\epsilon(d))}{\sum_{d'} \exp(-\epsilon(d'))} \quad (935)$$

$$= \exp(\epsilon_D^*(d)) \quad (936)$$

Consider the case where  $|D| = 1$ , i.e.  $d = d$  is a single variable. Then

$$\epsilon_D^*(d) = (-1)^{|D|} \ell(\xi_D^*) + (-1)^{|D-D|} \ell(d) \quad (937)$$

$$= -\ell(\xi_D^*) + \ell(d) \quad (938)$$

$$= -\ln P(d^*) + \ln P(d) \quad (939)$$

and therefore

$$\exp(\epsilon_D^*(d)) = P(d)/P(d^*) \quad (940)$$

which is clearly incorrect (???) **TODO:** figure out what's going on here. Either the book has a type in its for theorem 4.7, or I'm absolutely insane.

## Local Probabilistic Models (Ch. 5)

Table of Contents Local

Written by Brandon McKinzie

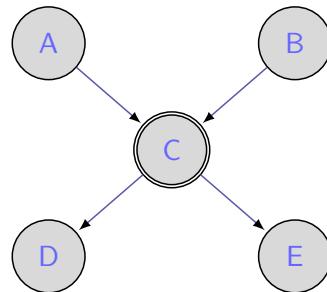
Koller and Friedman (2009). Local Probabilistic Models.

*Probabilistic Graphical Models: Principles and Techniques.*

**Deterministic CPDs.** When  $X$  is a deterministic function of its parents  $Pa_X$ :

$$P(x | Pa_X) = \begin{cases} 1 & x = f(Pa_X) \\ 0 & \text{otherwise.} \end{cases} \quad (941)$$

Consider the example below, where the double-line notation on C means that C is a deterministic function of A and B. What new conditional dependencies do we have?



Answer:  $(D \perp E | A, B)$ , which would not be true by d-separation alone. It only holds because C is a deterministic function of A and B. z

## Template-Based Representations (Ch. 6)

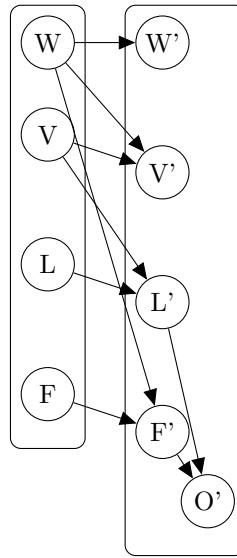
Table of Contents Local

Written by Brandon McKinzie

Koller and Friedman (2009). Template-Based Representations.

*Probabilistic Graphical Models: Principles and Techniques.*

In what follows, we will build on an example where a vehicle tries to track its true location ( $L$ ) using various sensor readings: velocity ( $V$ ), weather ( $W$ ), failure of sensor ( $F$ ), observed location ( $O'$ ) of the noisy sensor.



**Temporal Models.** We discretize time into slices of interval  $\Delta$ , and denote the ground random variables at time  $t \cdot \Delta$  by  $\mathcal{X}^{(t)}$ . We can simplify our formulation considerably by assuming a **Markovian system**: a dynamic system over template variables  $\mathcal{X}$  that satisfies the Markov assumption:

$$(\mathcal{X}^{(t+1)} \perp \mathcal{X}^{(0:(t-1))} \mid \mathcal{X}^{(t)}) \quad (942)$$

which allows us to define a more compact representation of the joint distribution from time 0 to T:

$$P(\mathcal{X}^{(0:T)}) = P(\mathcal{X}^{(0)}) \prod_{t=0}^{T-1} P(\mathcal{X}^{(t+1)} \mid \mathcal{X}^{(t)}) \quad (943)$$

One last simplifying assumption, to avoid having unique transition probabilities for each time  $t$ , is to assume a **stationary**<sup>205</sup> Markovian dynamic system, defined s.t.  $P(\mathcal{X}^{(t+1)} \mid \mathcal{X}^{(t)})$  is the same for all  $t$ .

---

<sup>205</sup>Also called time invariant or homogeneous.

**Dynamic Bayesian Networks** (6.2.2). Above, I've drawn the **2-time-slice Bayesian network** (2-TBN) for our location example. A 2-TBN is a conditional BN over  $\mathcal{X}'$  given  $\mathcal{X}_I$ , where  $\mathcal{X}_I \subseteq \mathcal{X}$  is a set of **interface variables**<sup>206</sup>. For each template variable  $X_i$ , the CPD  $P(X'_i | Pa_{X'_i})$  is a **template factor**. We can use the notion of the 2-TBN to define the more general **dynamic Bayesian network**:

A **dynamic Bayesian network** (DBN) is a pair  $\langle \beta_0, \beta_{\rightarrow} \rangle$ , where  $\beta_0$  is a Bayesian network over  $\mathcal{X}^{(0)}$ , representing the initial distribution over states, and  $\beta_{\rightarrow}$  is a 2-TBN for the process. For any  $T \geq 0$ , the unrolled Bayesian network is defined such that

- $p(X_i^{(0)} | Pa_{X_i^{(0)}})$  is the same as the CPD for the corresponding  $X_i$  in  $\beta_0$ .
- $p(X_i^{(t)} | Pa_{X_i^{(t)}})$  (for  $t > 0$ ) is the same as the CPD for the corresponding  $X_i'$  in  $\beta_{\rightarrow}$ .

**State-Observation Models** (6.2.3). Temporal models that, in addition to the Markov assumption (eq. 942), model the observation variables at time  $t$  as conditionally independent of the entire state sequence given the variables at time  $t$ :

$$(\mathbf{O}^{(t)} \perp \mathbf{X}^{(0:(t-1))}, \mathbf{X}^{(t+1:\infty)} | \mathbf{X}^{(t)}) \quad (944)$$

So basically a 2-TBN with the constraint that observation variables are leaves and only have parents in  $\mathbf{X}'$ . We now view our probabilistic model as consisting of 2 components: the *transition model*  $P(\mathbf{X}' | \mathbf{X})$ , and the *observation model*  $P(\mathbf{O} | \mathbf{X})$ . The two main architectures for such models are as follows:

- **Hidden Markov Models.** Defined as having a single state variable  $S$  and a single observation variable  $O$ . In practice, the transition model  $P(S' | S)$  is often assumed to be sparse (many possible transitions having zero probability). In such cases, one usually represents them visually as **probabilistic finite-state automaton**<sup>207</sup>.
- **Linear Dynamical Systems** (LDS) represent a system of one or more real-valued variables that evolve linearly over time, with some Gaussian noise. Such systems are often called **Kalman filters**, after the algorithm used to perform tracking. They can be viewed as a DBN with continuous variables and all dependencies are linear Gaussian<sup>208</sup>. A LDS is traditionally represented as a state-observation model, where both state and observation are vector-valued RVs, and the transition/observation models are encoded using matrices. More formally, for  $\mathbf{X}^{(t)} \in \mathbb{R}^n, O \in \mathbb{R}^m$ :

$$P(\mathbf{X}^{(t)} | \mathbf{X}^{(t-1)}) = \mathcal{N}(A\mathbf{X}^{(t-1)}; Q) \quad (945)$$

$$P(O^{(t)} | \mathbf{X}^{(t)}) = \mathcal{N}(H\mathbf{X}^{(t)}; R) \quad (946)$$

$$Q \in \mathbb{R}^{n \times n}$$

$$H \in \mathbb{R}^{m \times n}$$

---

<sup>206</sup>Interface variables are those variables whose values at time  $t$  can have a direct effect on the variables at time  $t + 1$ .

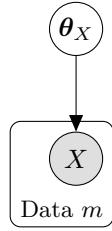
<sup>207</sup>FSA use the graphical notation where the nodes are the individual possible values in  $Val(S)$ , and the directed edges from some  $a$  to  $b$  have weight equal to  $P(S' = b | S = a)$ .

<sup>208</sup>Linear Gaussian: some var  $Z$  pointing to  $X$  denotes that  $X = \Lambda Z + \text{noise}$ , where noise  $\sim \mathcal{N}(\mu_x, \Sigma_x)$ .

**Template Variables and Template Factors** (6.3). It's convenient to view the world as being composed of a set of **objects**, which can be divided into a set of mutually exclusive and exhaustive *classes*  $\mathcal{Q} = Q_1, \dots, Q_k$ . Template **attributes** have a tuple of *arguments*, each of which is associated with a particular class of objects, which defines the set of objects that can be used to instantiate the argument in a given domain. Template attributes thus provide us with a “generator” for RVs in a given probability space. Formally,

An **attribute**  $A$  is a function  $A(U_1, \dots, U_k)$ , whose range is some set  $Val(A)$ , and where each argument  $U_i$  is a typed **logical variable** associated with a particular class  $Q[U_i]$ . The tuple  $U_1, \dots, U_k$  is called the **argument signature** of the attribute  $A$ , and denoted  $\alpha(A)$ .

**Plate Models** (6.4.1). The simplest example of a plate model is shown below. It describes multiple RVs generated from the same distribution  $\mathcal{D}$ .



This could be a plate model for a set of coin tosses sampled from a single coin. We have a set of  $m$  random variables  $X(d)$ , where  $d \in \mathcal{D}$ . Each  $X(d)$  is the random variable for the  $d$ th coin toss. We also explicitly model that the single coin for which the tosses are used is sampled from a distribution  $\theta_X$ , which takes on values  $[0, 1]$  and denotes the bias of the coin.

## Gaussian Network Models (Ch. 7)

Table of Contents Local

Written by Brandon McKinzie

Koller and Friedman (2009). Gaussian Network Models.

*Probabilistic Graphical Models: Principles and Techniques.*

**Multivariate Gaussians.** Here I'll give two forms of the familiar density function, followed by some comments and terminology.

$$p(\mathbf{x}) = \frac{1}{(2\pi)^{n/2}\sqrt{\det \Sigma}} \exp \left\{ -\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^T \Sigma^{-1} (\mathbf{x} - \boldsymbol{\mu}) \right\} \quad (947)$$

$$p(\mathbf{x}) \propto \exp \left\{ -\frac{1}{2} \mathbf{x}^T J \mathbf{x} + (\mathbf{J}\boldsymbol{\mu})^T \mathbf{x} \right\} \quad \text{where } J \triangleq \Sigma^{-1} \quad (948)$$

- The **standard Gaussian** is defined as  $\mathcal{N}(\mathbf{0}, \mathbf{I})$ .
- $\Sigma$  must be *positive definite*<sup>209</sup>:  $\forall \mathbf{x} \neq 0, \mathbf{x}^T \Sigma \mathbf{x} > 0$ . Recall that  $\Sigma_{i,j} = \text{Cov}[x_i, x_j] = \mathbb{E}[x_i x_j] - \mu_i \mu_j$ .

The two operations we usually want to perform on a Gaussian are (1) computing marginals, and (2) conditioning the distribution on an assignment of some subset of the variables. For (1), its easier to use the standard form of  $p(\mathbf{x})$ , whereas for (2) it is easier to use the information form (the one using  $J$ ).

Multivariate Gaussians are also special because we can easily determine whether two  $x_i$  and  $x_j$  are independent:  $x_i \perp x_j$  IFF  $\Sigma_{i,j} = 0$ <sup>210</sup>. For conditional independencies, the information matrix  $J$  is easier to work with:  $(x_i \perp x_j \mid \{x\}_{k \notin \{i,j\}})$  IFF  $J_{i,j} = 0$ . This condition is also how we defined pairwise independencies in a Markov network, which leads to the awesome realization:

*We can view the information matrix  $J$  as directly defining a minimal I-map Markov network for [multivariate Gaussian]  $p$ , whereby any entry  $J_{i,j} \neq 0$  corresponds to an edge  $x_i - x_j$  in the network.*

---

<sup>209</sup>Most definitions of p.d. also require that the matrix be symmetric. I also think it helps to view positive definite from the operator perspective: A linear operator  $T$  is positive definite if  $T$  is self-adjoint and  $\langle T(x), x \rangle > 0$ . In other words, “positive definite” means that the result of applying the matrix/operator to any nonzero  $\mathbf{x}$  will always have a positive component along the original direction  $\hat{\mathbf{x}}$ .

<sup>210</sup>This is not true in general – just for multivariate Gaussians!

## Variable Elimination (Ch. 9)

Table of Contents Local

Written by Brandon McKinzie

Koller and Friedman (2009). Variable Elimination.

*Probabilistic Graphical Models: Principles and Techniques.*

**Analysis of Exact Inference.** The focus of this chapter is the conditional probability query,

$$\Pr [\mathbf{Y} | \mathbf{E} = \mathbf{e}] = \frac{\Pr [\mathbf{Y}, \mathbf{e}]}{\Pr [\mathbf{e}]} \quad (9.1)$$

Ideally, we want to obtain all instantiations  $\Pr [\mathbf{y} | \mathbf{e}]$  of equation 9.1. Let  $\mathbf{W} = \mathcal{X} - \mathbf{Y} - \mathbf{E}$  be the RVs that are neither query nor evidence. Then

$$\Pr [\mathbf{y} | \mathbf{e}] = \frac{\sum_{\mathbf{w}} \Pr [\mathbf{y}, \mathbf{e}, \mathbf{w}]}{\sum_{\mathbf{y}} \Pr [\mathbf{y}, \mathbf{e}]} \quad (9.3)$$

and note that, by computing all instantiations for the numerator first, we can reuse them to obtain the denominator. We can formulate the inference problem as a decision problem, which we will call *BNPrDP*, defined as follows:

*Given:* Bayesian network  $\mathcal{B}$  over  $\mathcal{X}$ , a variable  $X \in \mathcal{X}$ , and a value  $x \in \text{Val}(X)$ .

*Decide:* whether  $\Pr_{\mathcal{B}} [X = x] > 0$ .

*BNPrDP*

**Thm 9.1:** The decision problem *BNPrDP* is  $\mathcal{NP}$ -complete. **Proof:**

1. **BNPrDP is in  $\mathcal{NP}$ :** Guess a full assignment  $\xi$  to the network<sup>211</sup>. If the guess is successful, where success is defined as  $(X = x) \in \xi$  and  $P(\xi) > 0$ , then we know that  $P(X = x) > 0$ <sup>212</sup>. Computing  $P(\xi)$  is linear in the number of factors for a BN, since we just multiply them together.
2. **BNPrDP is  $\mathcal{NP}$ -hard.** We show this by proving that we can solve 3-SAT (which is  $\mathcal{NP}$ -hard) by transforming inputs to 3-SAT to inputs of BNPrDP in polynomial time. Given any 3-SAT formula  $\phi$ , we can create a BN  $B_\phi$  with some special variable  $X$  s.t.  $\phi$  is satisfiable IFF  $P_{B_\phi}(X = x^1) > 0$ . You can easily build such a network by having a node  $Q_i$  for each binary RV  $q_i$ , and a node  $C_i$  for each of the clauses that's a deterministic function of its parents (up to 3 Q nodes). Then, the node X is a deterministic function of its parents, which are chains of AND gates along the  $C_i$ . Since each node has at most 3 parents, we can ensure that construction is bounded by polynomial time in the length of  $\phi$ .

<sup>211</sup> Apparently the time it takes to generate a guess is irrelevant.

<sup>212</sup>This is true because  $P(x)$  can be decomposed as  $P(\xi) + \sum P(\dots) \geq P(\xi)$ .

**Analysis of Approximate Inference.** Consider a specific query  $P(\mathbf{y} \mid \mathbf{e})$ , where we focus on a particular assignment  $\mathbf{y}$ . Let  $\rho$  denote some approximate answer, whose accuracy we wish to evaluate relative to the correct probability. We can use the **relative error** to estimate the quality of the approximation: *An estimate  $\rho$  has relative error  $\epsilon$  if:*

$$\frac{\rho}{1 + \epsilon} \leq P(\mathbf{y} \mid \mathbf{e}) \leq \rho(1 + \epsilon) \quad (949)$$

Unfortunately, the task of finding some approximation  $\rho$  with relative error  $\epsilon$  is also  $\mathcal{NP}$ -hard. Furthermore, even if we relax this metric by using absolute error instead, we end up finding that **in the case where we have evidence, approximate inference is no easier than exact inference, in the worst case.**

**Variable Elimination.** Basically, a dynamic programming approach for performing exact inference. Consider the simple BN  $X_1 \rightarrow \dots \rightarrow X_n$ , where each variable can take on  $k$  possible values. The dynamic programming approach for computing  $P(X_n)$  involves computing

$$P(X_{i+1}) = \sum_{x_i} P(X_{i+1} \mid x_i) P(x_i) \quad (950)$$

$n - 1$  times, starting with  $i = 1$ , all the way up to  $i = n - 1$ , reusing the previous computation at each step, with total cost  $\mathcal{O}(nk^2)$ . So for this simple network, even though the size of the joint is  $k^n$  (exponential in  $n$ ), we can do inference in linear time.

First, we formalize some basic concepts before defining the algorithm.

**Factor marginalization:**

Let  $\mathbf{X}$  be a set of variables, and  $Y \notin \mathbf{X}$  a variable. Let  $\phi(\mathbf{X}, Y)$  be a factor. Define the **factor marginalization** of  $Y$  in  $\phi$ , denoted  $\sum_Y \phi$ , to be a factor  $\psi$  over  $\mathbf{X}$  such that:

$$\psi(\mathbf{X}) = \sum_Y \phi(\mathbf{X}, Y)$$

The key observation that's easy to miss is that *we're only summing entries in the table where the values of  $\mathbf{X}$  match up*. One useful rule for exchanging factor product and summation: If  $X \notin \text{Scope}[\phi_1]$ , then

$$\sum_X (\phi_1 \cdot \phi_2) = \phi_1 \cdot \sum_X \phi_2 \quad (9.6)$$

So, when computing some marginal probability, the main idea is to group factors together and compute expressions of the form

$$\sum_{\mathbf{Z}} \prod_{\phi \in \Phi} \phi \quad (951)$$

where  $\Phi$  is the set of all factors  $\phi$  for which  $\mathbf{Z} \in \text{Scope}[\phi]$ . This is commonly called the **sum-product** inference task. The full algorithm for sum-product variable elimination, which is an instantiation of the sum-product inference task, is illustrated below.

```

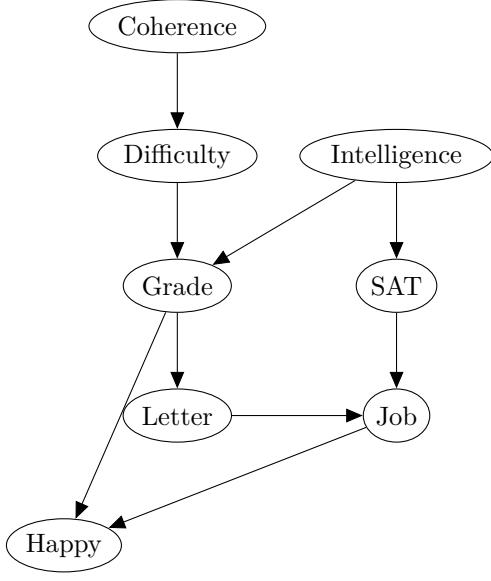
Procedure Sum-Product-VE (
     $\Phi$ , // Set of factors
     $\mathbf{Z}$ , // Set of variables to be eliminated
     $\prec$  // Ordering on  $\mathbf{Z}$ 
)
    Let  $Z_1, \dots, Z_k$  be an ordering of  $\mathbf{Z}$  such that
     $Z_i \prec Z_j$  if and only if  $i < j$ 
    for  $i = 1, \dots, k$ 
         $\Phi \leftarrow \text{Sum-Product-Eliminate-Var}(\Phi, Z_i)$ 
         $\phi^* \leftarrow \prod_{\phi \in \Phi} \phi$ 
    return  $\phi^*$ 

Procedure Sum-Product-Eliminate-Var (
     $\Phi$ , // Set of factors
     $Z$  // Variable to be eliminated
)
     $\Phi' \leftarrow \{\phi \in \Phi : Z \in \text{Scope}[\phi]\}$ 
     $\Phi'' \leftarrow \Phi - \Phi'$ 
     $\psi \leftarrow \prod_{\phi \in \Phi'} \phi$ 
     $\tau \leftarrow \sum_{\mathbf{Z}} \psi$ 
    return  $\Phi'' \cup \{\tau\}$ 

```

This is what we use to compute the marginal probability  $P(\mathbf{X})$  where  $\mathbf{X} = \mathcal{X} - \mathbf{Z}$ . To compute conditional queries of the form  $P(\mathbf{Y} | \mathbf{E} = \mathbf{e})$ , simply replace all factors whose scope overlaps with  $\mathbf{E}$  with their reduced factor (see chapter 4 notes for definition) to get the unnormalized  $\phi^*(\mathbf{Y})$  (the numerator of  $P(\mathbf{Y} | \mathbf{e})$ ). Then divide by  $\sum_{\mathbf{y}} \phi^*$  to obtain the final result.

**Example.** We will work through computing  $P(\text{Job})$  for the BN below.



Due to Happy being a child of Job,  $P(J)$  actually requires using all factors in the graph. Below shows how VE with elimination ordering  $C, D, I, H, G, S, L$  progressively simplifies the equation for computing  $P(J)$ .

$$P(J) = \sum P(J | s, \ell) P(s | i) P(\ell | g) P(g | d, i) P(d | c) P(h | J, g) P(c) P(i) \quad (952)$$

$$= \sum_{c,d,i,h,g,s,\ell} (\textcolor{red}{P(c)P(d | c)}) \cdot P(g | d, i) P(i) P(s | i) P(h | J, g) P(\ell | g) P(J | s, \ell) \quad (953)$$

$$= \sum_{\textcolor{red}{d},i,h,g,s,\ell} (\tau_1(d) P(g | d, i)) \cdot P(i) P(s | i) P(h | J, g) P(\ell | g) P(J | s, \ell) \quad (954)$$

$$= \sum_{i,h,g,s,\ell} (\tau_2(g, i) P(i) P(s | i)) \cdot P(h | J, g) P(\ell | g) P(J | s, \ell) \quad (955)$$

$$= \sum_{\textcolor{red}{h},g,s,\ell} (P(h | J, g)) \cdot \tau_3(g, s) P(\ell | g) P(J | s, \ell) \quad (956)$$

$$= \sum_{\textcolor{red}{g},s,\ell} (\tau_4(g, J) \tau_3(g, s) P(\ell | g)) \cdot P(J | s, \ell) \quad (957)$$

$$= \sum_{\textcolor{red}{s},\ell} \tau_5(J, \ell, s) \cdot P(J | s, \ell) \quad (958)$$

$$= \sum_{\ell} \tau_6(J, \ell) \quad (959)$$

where red indicates the focus of the given step in the VE algorithm.

## Clique Trees (Ch. 10)

Table of Contents Local

Written by Brandon McKinzie

Koller and Friedman (2009). Clique Trees.

*Probabilistic Graphical Models: Principles and Techniques.*

**Cluster Graphs.** A graphical flowchart of the factor-manipulation process that will be relevant when we discuss message passing. Each node is a *cluster*, which is associated with a subset of variables. Formally,

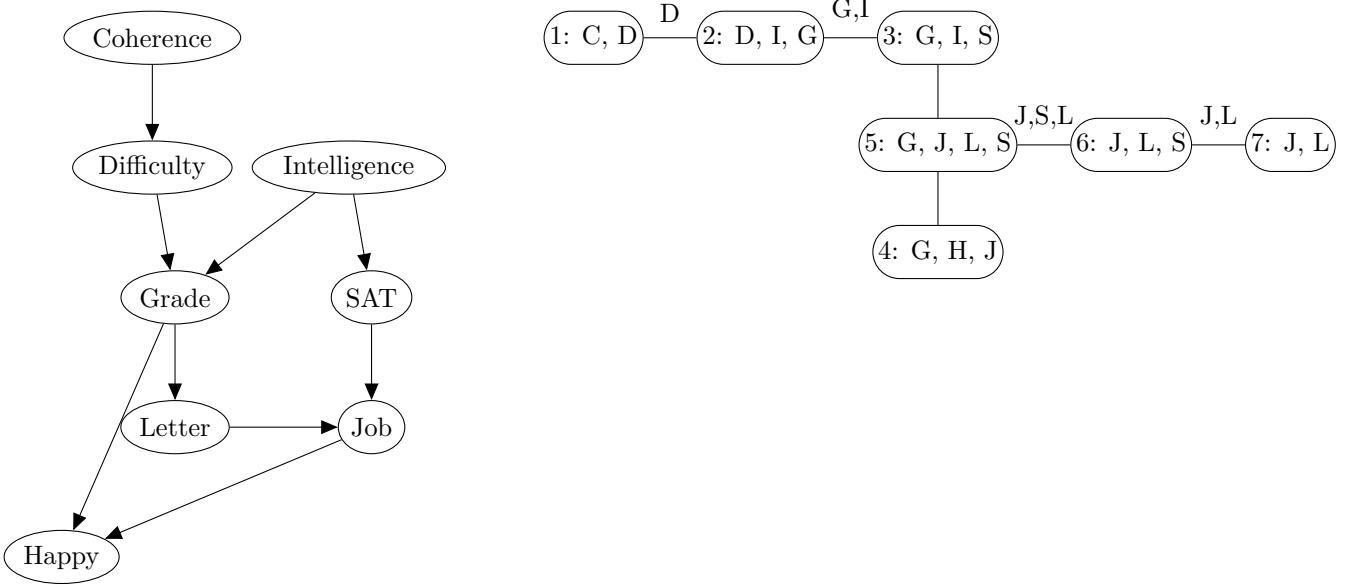
A **cluster graph**  $\mathcal{U}$  for a set of factors  $\Phi$  over  $\mathcal{X}$  is an undirected graph. Each node  $i$  is associated with a subset  $\mathbf{C}_i \subseteq \mathcal{X}$ . Each factor  $\phi \in \Phi$  must be associated with a cluster  $\mathbf{C}_i$ , denoted  $\alpha(\phi)$ , such that  $\text{Scope}[\phi] \subseteq \mathbf{C}_i$  (**family-preserving**). Each edge between a pair of clusters  $\mathbf{C}_i$  and  $\mathbf{C}_j$  is associated with a **sepset**  $S_{i,j} \subseteq \mathbf{C}_i \cap \mathbf{C}_j$ .

Recall that each step of variable elimination involves creating a factor  $\psi_i$  by multiplying a group of factors<sup>213</sup>. Then, denoting the variable we are eliminating at this step as  $Z$ , we obtain another factor  $\tau_i$  that's the factor marginalization of  $Z$  in  $\psi_i$  (denoted  $\sum_Z \psi_i$ ). An execution of variable elimination defines a cluster graph: we have a cluster for each of the  $\psi_i$ , defined as  $\mathbf{C}_i = \text{Scope}[\psi_i]$ . We draw an edge between  $\mathbf{C}_i$  and  $\mathbf{C}_j$  if the **message**  $\tau_i$  is used in the computation of  $\tau_j$ .

Consider when we applied variable elimination to the student graph network below, to compute  $P(J)$ . Elimination ordering  $C, D, I, H, G, S, L$ .

---

<sup>213</sup>All whose scope contains the variable we are currently trying to eliminate.



**Clique Trees.** Since VE uses each intermediate  $\tau_i$  at most once, the cluster graph induced by an execution VE is necessarily a *tree*, and it also defines a directionality: the direction of the message passing (left-to-right in the above illustration). All the messages flow toward a single cluster where the final result is computed – the **root** of the tree; we say the messages “flow up” to the root. Furthermore, for cluster trees induced by VE, the scope of each message (edge)  $\tau_i$  is *exactly*  $C_i \cap C_j$ , not just a subset<sup>214</sup>.

Let  $\Phi$  be a set of factors over  $\mathcal{X}$ . A cluster tree over  $\Phi$  that satisfies the running intersection property is called a **clique tree**. For clique trees, the clusters are also called **cliques**.

**Message Passing: Sum Product.** Previously we saw how an execution of VE can be illustrated with a clique tree. We now go the other direction – given a clique tree, we show how it can be used for variable elimination. Given a clique tree representation of some BN, we can use it to guide us along an execution of VE to compute any marginal we’d like. First, before any run, we generate the set of **initial potentials**  $\psi_i$  associated with each clique  $C_i$  in the tree, defined as just the multiplication of the initial factors associated with the clique. We define the root of the tree as any clique containing the variable whose marginal we want to compute (we pick arbitrarily). Starting from the leaves and moving toward the root, we pass messages along from clique to clique. A clique is *ready* to send a message when it has received a message from all of its downstream neighbors. The message from  $C_i$  to [a neighbor]  $C_j$  is computed using the **sum-product message passing** computation:

---

<sup>214</sup>This follows from the **running intersection property**, which is satisfied by any cluster tree that’s defined by variable elimination. It’s defined as, if any variable  $X$  is in both cluster  $C_i$  and  $C_j$ , then  $X$  is also in every cluster in the (unique) path in the tree between  $C_i$  and  $C_j$ .

$$\delta_{i \rightarrow j} = \sum_{C_i - S_{i,j}} \psi_i \cdot \prod_{k \in (Nb_i - \{j\})} \delta_{k \rightarrow i} \quad (10.2)$$

where the summation is simply over the variables in  $C_i$  that aren't passed along to  $C_j$ , and the product is over all messages that  $C_i$  received. Stated even simpler, we multiply all the incoming messages by our initial potential, then sum out all variables except those in  $S_{i,j}$ . When the root clique has received all messages, it multiplies them with its own initial potential, resulting in a factor called the **beliefs**,  $\beta_r(C_r)$ . It represents

$$\tilde{P}_\Phi(C_r) = \sum_{\mathcal{X} - C_r} \prod_{\phi} \phi \quad (960)$$

where, to be clear, the product is over all  $\phi$  in the graph.

Below is a more compact summary of all of this, showing the procedure for computing *all* final factors (belief)  $\beta_i$  for some marginal probability query on the variables in  $C_r$  *asynchronously*.

### Algorithm 10.2: Sum-Product Belief Propagation

1. For each clique  $C_i$ , compute its initial potential:

$$\psi_i(C_i) \leftarrow \prod_{\phi_j : \alpha(\phi_j) = i} [\phi_j]$$

2. While  $\exists i, j$  such that  $i$  is ready to transmit to  $j$ , compute:

$$\delta_{i \rightarrow j} \leftarrow \sum_{C_i - S_{i,j}} \psi_i \cdot \prod_{k \in (Nb_i - \{j\})} \delta_{k \rightarrow i}$$

3. Then, compute each belief factors  $\beta_i$  by multiplying the initial potential  $\psi_i$  by the incoming messages to  $C_i$ :

$$\beta_i \leftarrow \psi_i \cdot \prod_{k \in Nb_{C_i}} \delta_{k \rightarrow i}$$

4. Return the set of beliefs  $\{\beta_i\}$ , where

$$\beta_i = \sum_{\mathcal{X} - C_i} \tilde{P}_\Phi(\mathcal{X}) = \tilde{P}_\Phi(C_i) \quad (961)$$

The SP Belief Propagation algorithm above is also called **clique tree calibration**. A clique tree  $\mathcal{T}$  is *calibrated* if all pairs of adjacent cliques are calibrated. A calibrated clique tree satisfies the following property for what we'll now call the **clique beliefs**,  $\beta_i$ , and the **sepset beliefs**,  $\mu_{i,j}$  over  $S_{i,j}$ :

$$\mu_{i,j}(S_{i,j}) \triangleq \sum_{C_i - S_{i,j}} \beta_i = \sum_{C_j - S_{i,j}} \beta_j \quad (962)$$

$$\mu_{i,j} = \tilde{P}_\Phi(S_{i,j})$$

*The main advantage of the clique tree algorithm is that it computes the posterior probability of all variables in a graphical model using only twice the computation<sup>215</sup> of the upward pass in the same tree.*

---

<sup>215</sup>The algorithm is equivalent to doing one upward pass, one downward pass.

We can also show that  $\mu_{i,j} = \delta_{j \rightarrow i} \delta_{i \rightarrow j}$ , which then allows us to derive:

$$\tilde{P}_\Phi(\mathcal{X}) = \frac{\prod_{i \in \mathcal{V}_T} \beta_i}{\prod_{ij \in \mathcal{E}_T} \mu_{i,j}} \quad (10.10)$$

In other words, the clique and sepset beliefs provide a **reparameterization** of the unnormalized measure, a property called the **clique tree invariant**.

**Message Passing: Belief Update.** We now discuss an alternative message passing approach that is mathematically equivalent but intuitively different than the sum-product approach. First, we introduce some new definitions.

#### Factor Division:

Let  $\mathbf{X}$  and  $\mathbf{Y}$  be disjoint sets of variables, and let  $\phi_1(\mathbf{X}, \mathbf{Y})$  and  $\phi_2(\mathbf{Y})$  be two factors.

We define the division of  $\phi_1$  and  $\phi_2$  as a factor  $\psi$  with scope  $\mathbf{X}, \mathbf{Y}$  as follows:

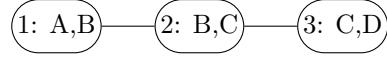
Define  $0/0 = 0$

$$\psi(\mathbf{X}, \mathbf{Y}) \triangleq \frac{\phi_1(\mathbf{X}, \mathbf{Y})}{\phi_2(\mathbf{Y})}$$

Looking back at equation 10.2, we can now see that another way to write  $\delta_{i \rightarrow j}$  is

$$\delta_{i \rightarrow j} = \frac{\sum_{C_i - S_{i,j}} \beta_i}{\delta_{j \rightarrow i}} \quad (10.13)$$

Now, consider the clique tree below for the simple Markov network A-B-C-D:



If we assigned  $C_2$  as the root, then our previous approach would compute  $\delta_{2 \rightarrow 1}$  as  $\sum_C \psi_2 \cdot \delta_{3 \rightarrow 2}$ . Alternatively, we can use equation 10.13 to realize this is equivalent to dividing  $\beta_2$  by  $\delta_{1 \rightarrow 2}$  and marginalizing out  $C$ . This observation motivates the algorithm below, which allows us to execute message passing in terms of the clique and sepset beliefs, without having to remember the initial potentials  $\psi_i$  or explicitly compute the messages  $\delta_{i \rightarrow j}$ .

#### Algorithm 10.3: Belief-Update Message Passing

1. For each clique  $C_i$ , set its initial belief  $\beta_i$  to its initial potential  $\psi_i$ . For each edge in  $\mathcal{E}_T$ , set  $\mu_{i,j} = 1$ .
2. While there exists an uninformed<sup>216</sup> clique in  $\mathcal{T}$ , select any edge in  $\mathcal{E}_T$ , and compute

---

<sup>216</sup>A clique is informed once it has received informed messages from all of its neighbors. An informed message is one that has been sent by taking into account information from all of the sending cliques' neighbors (aside from the receiving clique of that message, of course).

$$\sigma_{i \rightarrow j} \leftarrow \sum_{C_i - S_{i,j}} [\beta_i] \quad (963)$$

$$\beta_j \leftarrow \beta_j \cdot \frac{\sigma_{i \rightarrow j}}{\mu_{i,j}} \quad (964)$$

$$\mu_{i,j} \leftarrow \sigma_{i \rightarrow j} \quad (965)$$

3. Return the resulting set of informed beliefs  $\{\beta_i\}$ .

At convergence,  $\sigma_{i \rightarrow j} = \mu_{i,j} = \sigma_{j \rightarrow i}$ .

## Inference as Optimization (Ch. 11)

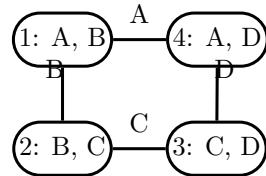
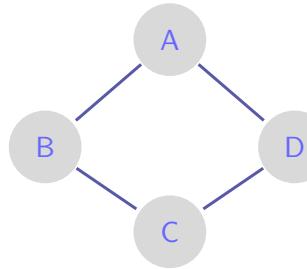
Table of Contents Local

Written by Brandon McKinzie

Koller and Friedman (2009). Inference as Optimization.

*Probabilistic Graphical Models: Principles and Techniques.*

**Propagation-Based Approximation.** We can use a general-purpose cluster graph rather than the more restrictive clique tree (needed to guarantee exact inference) for approximate inference methods. Consider the simple Markov network below on the left.



The clique tree for this network, which can be used for exact inference, has two cliques ABD and BCD and messages are passed between them consisting of  $\tau(B, D)$ . Suppose that, instead, we set up 4 clusters corresponding to each of the initial potentials, shown as the cluster graph above on the right. We can still apply belief propagation here, but due to it now having loops (as opposed to before when we only had trees), *the process may not converge*.

## Parameter Estimation (Ch. 17)

Table of Contents Local

Written by Brandon McKinzie

Koller and Friedman (2009). Parameter Estimation.

*Probabilistic Graphical Models: Principles and Techniques.*

**Maximum Likelihood Estimation** (17.1). In this chapter, assume the network structure is fixed and that our data set  $\mathcal{D}$  consists of fully observed instances of the network variables:  $\mathcal{D} = \{\xi[1], \dots, \xi[M]\}$ . We begin with the simplest learning problem: parameter learning for a single variable. We want to estimate the probability, denoted via the *parameter*  $\theta$ , with which the flip of a thumbtack will land heads or tails. Define the **likelihood function**  $L(\theta : \mathbf{x})$  as the probability of observing some sequence of outcomes  $\mathbf{x}$  under the parameter  $\theta$ . In other words, it is simply  $P(\mathbf{x} : \theta)$ , but interpreted as a *function of  $\theta$* . For our simple case, where  $\mathcal{D}$  consists of  $M$  thumbtack flip outcomes,

$$L(\theta : \mathcal{D}) = \theta^{M[1]}(1 - \theta)^{M[0]} \quad (966)$$

where  $M[1]$  denotes the number of outcomes in  $\mathcal{D}$  that were heads. Since it's easier to maximize a logarithm, and since it yields the same optimal  $\hat{\theta}$ , optimize the **log-likelihood** to obtain:

$$\hat{\theta} = \arg \max_{\theta} \ell(\theta : \mathcal{D}) = \arg \max_{\theta} [M[1] \log \theta + M[0] \log(1 - \theta)] \quad (967)$$

$$= \frac{M[1]}{M[1] + M[0]} \quad (968)$$

Note that MLE has the *disadvantage* that it can't communicate confidence of an estimate<sup>217</sup>.

We now provide the more general formal definitions for MLE.

- We are given a *training set*  $\mathcal{D}$  containing  $M$  (IID) instances of a set of random variables  $\mathcal{X}$ , where the samples of  $\mathcal{X}$  are drawn from some unknown distribution  $P^*(\mathcal{X})$ .
- We are given a **parametric model**, defined by a function  $P(\xi; \boldsymbol{\theta})$ , where  $\xi$  is an instance of  $\mathcal{X}$ , and we want to estimate its parameters  $\boldsymbol{\theta}$ <sup>218</sup>. The model also defines the space of legal parameter values  $\Theta$ , the **parameter space**.
- We then define the **likelihood function**  $L(\boldsymbol{\theta} : \mathcal{D}) = \prod_m P(\xi[m] : \boldsymbol{\theta})$ .

<sup>217</sup>We get the same result (0.3) if we get 3 heads out of 10 flips, as we do for getting 300 heads out of 1000 flips; yet, the latter experiment should include a higher degree of confidence

<sup>218</sup>We also have the constraint that  $P(\xi; \boldsymbol{\theta})$  must be a valid distribution (nonnegative and sums to 1 over all possible  $\xi$ )

We can often simplify the likelihood function to simpler terms, like our  $M[0]$  and  $M[1]$  values in the thumbtack example. These are called the **sufficient statistics**, defined as functions of the data that summarize the relevant information for computing the likelihood. Formally,

*A function  $\tau(\xi) : \xi \rightarrow \mathbb{R}^\ell$  (for some  $\ell$ ) is a **sufficient statistic** if for any two data sets  $\mathcal{D}$  and  $\mathcal{D}'$ , we have that*

$$\left[ \sum_{\xi[m] \in \mathcal{D}} \tau(\xi[m]) = \sum_{\xi'[m] \in \mathcal{D}'} \tau(\xi'[m]) \right] \implies \left[ L(\boldsymbol{\theta} : \mathcal{D}) = L(\boldsymbol{\theta} : \mathcal{D}') \right] \quad (969)$$

*We often informally refer to the tuple  $\sum_{\xi[m] \in \mathcal{D}} \tau(\xi[m])$  as the **sufficient statistics** of the data set  $\mathcal{D}$ .*

**MLE for Bayesian Networks – Simple Example.** We now move on to estimating parameters  $\boldsymbol{\theta}$  for the simple BN  $X \rightarrow Y$  for two binary RVs  $X$  and  $Y$ . Our parameters  $\boldsymbol{\theta}$  are the individual probabilities of  $P(X)$  and  $P(Y | X)$  (6 total). Since BNs have the nice property that their joint probability decomposes into a product of probabilities, just like how the likelihood function is a product of probabilities, we can write the likelihood function as a product of the individual *local* probabilities:

$$L(\boldsymbol{\theta} : \mathcal{D}) = \left( \prod_m P(x[m] : \boldsymbol{\theta}_X) \right) \left( \prod_m P(y[m] | x[m] : \boldsymbol{\theta}_{Y|X}) \right) \quad (970) \quad \text{decomposability of the likelihood function}$$

which can be decomposed even further by e.g. differentiating products over  $x[m] : x[m] = x^0$  etc. Just as we used  $M[0]$  in the thumbtack example to count the number of instances with a certain value, we can use the same idea for the general case.

*Let  $Z$  be some set of RVs, and  $z$  be some instantiation to them. We define  $M[z]$  to be the number of entries in data set  $\mathcal{D}$  that have  $Z[m] = z$ :*

$$M[z] = \sum_m \mathbb{1}\{Z[m] = z\} \quad (971)$$

**Global Likelihood Decomposition.** We now move to the more general case of computing the likelihood for BN with structure  $\mathcal{G}$ .

$$L(\boldsymbol{\theta} : \mathcal{D}) = \prod_m P_{\mathcal{G}}(\xi[m] : \boldsymbol{\theta}) \quad (972) \quad \text{global decomposition of the likelihood}$$

$$= \prod_m \prod_i P(x_i[m] | Pa_{X_i}[m] : \boldsymbol{\theta}) \quad (973)$$

$$= \prod_i \left[ \prod_m P(x_i[m] | Pa_{X_i}[m] : \boldsymbol{\theta}_{X_i|Pa_{X_i}}) \right] \quad (974)$$

$$= \prod_i L_i(\boldsymbol{\theta}_{X_i|Pa_{X_i}} : \mathcal{D}) \quad (975)$$

where  $L_i$  is the **local likelihood function** for  $X_i$ . Assuming these are each disjoint sets of parameters from one another, it implies that  $\hat{\boldsymbol{\theta}} = \langle \hat{\boldsymbol{\theta}}_{X_1|Pa_{X_1}}, \dots, \hat{\boldsymbol{\theta}}_{X_n|Pa_{X_n}} \rangle$

## Partially Observed Data (Ch. 19)

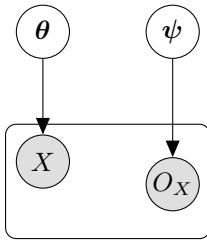
Table of Contents Local

Written by Brandon McKinzie

Koller and Friedman (2009). Partially Observed Data.

*Probabilistic Graphical Models: Principles and Techniques.*

**Likelihood of Data and Observation Models** (19.1.1). Consider the simple example of flipping a thumbtack, but occasionally the thumbtack rolls off the table. We choose to ignore the tosses for which the thumbtack rolls off. Now, in addition to the random variable  $X$  giving the flip outcome, we have the *observation variable*  $O_X$ , which tells us whether we observed the value of  $X$ .



The illustration above is a plate model where we choose a thumbtack sampled with bias  $\theta$  and repeat some number of flips with that same thumbtack. We also sample the random variable  $O_X$  that has probability of observation sampled from  $\psi$  and fixed for all the experiments we do. This leads to the following definition for the observability model.

Let  $\mathbf{X} = \{X_1, \dots, X_n\}$  be some set of RVs, and let  $O_{\mathbf{X}} = \{O_{X_1}, \dots, O_{X_n}\}$  be their **observability variable**. The **observability model** is a joint distribution

$$P_{\text{missing}}(\mathbf{X}, O_{\mathbf{X}}) = P(\mathbf{X}) \cdot P_{\text{missing}}(O_{\mathbf{X}} \mid \mathbf{X})$$

so that  $P(\mathbf{X})$  is parameterized by  $\theta$  and  $P_{\text{missing}}(O_{\mathbf{X}} \mid \mathbf{X})$  is parameterized by  $\psi$ . We define a new set of RVs  $\mathbf{Y} = \{Y_1, \dots, Y_n\}$  where  $\text{Val}(Y_i) = \text{Val}(X_i) \cup \{?\}$ . The actual observation  $\mathbf{Y}$  is a deterministic function of  $\mathbf{X}$  and  $O_{\mathbf{X}}$ :

$$Y_i = \begin{cases} X_i & O_{X_i} = o^1 \\ ? & O_{X_i} = o^0 \end{cases} \quad (976)$$

For our simple model above, we have

$$P(Y = 1) = \theta\psi \quad (977)$$

$$P(Y = 0) = (1 - \theta)\psi \quad (978)$$

$$P(Y = ?) = (1 - \psi) \quad (979)$$

$$L(\theta, \psi; \mathcal{D}) = \theta^{M[1]} (1 - \theta)^{M[0]} \psi^{M[1] + M[0]} (1 - \psi)^{M[?]} \quad (980)$$

The main takeaway is to understand that **when we have missing data, the data-generation process involves two steps:** (1) generate data by sampling from the model, then (2) determine which values we get to observe and which ones are hidden from us.

**The Likelihood Function** (19.1.3). Assume we have a BN network  $\mathcal{G}$  over a set of variables  $\mathbf{X}$ . In general, each instance has a different set of observed variables. Denote by  $\mathbf{O}[m]$  and  $\mathbf{o}[m]$  the observed vars and their values in the  $m$ 'th instance, and by  $\mathbf{H}[m]$  the missing (or hidden) vars in the  $m$ 'th instance.

# INFORMATION THEORY, INFERENCE, AND LEARNING ALGORITHMS

## CONTENTS

8.1	Introduction to Information Theory (Ch. 1) . . . . .	325
8.2	Probability, Entropy, and Inference (Ch. 2) . . . . .	327
8.2.1	More About Inference (Ch. 3 Summary) . . . . .	329
8.3	The Source Coding Theorem (Ch. 4) . . . . .	331
8.3.1	Data Compression and Typicality . . . . .	333
8.3.2	Further Analysis and Q&A . . . . .	335
8.4	Monte Carlo Methods (Ch. 29) . . . . .	337
8.5	Variational Methods (Ch. 33) . . . . .	339

## Introduction to Information Theory (Ch. 1)

Table of Contents Local

*Written by Brandon McKinzie*

[Note: Skipping most of this chapter since it's mostly introductory material.]

**Preface.** For ease of reference, some common quantities we will frequently be using:

- **Binomial distribution.** Let  $r$  denote the number of successful trials out of  $N$  total trials. Let  $f$  denote the probability of success for a single trial.

$$\Pr[r \mid f, N] = \binom{N}{r} f^r (1-f)^{N-r} \quad \mathbb{E}[r] = Nf \quad \text{Var}[r] = Nf(1-f) \quad (981)$$

- **Stirling's Approximation.**

$$x! \simeq x^x e^{-x} \sqrt{2\pi x} \Leftrightarrow \ln x! \simeq x \ln x - x + \frac{1}{2} \ln 2\pi x \quad (982) \quad \text{Recall that } \log_b x = \frac{\log_a x}{\log_a b}$$

$$\ln \binom{N}{r} \simeq r \ln \frac{N}{r} + (N-r) \ln \frac{N}{N-r} \quad (983)$$

- **Binary Entropy Function** and its relation with Stirling's approximation.

$$H_2(x) \triangleq x \lg \frac{1}{x} + (1-x) \lg \frac{1}{1-x} \quad (984)$$

$$\lg \binom{N}{r} \simeq NH_2(r/N) \quad (985)$$

**Perfect communication over an imperfect, noisy communication channel (1.1).** We want to make an encoder-decoder architecture, of the general form in the figure below, to achieve reliable communication over a noisy channel.

**Information theory** is concerned with the theoretical limitations and potentials of such systems. Let's explore some examples for the case of the **binary symmetric channel**<sup>219</sup>:

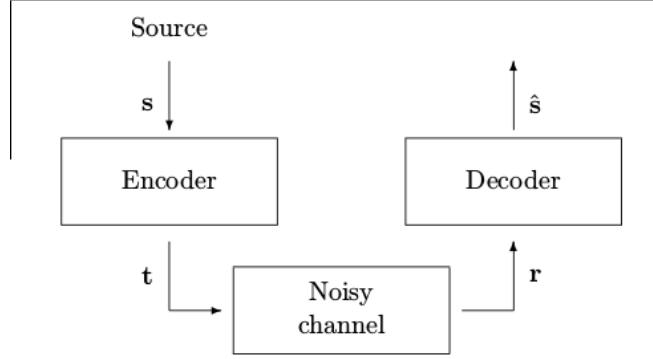
- **Repetition codes.** Let  $R_N$  denote the repetition code that repeats each bit in the message  $N$  times<sup>220</sup>. We model the channel as “adding<sup>221</sup>” a sparse noise vector  $\mathbf{n}$  to the encoded message  $\mathbf{t}$ , so  $\mathbf{r} := \mathbf{n} + \mathbf{t}$ . What is the optimal way of decoding  $\mathbf{r}$ ?

$$\hat{s}_i \leftarrow \arg \max_{s_i} \Pr[s_i \mid \mathbf{r}_{i:i+n}] = \arg \max_{s_i} \Pr[\mathbf{r}_{i:i+N} \mid s_i] \Pr[s_i] \quad (986)$$

<sup>219</sup>A binary symmetric channel transmits each bit correctly with probability  $(1 - f)$  and incorrectly with probability  $f$ , where  $f$  is assumed to be small.

<sup>220</sup>So  $R_2$  would encode 101 as 110011.

<sup>221</sup>We add in modulo 2, which is NOT the same as binary arithmetic (no carry). Addition modulo 2 is the same as doing XOR.

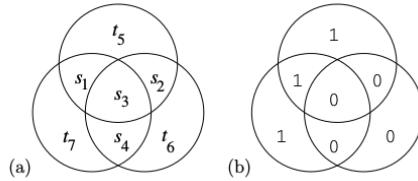


We see that we must make assumptions about the prior probability  $\Pr[s_i]$ . It is common to assume all possible values of  $s_i$  (0 or 1 in this case) are equally probable. It is useful to observe the **likelihood ratio**,

$$\frac{\Pr[r_{i:i+N} | s_i = 1]}{\Pr[r_{i:i+N} | s_i = 0]} = \prod_{n=i}^{i+N-1} \frac{\Pr[r_n | t_n = 1]}{\Pr[r_n | t_n = 0]} = \prod_{n=i}^{i+N-1} \begin{cases} \gamma & \text{if } r_n = 1 \\ \gamma^{-1} & \text{if } r_n = 0 \end{cases} \quad (987)$$

where we've defined  $\gamma := (1 - f)/f$ , with  $f$  being the probability of a bit getting flipped by the channel. We want to assign  $\hat{s}_i$  to the most likely **hypothesis** out of the possible  $s_i$ . If the likelihood ratio [for the two hypotheses] is greater than 1, we choose  $\hat{s}_i = 1$ , else we choose  $\hat{s}_i = 0$ .

- **Block Codes - the (7, 4) Hamming Code.** Although, by increasing the number of repetitions per bit  $N$  for our  $R_N$  repetition code can decrease the error-per-bit probability  $p_b$ , we incur a substantial decrease in the *rate of information transfer* – a factor of  $1/N$ . The **(7, 4) Hamming Code** tries to improve this by encoding *blocks* of bits at a time (instead of per-bit). It is a *linear block code* – it encodes each 4-bit block into a 7-bit block, where the additional 3 bits are linear functions of the original  $K = 4$  bits. The (7, 4) Hamming code has each of the extra 3 bits act as parity checks. It's easiest to show this with the illustration below:



which encodes 1000 as 1000101.

## Probability, Entropy, and Inference (Ch. 2)

Table of Contents Local

Written by Brandon McKinzie

[Note: Skipping most of this chapter since it's mostly introductory material.]

**Notation.** Some of the notation this author seems to use a lot.

- **Ensemble  $X$ :** a triple  $(x, \mathcal{A}_X, \mathcal{P}_X)$ , where the outcome  $x$  is the value of a R.V. which can take on one of a set of possible values (“alphabet”),  $\mathcal{A}_X = \{a_1, \dots, a_i, \dots, a_I\}$ , having probabilities  $\mathcal{P}_X = \{p_1, \dots, p_i, \dots, p_I\}$ . Note that this doesn’t appear technically consistent with how the author actually *uses* the term ensemble – in practice, he actually means “ $X$  is the set of all possible triples  $(x, \mathcal{A}_X, \mathcal{P}_X), \forall x \in \mathcal{A}_X$ , or something like that. He uses it to casually refer to the space of possibilities.

**Forward Probabilities and Inverse Probabilities.** Both of these involve a **generative model** of the data. In a *forward* probability problem, we want find the PDF, expectation, or some other function of a quantity that depends on the data/is *produced* by the generative process. For example, we can model a series of  $N$  coin flips as “producing” the quantity  $n_H$ , denoting the number of heads. In an *inverse* probability problem, we want to compute conditional probabilities on one or more of the *unobserved variables* in the process, *given* the observed variables.

**The Likelihood Principle.** For a generative model on data  $d$  given parameters  $\theta$ ,  $\Pr[d | \theta]$ , and having observed a particular outcome  $d_1$ , all inferences and predictions should depend only on the function  $\Pr[d_1 | \theta]$ .

### Entropy and Related Quantities.

- **Shannon information content** of an outcome  $x$ :

$$h(x) \triangleq \lg \left( \frac{1}{P(x)} \right) \quad (988)$$

They mention the example of the unigram probabilities for each character in a document. For example,  $p(z) = 0.007$ , which has information content of 10.4 bits<sup>222</sup>.

---

<sup>222</sup>Intuition digression: Recall from CS how to compute the number of bits required to represent  $N$  unique values (answer:  $\lg(N)$ ). Similarly, a probability of e.g 1/8 can be interpreted as “one of 8 possible outcomes”, meaning that  $\lg(1/(1/8))=3$  bits are needed to encode all possible outcomes. Similarly, one could interpret  $p(z)=0.007$  as “belonging to 7 of 1000 possible results”. I guess in some strange world you can then say that there are  $1000/7 \approx 142.86$  evenly-proportioned events like this (how do you even word this) and it would take

- **Entropy of Ensemble  $X$ .** Defined to be the average Shannon information content of an outcome.

$$H(X) \triangleq \sum_{x \in \mathcal{A}_X} \Pr[x] \lg \left( \frac{1}{\Pr[x]} \right) \quad 0 \times \lg \left( \frac{1}{0} \right) \triangleq 0 \quad (989)$$

$$H(X) \leq \lg(|\mathcal{A}_X|) \quad \text{with equality iff } p_i = \frac{1}{|\mathcal{A}_X|} \forall i \quad (990)$$

- **Decomposability of the Entropy.** For any probability distribution  $\mathbf{p} = \{p_1, p_2, \dots, p_I\}$  and  $m$  (where  $1 \leq m \leq I$ ):

$$H(\mathbf{p}) = H(\Sigma_{1:m}, \Sigma_{m+1:I}) \quad (991)$$

$$+ \Sigma_{1:m} H\left(\frac{p_1}{\Sigma_{1:m}}, \dots, \frac{p_m}{\Sigma_{1:m}}\right) \quad (992)$$

$$+ \Sigma_{m+1:I} H\left(\frac{p_{m+1}}{\Sigma_{m+1:I}}, \dots, \frac{p_I}{\Sigma_{m+1:I}}\right) \quad (993)$$

where I've let  $\Sigma_{1:m} := \sum_{i=1}^m p_i$ .

- **Kullback-Leibler Divergence** between two probability distributions  $P(X)$  and  $Q(X)$  that are defined over the same alphabet  $\mathcal{A}_X$ :

$$D_{KL}(P||Q) = \sum_x P(x) \lg \frac{P(x)}{Q(x)} \quad (994)$$

$$D_{KL}(P||Q) \geq 0 \quad [\text{Gibb's Inequality}] \quad (995)$$

where, in the words of the author, “**Gibb’s inequality is probably the most important inequality in this book**”.

- **Convex functions and Jensen’s Inequality.** A function  $f(x)$  is convex over the interval  $[x = a, x = b]$  if every chord of the function lies above the function. That is,  $\forall x_1, x_2 \in [a, b]$  and  $0 \leq \lambda \leq 1$ :

$$f(\lambda x_1 + (1 - \lambda)x_2) \leq \lambda f(x_1) + (1 - \lambda)f(x_2) \quad (996)$$

$$f(\mathbb{E}[x]) \leq \mathbb{E}[f(x)] \quad [\text{Jensen's Inequality}] \quad (997)$$

and we say  $f$  is **strictly convex** if,  $\forall x_1, x_2 \in [a, b]$ , we get equality for only  $\lambda = 0$  and  $\lambda = 1$ .

$\log(142.86)=10.4$  bits to encode all of them. Low-probability events (such as a character being  $z$ ) have high information content.

Perhaps a better way to think of this is explained on the wikipedia page:

When the content of a message is known a priori with certainty, with probability of 1, there is no actual information conveyed in the message. Only when the advance knowledge of the content of the message by the receiver is less than 10% certain does the message actually convey information. Accordingly, the amount of self-information contained in a message conveying content informing an occurrence of event,  $\omega_n$ , depends only on the probability of that event.

---

### 8.2.1 MORE ABOUT INFERENCE (CH. 3 SUMMARY)

---

Because this is too short to have as its own chapter...

**A first inference problem.** Author explores a particle decay problem of finding  $\Pr[\lambda | \{x\}]$  where  $\lambda$  is the characteristic decay length and  $\{x\}$  is our collection of observed decay distances. Plotting the likelihood  $\Pr[\{x\} | \lambda]$  as a function of  $\lambda$  for any given  $x \in \{x\}$  shows each has a *peak* value. The kicker is the interpretation: if each measurement  $x \in \{x\}$  is independent, the total likelihood is the product of all the individual likelihoods, which can be interpreted as updating/narrowing the interval  $[\lambda_a, \lambda_b]$  within which  $\Pr[\{x\} | \lambda]$  (as a function of  $\lambda$ ) peaks. In the words of the author's mentor:

*what you know about  $\lambda$  after the data arrive is what you knew before ( $\Pr[\lambda]$ ) and what the data told you ( $\Pr[\{x\} | \lambda]$ )*

We update our beliefs regarding the distribution of  $\lambda$  as we collect data.

#### Lessons learned from problems.

- (3.8) The classic Monty Hall problem. **Be careful defining probabilities after collecting data.** My blunder: when using Bayes' theorem to get the probability of the prize being behind door  $i \in \{1, 2\}$  after the host opens door 3, I failed to take into account that we chose door 1 while I was computing the evidence (the marginal probability that the host opened door 3)<sup>223</sup>. Quite embarrassing.
- (3.9) Monty Hall problem, but an earthquake opens door 3. Although I correctly answered that, in this case, both hypotheses (the prize being behind door 1 or door 2) are equiprobable, I still failed to account for the subtle fact that the earthquake could've opened *multiple* doors. The lesson here is **always write down the probability of everything**, which just so happens to be suggested by the solution for this problem, too.

So, why did I still get the answer correct? The reason is because enumeration of probabilities wasn't necessary at all, you just needed to realize that the likelihood for the two remaining hypotheses ( $\mathcal{H}_1$  and  $\mathcal{H}_2$ ) were the *same* – the probability of observing the earthquake open door 3 and the prize not being revealed was the same for the case of the prize being behind door 1 or door 2. So maybe the real lesson here is **determine whether calculations are even needed in order to solve the given problem**, which luckily I've had drilled in my head for years from studying physics.

- (3.15) Another biased coin variant. One of the best examples I've seen for favoring Bayesian methods over frequentist methods. Also, made use of the beta function:

$$\int_0^1 p^x (1-p)^y dp = \frac{\Gamma(x+1)\Gamma(y+1)}{\Gamma(x+y+2)} = B(x+1, y+1) \quad (998)$$

---

<sup>223</sup>Note that, while important to recognize and understand, I could've avoided this pitfall entirely by just ignoring the evidence during calculations and normalizing after, since the evidence can be determined solely by the normalization constraint.

where  $B$  is the **beta function**, which is defined by the LHS.

## The Source Coding Theorem (Ch. 4)

Table of Contents Local

Written by Brandon McKinzie

**Overview.** We will be examining the two following assertions:

1. The Shannon information content is a sensible measure of the information content of a given outcome  $x = a_i$ :

$$h(x = a_i) \triangleq \lg \frac{1}{p_i} \quad (999)$$

2. The entropy of an ensemble  $X$  is a sensible measure of the ensemble's average information content.

$$H(X) = \sum_i p_i \lg \frac{1}{p_i} \quad (1000)$$

**The weighing problem.** We are instructed to ponder the following problem.

*You're given 12 balls, all equal in weight except for one that is either heavier or lighter. You're also given a two-pan balance. Your task is to determine which ball is the odd ball, and in as few uses of the balance as possible. Note: each use of the balance must place the same number of balls on either side.*

An interesting observation is to consider the number of possible outcomes of the weighing process. Each outcome can be one of three possibilities: equal weight, left heavier, or left lighter. After  $N$  such weighings, the number of unique possible weighing result sequences is  $3^N$ . Note that there are  $12 \times 2 = 24$  unique final answers for our task (identifying which is the odd ball, and whether it is heavier or lighter). Therefore, since we are seeking a procedure to identify which of the 24 options is the correct option with 100% accuracy, we require our weighing procedure to take on *at least* 24 unique possible results. Since  $N = 3$  weighings corresponds to  $3^N = 27$  possible outcomes,  $N = 3$  is a *lower bound* on the number of weighings our approach will involve. It is impossible to guarantee a correct answer for  $N < 3$  weighings<sup>224</sup>.

Things I didn't consider until reading the solution:

- It's actually *not* optimal, upon observing both sides equal, to subsequently use only the balls not involved in that measurement. My initial reaction to this was "why? we already know the oddball is not any of the balls just measured, since the outcome was equal."

---

<sup>224</sup>Finally, it should be clear that, regardless of our approach, the final weighing will involve 2 balls, since we have to identify which is the oddball AND whether it is heavier/lighter AND the number of balls on the left of the scale must be the same as the right of the scale for every weighing.

The response to this reaction is: “yes, *exactly*, and we must use that information to be able to discern in the future whether, e.g., a measurement of “left side heavier” means the oddball is on the left and heavy, or if it’s on the right and light – *it’s useful to know that a given side of the scale does not contain the oddball before a measurement*.

- More generally, it’s also not optimal to greedily search for solutions that eliminate the highest number of possibilities *in any given single step*. Another way of thinking about this is that it’s undesirable for the  $i$ th measurement outcome to cause any of the 3 possible measurement outcomes to be impossible at the next stage.
- I focused a disproportionate amount of thought on handling the equal-weight measurement outcome, for whatever reason. I probably would’ve arrived at the solution faster if I’d actually thought about how my strategies would’ve handled some outcome being “left heavier” and *then considered what that strategy would put on the scale at the next step*, where the italics denote what would’ve illuminated the fatal flaw in all my approaches.

**Guessing Games.** What’s the smallest number of yes/no questions needed to identify an integer  $x$  between 0 and 63? Although it was obvious to me that the solution is to successively halve the possible values of  $x$ , I found it interesting that **you can write down the list of questions independent of the answers at each step** using a basic application of modular arithmetic. In other words, you can specify the full decision tree of  $N$  nodes with just  $\lg N$  questions. Nice. Also, recognize that the Shannon information content for any single outcome is  $\lg \frac{1}{0.5} = 1$  bit, and thus the total Shannon information content (for our predefined 6 questions) is 6 bits, which is not-coincidentally the number of possible values that  $x$  could be before we ask any questions.

In general, if an outcome  $x$  has Shannon information content  $h(x)$  number of bits, I like to interpret that as “learning the result  $x$  eliminates  $2^{h(x)}$  possibilities for the final result.” The battleship example follows this interpretation well. Stated another way (in the author’s words):

*The Shannon information content can be intimately connected to the size of a file that encodes the outcomes of a random experiment.*

---

### 8.3.1 DATA COMPRESSION AND TYPICALITY

---

**Data Compression.** A **lossy compressor** compresses some files, but maps some files to the *same* encoding. We introduce a parameter  $\delta$  that describes the risk (aggressiveness of our compression) we are taking with a given compression method:  $\delta$  is the probability that there will be no name for an outcome<sup>225</sup>  $x$ .

#### The smallest $\delta$ -sufficient subset

If  $S_\delta$  is the smallest subset of  $\mathcal{A}_X$  satisfying

$$\Pr [x \in S_\delta] \geq 1 - \delta \quad (1001)$$

then  $S_\delta$  is the smallest  $\delta$ -sufficient subset. It can be constructed by ranking the elements of  $\mathcal{A}_X$  in order of decreasing probability and adding successive elements starting from the most probable elements until the total probability is  $\geq (1 - \delta)$ .

- **Raw bit content** of  $X$ :  $H_0(X) \triangleq \lg |\mathcal{A}_X|$ . A lower bound for the number of binary questions that are always *guaranteed* to identify an outcome from the ensemble  $X$  – it simply maps each outcome to a constant-length binary string.
- **Essential bit content** of  $X$ :  $H_\delta(X) \triangleq \lg |S_\delta|$ . A compression code can be made by assigning a binary string of  $H_\delta$  bits to each element of the smallest sufficient subset.

Finally, we can now state **Shannon's source coding theorem**: Let  $X$  be an ensemble for the random variable  $x$  with entropy  $H(X) = H$  bits, and let  $X^N$  denote a sequence of identically distributed (but not necessarily independent<sup>226</sup>) of random variables/ensembles,  $(X_1, X_2, \dots, X_N)$ .

$$(\exists N_0 \in \mathbb{Z}^+) (\forall N > N_0) : \quad \left| \frac{1}{N} H_\delta(X^N) - H \right| < \epsilon \quad (0 < \delta < 1) \quad (\epsilon > 0) \quad (1002)$$

which, in English, can be read: *N i.i.d. random variables each with entropy  $H(X)$  can be compressed into more than  $NH(X)$  bits with negligible risk of information loss, as  $N \rightarrow \infty$ ; conversely if they're compressed into fewer than  $NH(X)$  bits it is virtually certain that information will be lost.*

---

<sup>225</sup>More specifically, if there is some subset,  $\{a\}$ , of unique values that  $x$  can take on but our compression method discards/ignores, then we say  $\delta = \sum_i p(x = a_i)$ .

<sup>226</sup>Actually, before the actual theorem statement, the author mentions we are now concerned with “string of  $N$  i.i.d. random variables from a single ensemble  $X$ .” It’s probably fair to assume this is true for the quantities in the theorem, but I’m leaving this note here as a reminder.

**Typicality.** The reason that large  $N$  in equation 1002 corresponds to larger potential for better compression is that the subset of likely results for a string of outcomes becomes more and more concentrated relative to the number of possible sequences as  $N$  increases<sup>227</sup>. I just realized this is for the same fundamental reasons that entropy exists in thermodynamics – *there are just more ways to exist in a high entropy state than otherwise*. The author showed  $\binom{N}{r}$  as a function of  $r$  (the number of 1s in the  $N$ -bit string). For large  $N$ , this becomes almost comically concentrated near the center (like a delta function at  $N/2$ ) – see footnote for more details<sup>228</sup>.

This motivates the notion of **typicality** for [a string of length  $N$  from] an arbitrary ensemble  $X$  with alphabet  $\mathcal{A}_X$ . For large  $N$ , we expect to find  $p_i N$  occurrences of the outcome  $x = a_i$ . Hence the probability of such a string, and its information content, is roughly<sup>229</sup>

$$\Pr[\mathbf{x}]_{typ} = \Pr[x_1] \Pr[x_2] \cdots \Pr[x_N] \simeq p_1^{(p_1 N)} p_2^{(p_2 N)} \cdots p_I^{(p_I N)} \quad (1003)$$

$$h(\mathbf{x})_{typ} = \lg \frac{1}{\Pr[\mathbf{x}]_{typ}} \simeq N \sum_{i=1}^I p_i \lg \frac{1}{p_i} = NH(X) \quad (1004)$$

Accordingly, we define the typical elements (strings of length  $N$ ) of  $\mathcal{A}_X^N$  to be those elements that have probability close to  $2^{-NH}$ . We introduce a parameter  $\beta$  that defines what we mean by “close,” and define the set of typical elements as the **typical set**  $T_{N\beta}$ :

$$T_{N\beta} \triangleq \left\{ \mathbf{x} \in \mathcal{A}_X^N : \left| \frac{1}{N} \lg \frac{1}{P(\mathbf{x})} - H \right| < \beta \right\} \quad (1005)$$

It turns out that whatever value of  $\beta$  we choose, the  $T_{N\beta}$  contains almost all the probability as  $N$  increases.

---

<sup>227</sup>The author gave an example for a sequence of bits with probability of any given bit being 1 as 0.1. He showed how, although the *average* number of 1s in a sequence of  $N$  bits grew as  $\mathcal{O}(N)$ , the standard deviation of that average only grew as  $\sqrt{N}$ .

<sup>228</sup>The probability of getting a string with  $r$  1s follows a binomial distribution with mean  $Np_1$  and standard deviation  $\sqrt{Np_1(1-p_1)}$ . This results in an increasingly narrower distribution  $P(r)$  for larger  $N$ .

<sup>229</sup>We appear to be assuming that each outcome  $x$  in the string  $\mathbf{x}$  are i.i.d. (CONFIRMED)

---

### 8.3.2 FURTHER ANALYSIS AND Q&A

---

#### Proving the Source Coding Theorem.

- **Setup.** We will make use of the following:

– **Chebyshev's Inequalities:**

$$\Pr [x \geq \alpha] \leq \frac{\mathbb{E}[x]}{\alpha} \quad \text{and} \quad \Pr [(x - \mathbb{E}[x])^2 \geq \alpha] \leq \frac{\text{Var}[x]}{\alpha} \quad (1006)$$

where  $\alpha$  is a positive real number, and  $x$  is assumed non-negative in the first inequality<sup>230</sup>.

- **Weak Law of Large Numbers** (WLLN): Consider a sample  $h_1, \dots, h_N$  of  $N$  independent RVs all with common mean  $\bar{h}$  and common variance  $\sigma_h^2$ . Let  $x = \frac{1}{N} \sum_{n=1}^N h_n$  be their average. Then

$$\Pr [(x - \bar{h})^2 \geq \alpha] \leq \frac{\sigma_h^2}{\alpha N} \quad (1007)$$

which can be easily derived from Chebyshev's inequalities.

- **Proving ‘asymptotic equipartition’ principle**, i.e. that an outcome  $\boldsymbol{x}$  is almost certain to belong to the typical set, approaching probability 1 for large enough  $N$ . It is a simple application of the WLLN to the random variable

$$\frac{1}{N} \lg \frac{1}{\Pr[\boldsymbol{x}]} = \frac{1}{N} \sum_{n=1}^N \lg \frac{1}{x_n} = \frac{1}{N} \sum_{n=1}^N h(x_n) \quad (1008)$$

where  $\mathbb{E}[h(x_n)] = H(X)$  for all terms in the summation. Observe, then, that the definition of the typical set given in equation 1005 (squaring both sides) has the same form as the definition for the WLLN. Plugging in and rearranging yields

$$\Pr [\boldsymbol{x} \in T_{N\beta}] \geq 1 - \frac{\sigma^2}{\beta^2 N} \quad (1009)$$

where  $\sigma^2 \equiv \text{Var} \left[ \lg \frac{1}{P(x_n)} \right]$ . This proves the asymptotic equipartition principle. It will also be useful to recognize that for any  $\boldsymbol{x}$  in the typical set, we can rearrange equation 1005 to obtain

$$2^{-N(H+\beta)} < \Pr [\boldsymbol{x}] < 2^{-N(H-\beta)} \quad (1010)$$

- **Proof of SCT Part I.** Want to show that  $\frac{1}{N} H_\delta(X^N) < H + \epsilon$ . **TODO**
- **Proof of SCT Part II.** Want to show that  $\frac{1}{N} H_\delta(X^N) > H - \epsilon$ . **TODO**

---

<sup>230</sup>Notice how the two inequalities are technically the same.

Questions & Answers. Collection of my questions as I read through the chapter, which I answered upon completing it.

- **Q:** Why isn't the essential bit content of a string of  $N$  i.i.d. variables  $NH$  when  $\delta = 0$ ?
  - **A:** I'm not entirely sure how to answer this still, but it seems the question is confused. First off, the essential bit content approaches the raw bit content as  $\delta$  decreases to 0:  $H_\delta \rightarrow H_0$  as  $\delta \rightarrow 0$ . It's important to notice that both  $H_\delta$  and  $H_0$  define an entropy where all members ( $S_\delta^N$  for  $H_\delta$ ;  $\mathcal{A}_X^N$  for  $H_0$ ) are *equiprobable*. I remember asking this question wondering “what is the significance of  $H_\delta(X^N)$  approaching  $NH(X)$  (not a typo!) for tiny  $\delta$ ”. The answer: for larger  $N$ , more of the probability mass is concentrated in a relatively smaller region, *with elements of that region being roughly equiprobable*. The last part is what I didn't initially realize – that **allowing for tiny  $\delta$  combined with large  $N$  essentially makes it so that  $S_\delta^N \approx T_{N\beta}$** .
- **Q:** Why aren't the most probable elements necessarily in the typical set?
  - **A:** In the limit of  $N \rightarrow \infty$ , *they are*, since in that limit, all elements are in the typical set and they're equiprobable. However, in essentially any real case, we can imagine that some elements will be too unlikely to be found within the typical set, which necessarily requires that there exist elements with probability too high to be in the typical set. Remember that the typical set is basically *defined* such that all elements have probability within the range given in equation 1010.

## Monte Carlo Methods (Ch. 29)

Table of Contents Local

*Written by Brandon McKinzie*

**Overview.** The aims of Monte Carlo methods are to solve one or both of the following:

1. Generate samples  $\{\mathbf{x}^{(r)}\}_{r=1}^R$  from a given probability distribution  $\Pr[\mathbf{x}]$ .
2. Estimate expectations of function under  $\Pr[\mathbf{x}]$ , for example

$$\Phi = \mathbb{E}_{\mathbf{x} \sim \Pr[\mathbf{x}]} [\phi(\mathbf{x})] \equiv \int d^N \mathbf{x} \Pr[\mathbf{x}] \phi(\mathbf{x}) \quad (1011)$$

where it's assumed that  $\Pr[\mathbf{x}]$  is sufficiently complex that we can't evaluate such expectations by exact methods.

Note that we can concentrate on the first problem (sampling), since we can use it to solve the second problem (estimating an expectation) by using the random samples  $\{\mathbf{x}^{(r)}\}_{r=1}^R$  to give the estimator

$$\hat{\Phi} \equiv \frac{1}{R} \sum_r \phi(\mathbf{x}^{(r)}) \quad (1012)$$

where

$$\mathbb{E} [\hat{\Phi}] = \frac{1}{R} \sum_r \mathbb{E}_{\mathbf{x}^{(r)} \sim \Pr[\mathbf{x}]} [\phi(\mathbf{x}^{(r)})] = \Phi \quad (1013)$$

$$\lim_{R \rightarrow \infty} \text{Var} [\hat{\Phi}] = \lim_{R \rightarrow \infty} \frac{\sum_r (\phi(\mathbf{x}^{(r)}) - \Phi)^2}{R-1} = \frac{\sigma^2}{R} \equiv \frac{1}{R} \int d^N \mathbf{x} \Pr[\mathbf{x}] (\phi(\mathbf{x}) - \Phi)^2 \quad (1014)$$

**Importance Sampling.** A generalization of [naively] uniformly sampling  $\mathbf{x}$  in order to approximate equation 1011. We assume henceforth that we are able to evaluate (for now, a 1-D)  $\Pr[x]$  at any point  $x$  at least within a multiplicative constant; thus we can evaluate a function  $P^*(x)$  such that  $P(x) = P^*(x)/Z$ . We assume we have some simpler  $Q(x)$ , called the **sampler density**, from which we *can* generate samples from and evaluate up to a multiplicative constant,  $Q(x) = Q^*(x)/Z_Q$ . We construct an approximation for our estimator in equation 1012 via sampling from  $Q(x)$  and computing:

$$\hat{\Phi} = \frac{\sum_r w_r \phi(x^{(r)})}{\sum_r w_r} \quad \text{where} \quad w_r \equiv \frac{P^*(x^{(r)})}{Q^*(x^{(r)})} \quad (1015)$$

**Rejection Sampling.** In addition to the assumptions in importance sampling, we now also assume that we know the value of a constant  $c$  such that

$$\forall x : cQ^*(x) > P^*(x) \quad (1016)$$

1. Generate sample  $x$  from proposal density  $Q(x)$ , and evaluate  $cQ^*(x)$ .
2. Generate a uniformly distributed random variable  $u$  from the interval  $[0, cQ^*(x)]$ .
3. If  $u > P^*(x)$ , then reject  $x$ ; else, accept  $x$  and add it to our set of samples  $\{x^{(r)}\}$ .

**Metropolis-Hastings Method.** Proposal density  $Q$  now depends on the current state  $x^{(t)}$ .

1. Sample tentative next state  $x'$  from proposal density  $Q(x'; x^{(t)})$ .
2. Compute:

$$a = \frac{P^*(x')}{P^*(x^{(t)})} \frac{Q(x^{(t)}; x')}{Q(x'; x^{(t)})} \quad (1017)$$

3. If  $a \geq 1$ , the new state is accepted. Otherwise, the new state is accepted with probability  $a$ .
4. If accepted, we set  $x^{(t+1)} = x'$ . If rejected, we set  $x^{(t+1)} = x^{(t)}$ .

## Variational Methods (Ch. 33)

Table of Contents Local

*Written by Brandon McKinzie*

**Probability Distributions in Statistical Physics**<sup>231</sup>. Consider the common situation below, where the state vector  $\mathbf{x} \in \{-1, +1\}^N$ :

$$\Pr[\mathbf{x} | \beta, \mathbf{J}] = \frac{e^{-\beta E(\mathbf{x}; \mathbf{J})}}{Z(\beta, \mathbf{J})} \quad (1018)$$

$$E(\mathbf{x}; \mathbf{J}) = -\frac{1}{2} \sum_{m,n} J_{mn} x_m x_n - \sum_n h_n x_n \quad (1019)$$

$$Z(\beta, \mathbf{J}) = \sum_{\mathbf{x}} e^{-\beta E(\mathbf{x}; \mathbf{J})} \quad (1020)$$

Note that evaluating  $E(\mathbf{x}; \mathbf{J})$  for a given  $\mathbf{x}$  takes polynomial time in the number of spins  $N$ , and evaluating  $Z$  is  $\mathcal{O}(2^N)$ . Variational free energy minimization is a method for approximating the complex distribution  $P(\mathbf{x})$  by a simpler ensemble  $Q(\mathbf{x}; \boldsymbol{\theta})$  parameterized by adjustable  $\boldsymbol{\theta}$ .

**Variational Free Energy.** How do we find/evaluate  $Q$ ? The objective function chosen to measure the quality of the approximation is the **variational free energy**,  $\tilde{F}(\boldsymbol{\theta})$ :

$$\beta \tilde{F}(\boldsymbol{\theta}) = \sum_{\mathbf{x}} Q(\mathbf{x}; \boldsymbol{\theta}) \ln \frac{Q(\mathbf{x}; \boldsymbol{\theta})}{e^{-\beta E(\mathbf{x}; \mathbf{J})}} \quad (1021)$$

$$= \beta \mathbb{E}_{\mathbf{x} \sim Q}[E(\mathbf{x}; \mathbf{J})] - H(Q) \quad (1022)$$

$$= D_{KL}(Q || P) + \beta F \quad (1023)$$

where  $F \triangleq -\ln Z(\beta, \mathbf{J})$  is the true free energy. It's not immediately clear why this approximation  $Q$  is more tractable than  $P$  – for that we turn to an example below.

---

<sup>231</sup>Yay!

**Variational Free Energy Minimization for Spin Systems.** For the system with energy function given in equation 1019, consider the *separable* approximating distribution,

$$Q(\mathbf{x}; \mathbf{a}) = \frac{e^{\sum_n a_n x_n}}{Z_Q} = \frac{\prod_n e^{a_n x_n}}{\sum_{\mathbf{x}'} \prod_{n'} e^{a_n x'_{n'}}} = \frac{\prod_n e^{a_n x_n}}{\sum_{x'_1} e^{a_1 x'_1} \sum_{x'_2} e^{a_2 x'_2} \cdots \sum_{x'_n} e^{a_n x'_n}} \quad (1024)$$

To compute  $\tilde{F}$ , we must compute the mean energy and entropy under  $Q$ .

→ **Mean Energy.** Given our definition of  $Q$ , and the fact that each  $x_n = \pm 1$ , the mean value of any  $x_n$  is  $\bar{x} = \tanh(a_n) = 2q_n - 1$ , where  $q_n \equiv Q(x_n = 1)$ .

$$\mathbb{E}_{\mathbf{x} \sim Q} [E(\mathbf{x}; \mathbf{J})] = \sum_{\mathbf{x}} Q(\mathbf{x}; \mathbf{a}) \left[ -\frac{1}{2} \sum_{m,n} J_{mn} x_m x_n - \sum_n h_n x_n \right] \quad (1025)$$

$$= -\frac{1}{2} \sum_{m,n} J_{mn} \bar{x}_m \bar{x}_n - \sum_n h_n \bar{x}_n \quad (1026)$$

→ **Entropy.** Recall that if  $Q(\mathbf{x}; \mathbf{a}) = \prod_n Q(x_n; a_n)$  (i.e.  $Q$  is separable), that  $H(\mathbf{x}) = \sum_n H(x_n)$ , so we have

$$H_Q(\mathbf{x}) = \sum_{\mathbf{x}} Q(\mathbf{x}; \mathbf{a}) \ln \frac{1}{Q(\mathbf{x}; \mathbf{a})} \quad (1027)$$

$$= \sum_n Q(x_n = 1; a_n) \ln \frac{1}{Q(x_n = 1; a_n)} + Q(x_n = 0; a_n) \ln \frac{1}{Q(x_n = 0; a_n)} \quad (1028)$$

$$= \sum_n q_n \ln \frac{1}{q_n} + (1 - q_n) \ln \frac{1}{1 - q_n} \quad (1029)$$

So the variational free energy is given by

$$\beta \tilde{F}(\mathbf{a}) = \beta \mathbb{E}_{\mathbf{x} \sim Q} [E(\mathbf{x}; \mathbf{J})] - H_Q(\mathbf{x}) \quad (1030)$$

$$= \beta \left( -\frac{1}{2} \sum_{m,n} J_{mn} \bar{x}_m \bar{x}_n - \sum_n h_n \bar{x}_n \right) - \left( \sum_n q_n \ln \frac{1}{q_n} + (1 - q_n) \ln \frac{1}{1 - q_n} \right) \quad (1031)$$

Remember, our goal is to find parameters  $\mathbf{a}$  that minimize  $\tilde{F}(a)$ :

$$\beta \frac{\partial \tilde{F}}{\partial a_m} = 2 \left( \frac{\partial q_m}{\partial a_m} \right) \left[ -\beta \left( \sum_n J_{mn} \bar{x}_n + h_m \right) + a_m \right] \quad (1032)$$

which, when set to zero, yields

$$a_m \leftarrow \beta \left( \sum_n J_{mn} \bar{x}_n + h_m \right) \quad (1033)$$

$$\bar{x}_n = \tanh(a_n) \quad (1034)$$

# MACHINE LEARNING: A PROBABILISTIC PERSPECTIVE

## CONTENTS

9.1	Probability (Ch. 2) . . . . .	342
9.1.1	Exercises . . . . .	343
9.2	Generative Models for Discrete Data (Ch. 3) . . . . .	344
9.2.1	Exercises . . . . .	348
9.3	Gaussian Models (Ch. 4) . . . . .	349
9.4	Bayesian Statistics (Ch. 5) . . . . .	353
9.5	Linear Regression (Ch. 7) . . . . .	355
9.6	Generalized Linear Models and the Exponential Family (Ch. 9) . . . . .	358
9.7	Mixture Models and the EM Algorithm (Ch. 11) . . . . .	361
9.8	Latent Linear Models (Ch. 12) . . . . .	364
9.9	Markov and Hidden Markov Models (Ch. 17) . . . . .	366
9.10	Undirected Graphical Models (Ch. 19) . . . . .	368

## Probability (Ch. 2)

Table of Contents Local

Written by Brandon McKinzie

Kevin P. Murphy (2012). Probability.

*Machine Learning: A Probabilistic Perspective.*

**Continuous Random Variables** (2.2.5). Let  $X$  be a continuous RV. We usually want to know the probability that  $X$  lies in the interval  $a \leq X \leq b$ , which is given by

$$p(a < X \leq b) = p(X \leq b) - p(X \leq a) \quad (1035)$$

Define the **cumulative distribution function** (cdf)  $F(q) \triangleq p(X \leq q)$ , and the **probability density function** (pdf)  $f(x) \triangleq \frac{d}{dx} F(x)$ .

**Transformation of Random Variables** (2.6). In what follows, we have some  $x \sim p_x$  and  $y = f(x)$ . How should we think about  $p_y(y)$ ? For discrete RV  $x$ , we just sum up the probability mass for all  $x$  such that  $f(x) = y$ ,

$$p_y(y) = \sum_{x:f(x)=y} p_x(x) \quad (1036)$$

If  $\mathbf{x}$  is *continuous*, we must instead work with the cdf,

$$P_y(y) \triangleq P(Y \leq y) = P(f(X) \leq y) = P(X \in \{x \mid f(x) \leq y\}) \quad (1037)$$

If  $f$  is monotonic (and hence invertible), then we can derive the pdf  $p_y(y)$  from the pdf  $p_x(x)$  by taking derivatives as follows:

$$p_y(y) \triangleq \frac{d}{dy} P_y(y) = \frac{d}{dy} P_x(f^{-1}(y)) = \frac{dx}{dy} p_x(x) \quad (1038)$$

and, since we only work with integrals over densities (i.e. overall sign does not matter), it is convention to take the absolute value of  $\frac{dx}{dy}$  in the above formula. In the multivariate case, the Jacobian matrix  $[\mathbf{J}_{\mathbf{y} \rightarrow \mathbf{x}}]_{i,j} \triangleq \frac{\partial x_i}{\partial y_j}$  is used:

$$p_y(\mathbf{y}) = p_x(\mathbf{x}) \left| \det \mathbf{J}_{\mathbf{y} \rightarrow \mathbf{x}} \right| \quad (1039)$$

**Central Limit Theorem** (2.6.3). Consider  $N$  i.i.d. RVs each with arbitrary pdf  $p(x_i)$ , mean  $\mu$ , and covariance  $\sigma^2$ . Let  $S_N = \sum_i X_i$  denote the sum over the  $N$  RVs. The **central limit theorem** states that

$$\lim_{N \rightarrow \infty} S_N \sim \mathcal{N}(N\mu, N\sigma^2) \quad (1040)$$

$$\text{or equivalently} \quad \lim_{N \rightarrow \infty} \sqrt{N}(\bar{X} - \mu) \sim \mathcal{N}(0, \sigma^2) \quad (1041)$$

### 9.1.1 EXERCISES

#### Exercise 2.1

(a) [correct]  $P(\text{oneIsGirl} \mid \text{hasAtLeastOneBoy}) = 2/3$ . Use muh intuition.

(b) [correct]  $P(\text{otherIsGirl} \mid \text{weSawOneIsABoy}) = P(\text{childIsGirl}) = 1/2$ . The other being a girl/boy is independent of the fact that the other is a boy. All about how it is phrased, yo.

#### Exercise 2.2 - Variance of a sum

Show that  $\text{Var}[X + Y] = \text{Var}[X] + \text{Var}[Y] + 2\text{Cov}[X, Y]$ .

[correct] Math approach:

$$\text{Var}[X + Y] \triangleq \mathbb{E}[(X + Y) - \mathbb{E}[X] - \mathbb{E}[Y]]^2 \quad (1042)$$

$$= \mathbb{E}[((X - \mathbb{E}[X]) + (Y - \mathbb{E}[Y]))^2] \quad (1043)$$

$$= \text{Var}[X] + \text{Var}[Y] + 2\text{Cov}[X, Y] \quad (1044)$$

Intuition approach: If  $X$  and  $Y$  are uncorrelated, it is intuitive that their sum should have variance equal to the sum of the individual variances. If there is some linear dependence between the two, we'd expect it to either increase (if positively correlated) or decrease (if negatively correlated) the variance of their sum.

## Generative Models for Discrete Data (Ch. 3)

Table of Contents Local

Written by Brandon McKinzie

Kevin P. Murphy (2012). Generative Models for Discrete Data.

*Machine Learning: A Probabilistic Perspective*.

**Bayesian Concept Learning** (3.1). The interesting notion of learning a *concept*  $\mathcal{C}$ , such as “all prime numbers”, by only seeing positive examples  $x \in \mathcal{C}$ . How should we approach predicting whether a new test case  $\tilde{x}$  belongs to the concept  $\mathcal{C}$ ? Well, what we’re actually doing is as follows: Given an initial **hypothesis space**  $\mathcal{H}$  of concepts, we collect data  $\mathcal{D}$  that gradually narrows the down the subset of  $\mathcal{H}$  consistent with our data. We also need to address how we (humans) will weigh certain hypotheses differently even if they are both entirely consistent with the evidence. The Bayesian approach can be summarized as follows (no particular order):

- **Likelihood.** Probability of observing  $\mathcal{D}$  given a particular hypothesis  $h$ . In the simple but illustrative case where the data is sampled from a uniform distribution over the *extension*<sup>232</sup> of a concept (a.k.a. the *strong sampling assumption*). The probability here of sampling  $N$  items independently under hypothesis  $h$  is

$$p(\mathcal{D} | h) = \left[ \frac{1}{|\mathcal{H}|} \right]^N \quad (1045)$$

which elucidates how *the model favors the smallest hypothesis space consistent with the data*<sup>233</sup>.

- **Prior.** Using just the likelihood can mean we fall prey to contrived/overfitting hypotheses that basically just enumerate the data. The prior  $p(h)$  allows us to specify properties about how the learned hypothesis ought to look.
- **Posterior.** This is just the likelihood times the prior, normalized [to one]. In general, as we collect more and more data, the posterior tends toward a Dirac measure peaked at the MAP estimate:

$$p(h | \mathcal{D}) \rightarrow \delta_{\hat{h}^{MAP}}(h) \quad \text{where} \quad (1046)$$

$$\hat{h}^{MAP} = \arg \max_h p(h | \mathcal{D}) = \arg \max_h p(\mathcal{D} | h)p(h) = \arg \max_h [\log p(\mathcal{D} | h) + \log p(h)] \quad (1047)$$

---

<sup>232</sup>The extension of a concept is just the set of numbers that belong to it.

<sup>233</sup>A result commonly known as **Occam’s razor** or the **size principle**.

**The Beta-Binomial Model** (3.3). In the previous section we considered inferring some discrete distribution over integers; now we will turn to a continuous version. Consider the problem of inferring the probability that a coin shows up heads, given a series of observed coin tosses.

- **Likelihood.** As should be familiar, we'll model the outcome of each toss  $X_i \in \{1, 0\}$  indicating heads or tails with  $X_i \sim \text{Ber}(\theta)$ . Assuming the tosses are i.i.d, this gives us  $p(\mathcal{D} | \theta) \propto \theta^{N_1}(1 - \theta)^{N_0}$ , where  $N_1$  and  $N_0$  are the **sufficient statistics** of the data, since  $p(\theta | \mathcal{D})$  can be entirely modeled as  $p(\theta | N_1, N_0)$ .
- **Prior.** We technically just need a prior  $p(\theta)$  with support over the interval  $[0, 1]$ , but it would be convenient if the prior had the same form as the likelihood, i.e.  $p(\theta) \propto \theta^{\gamma_1}(1 - \theta)^{\gamma_2}$  for some hyperparameters  $\gamma_1$  and  $\gamma_2$ . This is satisfied by the **Beta distribution**<sup>234</sup>. This would also result in the posterior having the same form as the prior, meaning the prior is a **conjugate prior** for the corresponding likelihood.
- **Posterior.** As mentioned, the posterior corresponding to a Bernoulli/binomial likelihood with a beta prior is itself a beta distribution:

$$p(\theta | \mathcal{D}) \propto \text{Bin}(N_1 | \theta, N_0 + N_1) \text{Beta}(\theta | a, b) \propto \text{Beta}(\theta | N_1 + a, N_0 + b) \quad (1048)$$

By either reading off from a table or deriving via calculus, we can obtain the following properties for our Beta posterior:

$$\text{[mode]} \quad \hat{\theta}_{MAP} = \frac{a + N_1 - 1}{a + b + N - 2} \quad (1049)$$

$$\text{[mean]} \quad \bar{\theta} = \frac{a + N_1}{a + b + N} = \lambda m_1 + (1 - \lambda) \hat{\theta}_{MLE} \quad (1050)$$

where  $m_1 = \frac{a}{a+b}$  is the prior mean. The last form captures the notion that the posterior is a compromise between what we previously believed and what the data is telling us.

**The Dirichlet-Multinomial Model** (3.4). We now generalize further to inferring the probability that a die with  $K$  sides comes up as face  $k$ .

- **Likelihood.** As before, we assume a specific observed sequence  $\mathcal{D}$  of  $N$  dice roles. Assuming the data is i.i.d., the likelihood has the form

$$p(\mathcal{D} | \boldsymbol{\theta}) = \prod_{k=1}^K \theta_k^{N_k} \quad (1051)$$

which is the likelihood for the multinomial model up to an irrelevant constant factor.

- **Prior.** We'd like to find a conjugate prior for our likelihood, and we need it to have support over the  $K$ -dimensional *probability simplex*<sup>235</sup>. The Dirichlet distribution satisfies

---

<sup>234</sup>You may be wondering, why not the Bernoulli distribution? It trivially has the same form as the Bernoulli distribution, eh? Then, pause and actually think about what you're saying for five seconds. You want to model a prior on  $\theta$  with a Bernoulli distribution? You do realize that the support for a Bernoulli is in  $k \in \{0, 1\}$ . It's the opposite domain entirely. We want something that "looks" like a Bernoulli but is a distribution over  $\theta$ , NOT the value(s) of  $x$ .

<sup>235</sup>The  $K$ -dimensional probability simplex is the  $(K-1)$ th dimensional simplex determined by the unit vectors  $e_1, \dots, e_K \in \mathbb{R}^K$ , i.e. the set of vectors  $\mathbf{x}$  such that  $x_i \geq 0$  and  $\sum_i x_i = 1$ .

both of these and is defined as

$$\text{Dir}(\boldsymbol{\theta} \mid \boldsymbol{\alpha}) = \frac{1}{B(\boldsymbol{\alpha})} \prod_{k=1}^K \theta_k^{\alpha_k - 1} \quad (1052)$$

where  $\boldsymbol{\theta} \in S_K$  is a built-in assumption.

- **Posterior.** By construction, this will also be Dirichlet. Note that to derive the MAP estimator we must enforce the constraint that  $\sum_k \theta_k = 1$ , which can be done with a **Lagrange multiplier**. The constrained objective (the Lagrangian) is

$$\ell(\boldsymbol{\theta}, \lambda) = \sum_k N_k \log \theta_k + \sum_k (\alpha_k - 1) \log \theta_k + \lambda \left( 1 - \sum_k \theta_k \right) \quad (1053)$$

To get  $\hat{\theta}_{MAP}$ , we'd take derivatives w.r.t.  $\lambda$  and each  $\theta_k$ , do some substitutions and solve.

### Example: Language Models with Bag of Words

*Given a sequence of words, predict the most likely next word. Assume that the  $i$ th word  $X_i \in \{1, \dots, K\}$  (where  $K$  is the size of our vocab) is sampled indep from all others using a  $\text{Cat}(\boldsymbol{\theta})$  (multinoulli) distribution. This is the BoW model.*

My attempt: We need to derive the form of posterior predictive  $p(X_{i+1} \mid X_1, \dots, X_i)$  where  $\boldsymbol{\theta} \in S_K$ . First, I know that the posterior is

$$p(\boldsymbol{\theta} \mid X_1, \dots, X_i) \propto \text{Dir}(\boldsymbol{\theta} \mid \boldsymbol{\alpha}) P(X_1, \dots, X_i \mid \boldsymbol{\theta}) = \text{Dir}(\boldsymbol{\theta} \mid \alpha_1 + N_1, \dots, \alpha_K + N_K) \quad (1054)$$

and so I can derive the posterior predictive in the usual way, while also exploiting the fact that all  $X_i \perp X_j$ ,

$$p(X_{i+1} = k \mid X_1, \dots, X_i) = \int p(X_{i+1} \mid \boldsymbol{\theta}) p(\boldsymbol{\theta} \mid X_1, \dots, X_i) d\boldsymbol{\theta} \quad (1055)$$

$$= \int \theta_k p(\theta_k, \boldsymbol{\theta}_{-k} \mid X_1, \dots, X_i) d\theta_k d\boldsymbol{\theta}_{-k} \quad (1056)$$

$$= \int \theta_k p(\theta_k \mid X_1, \dots, X_i) d\theta_k \quad (1057)$$

$$= \mathbb{E} [\theta_k \mid X_1, \dots, X_i] \quad (1058)$$

$$= \frac{\alpha_k + N_k}{\sum_j \alpha_j + N_j} \quad (1059)$$

which also shows another nice example of Bayesian smoothing.

**Naive Bayes Classifiers** (3.5). For classifying vectors of discrete-valued features,  $\mathbf{x} \in \{1, \dots, K\}^D$ . Assumes features are conditionally independent given the class label:

$$p(\mathbf{x} \mid y = c, \boldsymbol{\theta}) = \prod_{j=1}^D p(x_j \mid y = c, \theta_{j,c}) \quad (1060)$$

with parameters  $\boldsymbol{\theta} \in \mathbb{R}^{D \times |\mathcal{Y}|}$ <sup>236</sup>. We can model the individual class-conditional densities with the multinoulli distribution. If we were modeling real-valued features, we could instead use a Gaussian distribution.

- **MLE fitting.** Our log-likelihood is

$$\log p(\mathcal{D} \mid \boldsymbol{\theta}) = \sum_c N_c \log \pi_c + \sum_i^N \sum_j^D \log p(x_j^{(i)} \mid y^{(i)}, \theta_{j,y^{(i)}}) \quad (1061)$$

where  $\pi_c = p(y = c)$  are the class priors<sup>237</sup>. We see that we can optimize the class priors separately from the others, and that they have the same form as the multinomial likelihood we used in the last section. We already know that the MLE for these are  $\hat{\pi}_c = N_c/N$  (remember this involves using a Lagrangian). Let's assume next that we're working in the case of binary features ( $x_j \mid c \sim \text{Ber}(\theta_{j,c})$ ). This results in MLE estimates  $\hat{\theta}_{j,c} = N_{j,c}/N_c$ .

---

<sup>236</sup>You could also generalize this to having some number of params  $N$  associated with each pairwise  $(j, c)$ . It's also important to recognize that this is the first model of this chapter where the input  $\mathbf{x}$  is a *vector*, and thus introduces pairwise parameters.

<sup>237</sup>We only see the class prior parameters here because they appear in the likelihood of generative classifiers, since  $p(x, y) = p(y)p(x \mid y)$ . We don't see the priors for  $\theta$  that aren't class priors because MLE is only concerned with maximizing the likelihood, not the posterior (which would contain those priors).

---

### 9.2.1 EXERCISES

---

#### MLE for uniform distribution (3.8)

The density function for the uniform distribution centered on 0 with width 2a is

$$p(x) = \frac{1}{2a} \mathbb{1}\{x \in [-a, a]\}$$

Remember this is for continuous  $x$ , and  $p(x)$  is interpreted as the derivative of the corresponding CDF  $P(x)$ .

- a. Given data  $x_1, \dots, x_n$ , find  $\hat{a}_{MLE}$ . So there are a few ways of doing this. We can get the answer pretty quick with intuition, and not-so-quick by grinding through the math. **Quick-n-easy:** If you were paying attention to the chapter, you'd instantly remember that MLE is all about finding the smallest hypothesis space consistent with the data. It should then be obvious that  $\hat{a}_{MLE} = \max |x_i|$ . **Slightly more rigorous.** We can also solve a constrained optimization problem, with constraint that  $a \geq \max |x_i|$ , since we must require our solution to assign nonzero probability to any of our observations.

$$\hat{a}_{MLE} = \arg \max_a \log p(x_1, \dots, x_n | a) + \lambda(a - \max |x_i|) \quad (1062)$$

The rest is mechanical: (1) take deriv wrt  $\lambda$ , (2) deriv wrt  $a$  and set to zero, (3) solve for  $a$  as a function of  $\lambda$ , (4) solve for  $\lambda$  by substituting previous step's results into constraint, yielding that  $\lambda \leq n/|x_{max}|$ , (5) plug result for  $\lambda$  into result from step 3 to obtain result that  $a \geq x_{max}$ . In the limit of many samples, the first term becomes more important in the optimization, which decreases as  $a$  increases, and so we choose the lowest value of  $a$  that satisfies the constraints:  $a := x_{max}$ .

- b. What probability would the model assign to a new data point  $x_{n+1}$  using  $\hat{a}$ . We are obviously meant to trivially answer that it will assign the density with  $a = \hat{a}$ . However, I take objection to this question, since it makes no sense to evaluate a density at a single point  $p(x)$ .  
c. Do you see any problem with this approach? Yes, it extremely overfits to the data. We'll assign zero probabilities to any points outside the interval of observations, and the same probability to anything in that interval. We should do a more Bayesian approach instead.

## Gaussian Models (Ch. 4)

Table of Contents Local

Written by Brandon McKinzie

Kevin P. Murphy (2012). Gaussian Models.

*Machine Learning: A Probabilistic Perspective.*

**Basics** (4.1). I'll be filling in the gaps that the book leaves out in its derivations, as a way of reviewing the relevant linear algebra/calculus/etc. For notation's sake, here how the author writes the MVN in D dimensions:

$$\mathcal{N}(\mathbf{x} \mid \boldsymbol{\mu}, \boldsymbol{\Sigma}) \triangleq \frac{1}{(2\pi)^{D/2} |\boldsymbol{\Sigma}|^{1/2}} \exp \left\{ -\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}) \right\} \quad (1063)$$

To get a better understanding, let's inspect the eigendecomposition of  $\boldsymbol{\Sigma}$ . We know that  $\boldsymbol{\Sigma}$  is positive definite<sup>238</sup>, and therefore the eigendecomposition  $\boldsymbol{\Sigma} = \mathbf{U}\boldsymbol{\Lambda}\mathbf{U}^T$  exists, where  $\mathbf{U}$  is an orthonormal matrix of eigenvectors. By the invertible matrix theorem, we therefore know that  $\mathbf{U}$  is invertible. Since  $\boldsymbol{\Sigma}$  is p.d., it's eigenvalues are all positive, and thus  $\boldsymbol{\Lambda}$  is also invertible. We can then apply the basic definition for an invertible matrix to write

$$\boldsymbol{\Sigma}^{-1} = \mathbf{U}\boldsymbol{\Lambda}^{-1}\mathbf{U}^T = \sum_{i=1}^D \frac{1}{\lambda_i} \mathbf{u}_i \mathbf{u}_i^T \quad (1064)$$

Decomposition of  $\boldsymbol{\Sigma}$ 

Remember, an orthonormal matrix satisfies  $\mathbf{U}^T = \mathbf{U}^{-1}$ .

where  $\mathbf{u}_i$  is the  $i$ th eigenvector and  $i$ th *column* of  $\mathbf{U}$ . We can use this to rewrite

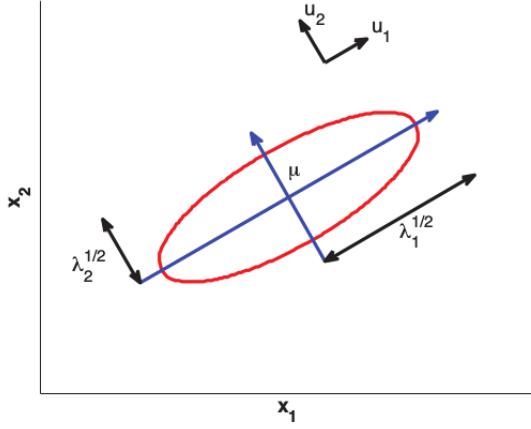
$$(\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}) = (\mathbf{x} - \boldsymbol{\mu})^T \left( \sum_i^D \frac{1}{\lambda_i} \mathbf{u}_i \mathbf{u}_i^T \right) (\mathbf{x} - \boldsymbol{\mu}) \quad (1065)$$

Rewriting in form of an ellipse.

$$= \sum_i^D \frac{1}{\lambda_i} (\mathbf{x} - \boldsymbol{\mu})^T \mathbf{u}_i \mathbf{u}_i^T (\mathbf{x} - \boldsymbol{\mu})^T \quad (1066)$$

$$= \sum_i^D \frac{1}{\lambda_i} y_i^2 \quad (1067)$$

<sup>238</sup>We know this because all covariance matrices of any random vector  $X$  are symmetric p.s.d., and the additional requirement that  $\boldsymbol{\Sigma}^{-1}$  exists means that it is p.d.



where  $y_i \triangleq \mathbf{u}_i^T(\mathbf{x} - \boldsymbol{\mu})$ . The fascinating insight is recalling that the equation for an ellipse in 2D is

$$\frac{y_1^2}{\lambda_1} + \frac{y_2^2}{\lambda_2} = 1 \quad (1068)$$

Hence we see that the contours of equal probability density of a Gaussian lie along ellipses, as illustrated above. The eigenvectors determine the orientation of the ellipse, and the eigenvalues determine how elongated it is.

**Maximum Entropy Derivation of the Gaussian** (4.1.4). Recall that the exponential family can be derived as the family of distributions  $p(\mathbf{x})$  that maximizes  $H(p)$  subject to constraints that the moments of  $p$  match some set of empirical moments  $F_k$  specified by us. It turns out that *the MVN is the distribution with maximum entropy subject to having a specified mean and covariance*. Consider the zero-mean MVN and its entropy<sup>239</sup>:

$$p(\mathbf{x}) = \frac{1}{Z} \exp \left\{ -\frac{1}{2} \mathbf{x}^T \boldsymbol{\Sigma}^{-1} \mathbf{x} \right\} \quad (1069)$$

$$h(p) = \frac{1}{2} \ln \left[ (2\pi e)^D \det \boldsymbol{\Sigma} \right] \quad (1070)$$

Let  $p = \mathcal{N}(\mathbf{0}, \boldsymbol{\Sigma})$  above and let  $q(\mathbf{x})$  be any density satisfying  $\int q(\mathbf{x}) x_i x_j d\mathbf{x} = \Sigma_{ij}$ . The result is based on the fact that  $h(q)$  must be less than or equal to  $h(p)$ . This can be shown by evaluating  $D_{KL}(q||p)$  and recalling that  $D_{KL}$  is always greater than or equal to zero.

---

<sup>239</sup>For derivation, see this wonderful answer on stackexchange.

**Gaussian Discriminant Analysis** (4.2). With generative classifiers, it is common to define the class-conditional density as a MVN,

$$p(\mathbf{x} \mid y=c, \boldsymbol{\theta}) = \mathcal{N}(\mathbf{x} \mid \boldsymbol{\mu}_c, \boldsymbol{\Sigma}_c) \quad (1071)$$

which results in a technique called (Gaussian) **discriminant**<sup>240</sup> **analysis**. If  $\boldsymbol{\Sigma}_c$  is diagonal, this is just a form of naive Bayes<sup>241</sup>. We classify some feature vector  $\mathbf{x}$  using the decision rule:

$$\hat{y}(\mathbf{x}) = \arg \max_c [\log p(y=c \mid \boldsymbol{\pi}) + \log p(\mathbf{x} \mid y=c, \boldsymbol{\theta})] \quad (1072)$$

If we have uniform prior over classes, we can classify a new test vector as follows:

$$\hat{y}(\mathbf{x}) = \arg \min_c (\mathbf{x} - \boldsymbol{\mu}_c)^T \boldsymbol{\Sigma}_c^{-1} (\mathbf{x} - \boldsymbol{\mu}_c) \quad (1073)$$

**Linear Discriminant Analysis** (LDA). The special case where all covariance matrices are the same,  $\boldsymbol{\Sigma}_c = \boldsymbol{\Sigma}$ . Now the quadratic term  $\mathbf{x}^T \boldsymbol{\Sigma}^{-1} \mathbf{x}$  is independent of  $c$  and thus is not important for classification. Instead we end up with the much simpler,

$$p(y=c \mid \mathbf{x}, \boldsymbol{\theta}) = \frac{e^{\beta_c^T \mathbf{x} + \gamma_c}}{\sum_{c'} e^{\beta_{c'}^T \mathbf{x} + \gamma_{c'}}} = \mathcal{S}(\boldsymbol{\eta})_c \quad (1074)$$

$$\boldsymbol{\beta}_c := \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}_c \quad (1075)$$

$$\gamma_c := -\frac{1}{2} \boldsymbol{\mu}_c^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}_c + \log \pi_c \quad (1076)$$

where  $\mathcal{S}$  is the familiar **softmax**. **TODO**: come back and compare/contrast LDA with multinomial logistic regression after reading chapter 8.

**Inference in Jointly Gaussian Distributions** (4.3). Suppose  $\mathbf{x} = (\mathbf{x}_1, \mathbf{x}_2)$  is jointly Gaussian with parameters

$$\boldsymbol{\mu} = \begin{pmatrix} \boldsymbol{\mu}_1 \\ \boldsymbol{\mu}_2 \end{pmatrix}, \quad \boldsymbol{\Sigma} = \begin{pmatrix} \boldsymbol{\Sigma}_{11} & \boldsymbol{\Sigma}_{12} \\ \boldsymbol{\Sigma}_{21} & \boldsymbol{\Sigma}_{22} \end{pmatrix} \quad (1077)$$

Then, the marginals are given by

$$p(\mathbf{x}_1) = \mathcal{N}(\mathbf{x}_1 \mid \boldsymbol{\mu}_1, \boldsymbol{\Sigma}_{11}) \quad (1078)$$

$$p(\mathbf{x}_2) = \mathcal{N}(\mathbf{x}_2 \mid \boldsymbol{\mu}_2, \boldsymbol{\Sigma}_{22}) \quad (1079)$$

and the posterior conditional is given by

---

<sup>240</sup>Don't confuse the word "discriminant" for "discriminative" – we are still in a generative model setting. See 8.6 for details on the distinction.

<sup>241</sup>This is easy to show. If diagonal, then  $p(\mathbf{x} \mid y)$  factorizes. We know the  $i$ th item in the product corresponds to  $p(x_i \mid c)$  by considering how simple it is to compute marginals for Gaussians with diagonal  $\boldsymbol{\Sigma}$ .

$$p(\mathbf{x}_1 \mid \mathbf{x}_2) = \mathcal{N}(\mathbf{x}_1 \mid \boldsymbol{\mu}_{1|2}, \boldsymbol{\Sigma}_{1|2})$$

$$\text{where } \boldsymbol{\mu}_{1|2} = \boldsymbol{\mu}_1 + \boldsymbol{\Sigma}_{12}\boldsymbol{\Sigma}_{22}^{-1}(\mathbf{x}_2 - \boldsymbol{\mu}_2) \quad (4.69)$$

$$\boldsymbol{\Sigma}_{1|2} = \boldsymbol{\Sigma}_{11} - \boldsymbol{\Sigma}_{12}\boldsymbol{\Sigma}_{22}^{-1}\boldsymbol{\Sigma}_{21} = \boldsymbol{\Lambda}_{11}^{-1}$$

where  $\boldsymbol{\Lambda} := \boldsymbol{\Sigma}^{-1}$ . Notice that the conditional covariance is a *constant* matrix independent of  $\mathbf{x}_2$ . The proof here relies on Schur complements (see Appendix).

**Information Form** (4.3.3). Thus far, we've been working in terms of the **moment parameters**  $\boldsymbol{\mu}$  and  $\boldsymbol{\Sigma}$ . We can also rewrite the MVN in terms of its **natural (canonical) parameters**  $\boldsymbol{\Lambda}$  and  $\boldsymbol{\xi}$ ,

$$\mathcal{N}_c(\mathbf{x} \mid \boldsymbol{\xi}, \boldsymbol{\Lambda}) = (2\pi)^{-D/2} |\boldsymbol{\Lambda}|^{\frac{1}{2}} \exp \left\{ -\frac{1}{2} (\mathbf{x}^T \boldsymbol{\Lambda} \mathbf{x} + \boldsymbol{\xi}^T \boldsymbol{\Lambda}^{-1} \boldsymbol{\xi} - 2\mathbf{x}^T \boldsymbol{\xi}) \right\} \quad (1080)$$

$$\text{where } \boldsymbol{\Lambda} \triangleq \boldsymbol{\Sigma}^{-1} \quad \text{and} \quad \boldsymbol{\xi} \triangleq \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu} \quad (1081)$$

where  $\mathcal{N}_c$  is how we'll denote “in canonical form.” The marginals and conditionals in canonical form are

$$p(\mathbf{x}_2) = \mathcal{N}_c(\mathbf{x}_2 \mid \boldsymbol{\xi}_2 - \boldsymbol{\Lambda}_{21}\boldsymbol{\Lambda}_{11}^{-1}\boldsymbol{\xi}_1, \boldsymbol{\Lambda}_{22} - \boldsymbol{\Lambda}_{21}\boldsymbol{\Lambda}_{11}^{-1}\boldsymbol{\Lambda}_{12}) \quad (1082)$$

$$p(\mathbf{x}_1 \mid \mathbf{x}_2) = \mathcal{N}_c(\mathbf{x}_1 \mid \boldsymbol{\xi}_1 - \boldsymbol{\Lambda}_{12}\mathbf{x}_2, \boldsymbol{\Lambda}_{11}) \quad (1083)$$

and we see that **marginalization is easy in moment form, while conditioning is easier in information form.**

**Linear Gaussian Systems** (4.4). Suppose we have a hidden variable  $\mathbf{x}$  and a noisy observation of it  $\mathbf{y}$ . Let's assume we have the following prior and likelihood:

$$p(\mathbf{x}) = \mathcal{N}(\mathbf{x} \mid \boldsymbol{\mu}_x, \boldsymbol{\Sigma}_x)$$

$$p(\mathbf{y} \mid \mathbf{x}) = \mathcal{N}(\mathbf{y} \mid \mathbf{A}\mathbf{x} + \mathbf{b}, \boldsymbol{\Sigma}_y) \quad (4.124)$$

which is an example of a **linear Gaussian system**  $\mathbf{x} \rightarrow \mathbf{y}$ .

**The Wishart Distribution** (4.5). We now dive into the distributions over the *parameters*  $\boldsymbol{\Sigma}$  and  $\boldsymbol{\mu}$ . First, we need some math prereqs out of the way. The **Wishart** is the generalization of the Gamma to pd matrices:

$$\text{Wi}(\boldsymbol{\Lambda} \mid \mathbf{S}, \nu) = \frac{1}{Z_{Wi}} |\boldsymbol{\Lambda}|^{(\nu-D-1)/2} \exp \left\{ -\frac{1}{2} \text{tr}(\boldsymbol{\Lambda} \mathbf{S}^{-1}) \right\} \quad (1084)$$

$$Z_{Wi} = 2^{\nu D/2} \Gamma_D(\nu/2) |\mathbf{S}|^{\nu/2} \quad (1085)$$

where

- $\nu$ : degrees of freedom
- $\mathbf{S}$ : scale matrix (a.k.a. scatter matrix). Basically empirical  $\boldsymbol{\Sigma}$ .
- $\Gamma_D$ : multivariate gamma function

A nice property is that if  $\boldsymbol{\Sigma}^{-1} \sim \text{Wi}(\mathbf{X}, \nu)$ , then  $\mathbf{Sigma} \sim \text{IW}(\mathbf{S}^{-1}, \nu + D + 1)$ , the **inverse Wishart** (multi-dim generalization of inv Gamma).

## Bayesian Statistics (Ch. 5)

Table of Contents Local

Written by Brandon McKinzie

Kevin P. Murphy (2012). Bayesian Statistics.

*Machine Learning: A Probabilistic Perspective.*

**MAP Estimation** (5.2.1). The most popular **point estimate** for parameters  $\theta$  is the posterior mode, aka the **MAP estimate**. However, there are many drawbacks:

- No measure of uncertainty (true for any point estimate).
- Using  $\theta_{MAP}$  for predictions is prone to overfitting.
- The mode is an atypical point.
- It's not invariant to reparameterization. Say two possible parameterizations  $\theta_1$  and  $\theta_2 = f(\theta_1)$ , where  $f$  is some deterministic function. In general, it is **not** the case that  $\hat{\theta}_2 = f(\hat{\theta}_1)$  under MAP.

**Bayesian Model Selection** (5.3). A model selection technique where we compute the best model  $m$  for data  $\mathcal{D}$  using the formulas,

$$p(m \mid \mathcal{D}) = \frac{p(\mathcal{D} \mid m)p(m)}{\sum_{m \in \mathcal{M}} p(\mathcal{D} \mid m)p(m)} \quad (1086)$$

$$p(\mathcal{D} \mid m) = \int p(\mathcal{D} \mid \theta)p(\theta \mid m)d\theta \quad (1087)$$

where the latter is the **marginal likelihood**<sup>242</sup> for model  $m$ . Note that this isn't anything new; we've usually just denoted it simply as  $p(\mathcal{D})$ , since typically  $m$  is specified beforehand. Although large models with many parameters can achieve higher likelihoods under MLE/MAP,  $p(\mathcal{D} \mid \hat{\theta}_m)$ , this is *not* necessarily the case with *marginal likelihood*, an effect known as the **Bayesian Occam's razor**. Below we give the marginal likelihoods for familiar models:

- Beta-Binomial:

$$p(\mathcal{D} \mid m=\text{BetaBinom}) = \binom{N}{N_1} \frac{B(a + N_1, b + N_0)}{B(a, b)}$$

where  $B$  is the Beta function.

- Dirichlet-Multinoulli:

$$p(\mathcal{D}) = \frac{B(N + \alpha)}{\alpha}$$

---

<sup>242</sup>Also called the integrated likelihood or the evidence.

**BIC Approximation** (5.3.2.4). The integral involved in computing  $p(\mathcal{D} \mid m)$  (henceforth denoted simply as  $p(\mathcal{D})$ ) can be intractable. The **Bayesian information criterion** (BIC) is a popular approximation:

$$\text{BIC} \triangleq \log p(\mathcal{D} \mid \hat{\boldsymbol{\theta}}) - \frac{1}{2} \text{dof}(\hat{\boldsymbol{\theta}}) \log N \approx \log p(\mathcal{D}) \quad (1088)$$

where

- $\text{dof}(\hat{\boldsymbol{\theta}})$  is the number of **degrees of freedom** in the model.
- $\hat{\boldsymbol{\theta}}$  is the MLE for the model.

BIC is also closely related to the **minimum description length** (MDL) principle and the **Akaike information criterion** (AIC).

**Hierarchical Bayes** (5.5). When defining our prior  $p(\boldsymbol{\theta} \mid \boldsymbol{\eta})$ , we have to of course specify the hyperparameters  $\boldsymbol{\eta}$  required by our choice of prior. The Bayesian approach for doing this is to put a prior on our prior! This situation can be represented as a directed graphical model, illustrated below.



This is an example of a **hierarchical Bayesian model**, also called a **multi-level model**.

**Bayesian Decision Theory** (5.7). Decision problems can be cast as games against nature, where natures selects a quantity  $y \in \mathcal{Y}$  unknown to us, and then generates an observation  $\mathbf{x} \in \mathcal{X}$  that we get to see. Our goal is to devise a **decision procedure** or **policy**  $\delta : \mathcal{X} \mapsto \mathcal{A}$  for generating an action  $a \in \mathcal{A}$  from observation  $\mathbf{x}$  that's deemed most compatible with the hidden state  $y$ . We define “compatible” via a loss function  $L(y, a)$ :

$$\delta(\mathbf{x}) = \arg \min_{a \in \mathcal{A}} \rho(a \mid \mathbf{x}) \quad (1089)$$

$$\text{where } \rho(a \mid \mathbf{x}) = \mathbb{E}_{p(y|\mathbf{x})} [L(y, a)] \quad (1090)$$

In this context, we call  $\rho$  the **posterior expected loss**, and  $\delta(x)$  the **Bayes estimator**. Some Bayes estimators for common loss functions are given below.

- **0-1 loss:**  $L(y, a) = \mathbb{1}\{y \neq a\}$ . Easy to show that  $\delta(x) = \arg \max_{y \in \mathcal{Y}} p(y \mid \mathbf{x})$ .

## Linear Regression (Ch. 7)

Table of Contents Local

Written by Brandon McKinzie

Kevin P. Murphy (2012). Linear Regression.

*Machine Learning: A Probabilistic Perspective.*

**Model Specification** (7.2). The linear regression model is a model of the form

$$p(y | \mathbf{x}, \boldsymbol{\theta}) = \mathcal{N}(y | \mathbf{w}^T \mathbf{x}, \sigma^2) \quad (1091)$$

**Maximum Likelihood Estimation** (least squares) (7.3). Most commonly, we'll estimate the parameters by computing the MLE, defined by

$$\hat{\boldsymbol{\theta}} \triangleq \arg \max_{\boldsymbol{\theta}} \log p(\mathcal{D} | \boldsymbol{\theta}) \quad (1092)$$

$$= \arg \min_{\boldsymbol{\theta}} [-\log p(\mathcal{D} | \boldsymbol{\theta})] \quad (1093)$$

$$= -\frac{1}{2\sigma^2} RSS(\mathbf{w}) - \frac{N}{2} \log(2\pi\sigma^2) \quad (1094)$$

$$RSS(\mathbf{w}) \triangleq \sum_{i=1}^N (y_i - \mathbf{w}^T \mathbf{x}_i)^2 \quad (1095)$$

where  $RSS(\mathbf{w})$  is the **residual sum of squares**. Notice that  $\boldsymbol{\theta} := (\mathbf{w}, \sigma^2)$ , but typically we're focused on estimating  $\mathbf{w}$ <sup>243</sup>.

**Derivation of the MLE** (7.3.1). We'll now denote the negative log likelihood as  $NLL(\mathbf{w})$  and drop constant terms that are irrelevant for the optimization task.

$$NLL(\mathbf{w}) = \frac{1}{2} \|\mathbf{y} - \mathbf{X}\mathbf{w}\|_2^2 = \frac{1}{2} \mathbf{w}^T (\mathbf{X}^T \mathbf{X}) \mathbf{w} - \mathbf{w}^T (\mathbf{X}^T \mathbf{y}) \quad (1096)$$

$$\text{where } \mathbf{X}^T \mathbf{X} = \sum_{i=1}^N \mathbf{x}_i \mathbf{x}_i^T \quad (1097)$$

$$\mathbf{X}^T \mathbf{y} = \sum_{i=1}^N \mathbf{x}_i y_i \quad (1098)$$

---

<sup>243</sup>Since our goal is typically to make future predictions  $\hat{y}(\mathbf{x}) = \mathbf{w}^T \mathbf{x}$ , rather than sampling  $y \sim p(y | \mathbf{x}, \boldsymbol{\theta})$ , we aren't concerned with estimating  $\sigma^2$ . We assume some variability, and the goal is focused on fitting the data to a straight line.

And the optimal  $\hat{\mathbf{w}}_{OLS}$  be found by taking the gradient, setting to zero, and solving for  $\mathbf{w}$  as usual:

$$\nabla \text{NLL}(\mathbf{w}) = \mathbf{X}^T \mathbf{X} \mathbf{w} - \mathbf{X}^T \mathbf{y} \quad (1099)$$

$$\mathbf{X}^T \mathbf{X} \mathbf{w} = \mathbf{X}^T \mathbf{y} \quad [\text{normal eq.}] \quad (1100)$$

$$\hat{\mathbf{w}}_{OLS} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y} \quad (1101)$$

**Geometric Interpretation** (7.3.2). Use the column-vector representation of  $\mathbf{X} \in \mathbb{R}^{N \times D}$ , where we assume  $N > D$ <sup>244</sup>. Our prediction can then be written

$$\hat{\mathbf{y}} = \mathbf{X} \mathbf{w} = w_1 \tilde{\mathbf{x}}_1 + \cdots + w_D \tilde{\mathbf{x}}_D \quad (1102)$$

i.e. a linear combination of the  $D$  column vectors  $\tilde{\mathbf{x}}_i \in \mathbb{R}^N$ . Crucially, observe that this means  $\hat{\mathbf{y}} \in \text{span}(\{\tilde{\mathbf{x}}\}_i^D)$  no matter what (a hard constraint by definition of our model). So, how do you minimize the residual norm  $\|\mathbf{y} - \hat{\mathbf{y}}\|$  given that  $\hat{\mathbf{y}}$  is restricted to a particular subspace? You require  $\mathbf{y} - \hat{\mathbf{y}}$  to be orthogonal to that subspace, of course<sup>245</sup>! Formally, this means  $\tilde{\mathbf{x}}_j^T (\mathbf{y} - \hat{\mathbf{y}}) = 0$ , for all  $1 \leq j \leq D$ . Equivalently,

$$\mathbf{X}^T (\mathbf{y} - \mathbf{X} \mathbf{w}) = \mathbf{0} \implies \hat{\mathbf{w}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y} \quad (1103)$$

$$\hat{\mathbf{y}} = \mathbf{X} \hat{\mathbf{w}} = \mathbf{P} \mathbf{y} \quad (1104)$$

$$\text{where } \mathbf{P} \triangleq \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \quad (1105)$$

Although this neat, I'm left unsatisfied since there appears to be no intuition of what the column vectors of  $\mathbf{X}$  really mean on a conceptual level.

**Robust Linear Regression** (7.4). Gaussians sensitive to outliers, since their log-likelihood penalizes deviations quadratically<sup>246</sup>. One way to achieve **robustness** to outliers is to instead use a distribution with **heavy tails**, such that they still allow for higher likelihoods of outliers, but they don't need to shift the whole distribution around to accommodate for them. One popular choice is the **Laplace distribution**,

$$p(y | \mathbf{x}, \mathbf{w}, b) = \text{Lap}(y | \mathbf{w}^T \mathbf{x}, b) \triangleq \frac{1}{2b} \exp \left\{ -\frac{1}{b} |y - \mathbf{w}^T \mathbf{x}| \right\} \quad (1106)$$

$$NLL(\mathbf{w}) = \sum_i |y_i - \mathbf{w}^T \mathbf{x}_i| \quad (1107)$$

Replacing Gaussian with heavy-tailed dists.

---

<sup>244</sup>This means we have more rows than columns, which means our column space cannot span all of  $\mathbb{R}^N$ .

<sup>245</sup>Consider that  $\mathbf{y} - \hat{\mathbf{y}}$  points from our prediction (which is in the constraint subspace)  $\hat{\mathbf{y}}$  to the true  $\mathbf{y}$  that we want to get closest to. Intuitively that means  $\hat{\mathbf{y}}$  is looking “straight out” at  $\mathbf{y}$ , in a direction orthogonal to the subspace that  $\hat{\mathbf{y}}$  lives in. Formally, we can write  $\mathbf{y} = (\mathbf{y}_{\parallel}, \mathbf{y}_{\perp})$ , where  $\mathbf{y}_{\parallel}$  is the component within  $\text{Col}(\mathbf{x})$ . The best we can do, then, is  $\hat{\mathbf{y}} := \mathbf{y}_{\parallel}$ .

<sup>246</sup>In other words, outliers initially get huge loss values, and the distribution shifts toward them to minimize loss (undesirably).

Goal: convert  $NLL$  to a form easier to optimize (linear). Let  $r_i \triangleq y_i - \mathbf{w}^T \mathbf{x}_i$  be the  $i$ 'th residual. The following steps show how we can convert this into a **linear program**:

$$r_i \triangleq r_i^+ - r_i^- \quad (r_i^+ \geq 0)(r_i^- \geq 0) \quad (1108)$$

$$\min_{\mathbf{w}, \mathbf{r}^+, \mathbf{r}^-} \sum_i (r_i^+ + r_i^-) \quad \text{s.t.} \quad \mathbf{w}^T \mathbf{x}_i + r_i^+ - r_i^- = y_i \quad (1109)$$

$$\min_{\boldsymbol{\theta}} \mathbf{f}^T \boldsymbol{\theta} \quad \text{s.t.} \quad \mathbf{A}\boldsymbol{\theta} \leq \mathbf{b}, \quad \mathbf{A}_{eq}\boldsymbol{\theta} = \mathbf{b}_{eq}, \quad \mathbf{1} \leq \boldsymbol{\theta} \leq \mathbf{u} \quad (1110)$$

where the last equation is the **standard form** of a LP.

**Ridge Regression** (7.5). We know that MLE can overfit by essentially memorizing the data. If, for example, we model 21 points with a degree-14 polynomial<sup>247</sup>, we get many large positive and negative numbers for our learned coefficients, which allow the curve to wiggle in just the right way to almost perfectly interpolate the data – **this is why we often regularize weights to have low absolute value**. This encourages smoother/less-wiggly curves. One way to do this is by using a zero-mean Gaussian prior on our weights:

$$p(\mathbf{w}) = \prod_j \mathcal{N}(w_j | 0, \tau^2) \quad (1111)$$

This makes our MAP estimation problem and solution take the form

$$J(\mathbf{w}) = \frac{1}{N} \sum_i^N (y_i - \mathbf{w}^T \mathbf{x}_i - w_0)^2 + \lambda \|\mathbf{w}\|_2^2 \quad (1112)$$

$$\hat{\mathbf{w}}_{ridge} = (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I}_D)^{-1} \mathbf{X}^T \mathbf{y} \quad (1113)$$

Doing MAP instead of  
MLE

Note that  $w_0$  is NOT  
regularized.

---

<sup>247</sup>To review polynomials, search “Lagrange interpolation” in your CS 70 notes.

## Generalized Linear Models and the Exponential Family (Ch. 9)

Table of Contents Local

Written by Brandon McKinzie

Kevin P. Murphy (2012). Generalized Linear Models and the Exponential Family.

*Machine Learning: A Probabilistic Perspective.*

**The Exponential Family** (9.2). Why is the exponential family important?

- It's the only family with **finite-sized sufficient statistics**<sup>248</sup>.
- It's the only family for which **conjugate priors exist**.
- It makes the **least set of assumptions** subject to some user-chosen constraints.
- It's at the core of GLMs and variational inference.

A *pdf or pmf*  $p(\mathbf{x} \mid \boldsymbol{\theta})$ , for  $\mathbf{x} \in \mathcal{X}^m$  and  $\boldsymbol{\theta} \in \Theta \subseteq \mathbb{R}^d$ , is said to be in the **exponential family** if it's of the form

$$p(\mathbf{x} \mid \boldsymbol{\theta}) = \frac{1}{Z(\boldsymbol{\theta})} h(\mathbf{x}) \exp\{\boldsymbol{\theta}^T \phi(\mathbf{x})\} = h(\mathbf{x}) \exp\left\{\boldsymbol{\theta}^T \phi(\mathbf{x}) - A(\boldsymbol{\theta})\right\} \quad (1114)$$

$$Z(\boldsymbol{\theta}) = \int_{\mathcal{X}^m} h(\mathbf{x}) \exp\{\boldsymbol{\theta}^T \phi(\mathbf{x})\} \quad (1115)$$

$$A(\boldsymbol{\theta}) = \log Z(\boldsymbol{\theta}) \quad (1116)$$

$h(\mathbf{x})$  is a scaling constant, often 1.

where  $\boldsymbol{\theta}$  are the **natural (canonical) parameters**<sup>249</sup>, and  $\phi(\mathbf{x})$  are the **sufficient statistics**.

Below are some (quick/condensed) examples showing the first couple steps in rewriting familiar distributions in exponential family form:

$$\text{[Bernoulli]} \quad \text{Ber}(x \mid \mu) = \mu^x (1 - \mu)^{1-x} = \exp\{x \log \mu + (1 - x) \log(1 - \mu)\} \quad (1117)$$

$$\text{[Multinoulli]} \quad \text{Cat}(x \mid \boldsymbol{\mu}) = \prod_k^K \mu_k^{x_k} = \exp\left\{\sum_{k=1}^{K-1} x_k \log(\mu_k / \mu_K) + \log \mu_K\right\} \quad (1118)$$

<sup>248</sup>Given certain regularity conditions.

<sup>249</sup>We often generalize this with  $\eta(\boldsymbol{\theta})$ , which maps whatever params  $\boldsymbol{\theta}$  we've chosen to the canonical params  $\eta(\boldsymbol{\theta})$ .

**Log Partition Function** (9.2.3). The derivatives of the log partition,  $A(\boldsymbol{\theta})$ , can be used to generate **cumulants**<sup>250</sup> for the sufficient statistics,  $\phi(\mathbf{x})$ .

$$\frac{\partial A(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}} = \mathbb{E} [\phi(\mathbf{x})] \quad (1119)$$

$$\nabla^2 A(\boldsymbol{\theta}) = \text{cov} [\phi(\mathbf{x})] \quad (1120)$$

**MLE for the Exponential Family** (9.2.4). The likelihood takes the form

$$p(\mathcal{D} | \boldsymbol{\theta}) = \left[ \prod_i^N h(\mathbf{x}^{(i)}) \right] g(\boldsymbol{\theta})^N \exp \left\{ \boldsymbol{\eta}(\boldsymbol{\theta})^T \phi(\mathcal{D}) \right\} \quad (1121)$$

It appears that we are denoting  $1/Z$  with  $g$  now.

$$\phi(\mathcal{D}) = \sum_i^N \phi(\mathbf{x}^{(i)}) = \begin{bmatrix} \sum_i \phi_1 & \cdots & \sum_i \phi_K \end{bmatrix} \quad (1122)$$

where I've denoted  $\sum_{i=1}^N \phi_k(\mathbf{x}^{(i)})$  as simply  $\sum_i \phi_k$ . The **Pitman-Koopman-Darmois theorem** states, given certain regularity conditions/constraints<sup>251</sup>, that the exponential family is the only family with *finite sufficient statistics* (dimensionality independent of the size of the data set). For example, in the above formula, we have  $K + 1$  sufficient statistics (+1 since we need to know the value of  $N$ ).

Consider a canonical<sup>252</sup> exponential family which also sets  $h(\cdot) = 1$ . The log-likelihood is

$$p(\mathcal{D} | \boldsymbol{\theta}) = \boldsymbol{\theta}^T \phi(\mathcal{D}) - N A(\boldsymbol{\theta}) \quad (1123)$$

Since  $-A(\boldsymbol{\theta})$  is concave<sup>253</sup> and the other term is linear in  $\boldsymbol{\theta}$ , the log-likelihood is concave and thus has a global maximum.

---

<sup>250</sup>The first and second cumulants are mean and variance.

<sup>251</sup>The wording is weird here. We mean “out of all families/distributions that *already* satisfy certain constraints that must be met, the exponential family is the only...”. For example, the uniform distribution has finite statistics and is not in the exponential family, but it does not meet the constraint that its support must be independent of the parameters, so it’s outside the scope of the theorem.

<sup>252</sup>Defined as those which satisfy  $\boldsymbol{\eta}(\boldsymbol{\theta}) = \boldsymbol{\theta}$ .

<sup>253</sup>We know  $-A$  is concave because  $A$  is convex. We know  $A$  is convex because  $\nabla^2 A$  is positive definite. Remember that any twice-differentiable multivariate function  $f$  is convex IFF its Hessian is pd for all  $\boldsymbol{\theta}$ . See sec 7.3.3 and 9.2.3 for more.

**Maximum Entropy Derivation of the Exponential Family** (9.2.6). Suppose we don't know which distribution  $p$  to use, but we do know the expected values of certain features or functions:

$$F_k \triangleq \mathbb{E}_{\mathbf{x} \sim p(\mathbf{x})} [f_k(\mathbf{x})] = \sum_{\mathbf{x}} f_k(\mathbf{x}) p(\mathbf{x}) \quad (1124)$$

The principle of **maximum entropy** or **maxent** says we should pick the distribution with maximum entropy, subject to the constraints that the moments of the distribution match the empirical moments of the specified functions  $f_k(\mathbf{x})$ . Treating  $p$  as a fixed length vector (i.e. assuming  $\mathbf{x}$  is discrete), we can take the derivative of our Lagrangian (entropy in units of nats with constraints) w.r.t. each "element"  $p_x = p(\mathbf{x})$  to find the optimal distribution.

$$J(p, \boldsymbol{\lambda}) = H(p) + \lambda_0 \left( 1 - \sum_{\mathbf{x}} p(\mathbf{x}) \right) + \sum_k \lambda_k \left( F_k - \sum_{\mathbf{x}} p(\mathbf{x}) f_k(\mathbf{x}) \right) \quad (1125)$$

$$\frac{\partial J}{\partial p(\mathbf{x})} = -1 - \log p(\mathbf{x}) - \lambda_0 - \sum_k \lambda_k f_k(\mathbf{x}) \quad (1126)$$

Setting this derivative to zero yields

$$p(\mathbf{x}) = \frac{1}{Z} \exp \left\{ - \sum_k \lambda_k f_k(\mathbf{x}) \right\} \quad (1127)$$

Thus the maxent distribution  $p(\mathbf{x})$  has the form of the exponential family, a.k.a. the **Gibbs distribution**.

## Mixture Models and the EM Algorithm (Ch. 11)

Table of Contents Local

Written by Brandon McKinzie

Kevin P. Murphy (2012). Mixture Models and the EM Algorithm.

*Machine Learning: A Probabilistic Perspective.*

**Latent Variable Models** (LVMs) (11.1). In this chapter, we explore directed GMs that have hidden/latent variables. Advantages of LVMs:

1. Often have fewer params.
2. Hidden vars can serve as a bottleneck (representation learning).

**Mixture Models** (11.2). The simplest form of LVM is where the hidden variables  $z_i \in \{1, \dots, K\}$  represent a discrete latent state. We use discrete prior  $p(z_i) = \text{Cat}(\boldsymbol{\pi}) = \pi_i$ , and likelihood  $p(\mathbf{x}_i | z=k) = p_k(\mathbf{x}_i)$ . A **mixture model** is defined by

$$p(\mathbf{x}_i | \boldsymbol{\theta}) = \sum_{k=1}^K \pi_k p_k(\mathbf{x}_i | \boldsymbol{\theta}) \quad (1128)$$

We call  $p_k$  the *kth base distribution*.

Some popular mixture models:

- **Mixture of Gaussians** (11.2.1). Each  $p_k = \mathcal{N}(\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$ . Given large enough  $K$ , a GMM can approximate any density defined on  $\mathbb{R}^D$ .
- **Mixture of Multinoullis** (11.2.2). Let  $\mathbf{x} \in \{0, 1\}^D$ . Each  $p_k(\mathbf{x}) = \prod_j^D \text{Ber}(x_j | \mu_{jk})$ .

Below, I derive the expectation and covariance of  $\mathbf{x}$  in this case<sup>254</sup>.

$$\mathbb{E} [\mathbf{x}] = \sum_{\mathbf{x}} \mathbf{x} p(\mathbf{x}) \quad (1129)$$

$$= \sum_{\mathbf{x}} \mathbf{x} \sum_k^K \pi_k p_k(\mathbf{x}) \quad (1130)$$

$$= \sum_{\mathbf{x}} \mathbf{x} \sum_k^K \pi_k \prod_j^D \text{Ber}(x_j | \mu_{jk}) \quad (1131)$$

$$= \sum_k^D \pi_k \sum_{x_1} \cdots \sum_{x_D} \mathbf{x} \prod_j^D \text{Ber}(x_j | \mu_{jk}) \quad (1132)$$

$$= \sum_k^D \pi_k \sum_{x_1} \text{Ber}(x_1 | \mu_{1k}) \cdots \sum_{x_D} \text{Ber}(x_D | \mu_{Dk}) \mathbf{x} \quad (1133)$$

$$= \sum_k^D \pi_k \boldsymbol{\mu}_k \quad (1134)$$

Next we want to find  $\text{cov}(\mathbf{x}) = \mathbb{E} [\mathbf{x}\mathbf{x}^T] - \mathbb{E} [\mathbf{x}] \mathbb{E} [\mathbf{x}]^T$ . I think the insight that makes finding the first term easiest is realizing that you only need to find the two cases,  $\mathbb{E} [x_i^2]$  and  $\mathbb{E} [x_i x_{j \neq i}]$ , where in this case

$$\mathbb{E} [x_i^2] = \sum_k \pi_k \mu_{ik} \quad (1135)$$

$$\mathbb{E} [x_i x_{j \neq i}] = \sum_k \pi_k \mu_{ik} \mu_{jk} \quad (1136)$$

$$\therefore \mathbb{E} [\mathbf{x}\mathbf{x}^T] = \sum_k \pi_k (\boldsymbol{\Sigma}_k + \boldsymbol{\mu}_k \boldsymbol{\mu}_k^T) \quad (1137)$$

where  $\boldsymbol{\Sigma}_k = \text{diag}(\mu_{jk}(1 - \mu_{jk}))$  is the covariance of  $\mathbf{x}$  under  $p_k$ . The fact that the mixture covariance matrix is now non-diagonal confirms that the mixture can capture correlations between variables  $x_i, x_{j \neq i}$ , unlike a single product-of-Bernoullis model.

**The EM Algorithm** (11.4). In LVMs, it's usually intractable to compute the MLE since we have to marginalize over hidden variables while satisfying constraints like positive definite covariance matrices, etc. The EM algorithm gets around these issues via a two-step process. The first step involves taking an expectation, where the expectation is over  $\mathbf{z} \sim p(\mathbf{z} | \mathbf{x}, \theta^{t-1})$  for each individual observed  $\mathbf{x}$ , where we use the current parameter estimates when sampling  $\mathbf{z}$  in the expectation. **This gives us an auxiliary likelihood that's a function of  $\theta$**  which will serve as a stand-in (in the 2nd step) for what we typically use as the likelihood in MLE. The second step is then just finding the optimal  $\theta^t$  over the auxiliary likelihood function from the first step. This iterates until convergence or some stopping condition.

---

<sup>254</sup>Shown in excruciating detail because I was unable to work through this in my head alone.

### Procedure: EM Algorithm

First, let's define our auxiliary function  $Q$  as

$$Q(\boldsymbol{\theta} \mid \boldsymbol{\theta}^{t-1}) = \mathbb{E}_{p(\mathbf{z} \mid \mathbf{x}, \boldsymbol{\theta}^{t-1})} [\ell_c(\boldsymbol{\theta}) \mid \mathcal{D}] \quad (1138)$$

$$\text{where } \ell_c(\boldsymbol{\theta}) = \sum_{i=1}^N \log p(\mathbf{x}^{(i)}, \mathbf{z}^{(i)} \mid \boldsymbol{\theta}) \quad (1139)$$

where, again, the “expectation” serves the purpose of determining the expected value of  $\mathbf{z}^{(i)}$  for each observed  $\mathbf{x}^{(i)}$ . It's somewhat of a misnomer to denote the expectation like this, since each  $\mathbf{z}^{(i)}$  is innately tied with its corresponding observation  $\mathbf{x}^{(i)}$ .

1. **E-Step:** Evaluate  $Q(\boldsymbol{\theta} \mid \boldsymbol{\theta}^{t-1})$  using (obviously) the previous parameters  $\boldsymbol{\theta}^{t-1}$ . This yields a function of  $\boldsymbol{\theta}$ . This gives us the expected log-likelihood function of  $\boldsymbol{\theta}$  for the observed data  $\mathcal{D}$ .

2. **M-Step:** Optimize  $Q$  w.r.t  $\boldsymbol{\theta}$  to get  $\boldsymbol{\theta}^t$ :

$$\boldsymbol{\theta}^t = \arg \max_{\boldsymbol{\theta}} Q(\boldsymbol{\theta} \mid \boldsymbol{\theta}^{t-1}) \quad (1140)$$

## Latent Linear Models (Ch. 12)

Table of Contents Local

Written by Brandon McKinzie

Kevin P. Murphy (2012). Latent Linear Models.

*Machine Learning: A Probabilistic Perspective.*

**Factor Analysis** (12.1). Whereas mixture models define  $p(z) = \text{Cat}(\boldsymbol{\pi})$  for a single hidden variable  $z \in \{1, \dots, K\}$ , **factor analysis** begins by instead using a *vector* of *real-valued* latent variables,  $\mathbf{z} \in \mathbb{R}^L$ . The simplest prior is  $\mathbf{z} \sim \mathcal{N}(\boldsymbol{\mu}_0, \boldsymbol{\Sigma}_0)$ . If  $\mathbf{x} \in \mathbb{R}^D$ , we can define

$$p(\mathbf{x}_i | \mathbf{z}_i, \boldsymbol{\theta}) = \mathcal{N}(\mathbf{W}\mathbf{z}_i + \boldsymbol{\mu}, \boldsymbol{\Psi}) \quad (1141)$$

where

- $\mathbf{W} \in \mathbb{R}^{D \times L}$ : **factor loading matrix**.
- $\boldsymbol{\Psi} \in \mathbb{R}^{D \times D}$  is a diagonal covariance matrix, since we want to “force”  $\mathbf{z}$  to explain the correlation<sup>255</sup>. If  $\boldsymbol{\Psi} = \sigma^2 \mathbf{I}$ , we get **probabilistic PCA**.

Summaries of key points regarding FA:

- **Low-rank parameterization of MVN.** FA can be thought of as specifying  $p(\mathbf{x})$  using a small number of parameters. [math] yields that

$$\text{cov}(\mathbf{x}) = \mathbf{W}\mathbf{W}^T + \boldsymbol{\Psi}$$

which has  $\mathcal{O}(LD)$  params (remember  $\boldsymbol{\Psi}$  is diagonal) instead of the usual  $\mathcal{O}(D^2)$ .

- **Unidentifiability:** The params of an FA model are unidentifiable.

**Classical PCA** (12.2.1). Goal: find an orthogonal set of  $L$  linear basis vectors  $\mathbf{w}_j \in \mathbb{R}^D$ , and the scores  $\mathbf{z}_i \in \mathbb{R}^L$ , such that we minimize the average **reconstruction error**:

$$J(\mathbf{W}, \mathbf{Z}) = \frac{1}{N} \sum_{i=1}^N \|\mathbf{x}_i - \hat{\mathbf{x}}_i\|^2, \quad \text{where } \hat{\mathbf{x}}_i := \mathbf{W}\mathbf{z}_i \quad (1142)$$

$$= \|\mathbf{X} - \mathbf{W}\mathbf{Z}^T\|_F^2 \quad (1143)$$

where  $\mathbf{W} \in \mathbb{R}^{D \times L}$  is orthonormal. **Solution:** assign each column  $\mathbf{W}_{:, \ell}$  to the eigenvector with  $\ell$ 'th largest eigenvalue of  $\hat{\boldsymbol{\Sigma}} = \frac{1}{N} \sum_i^N \mathbf{x}_i \mathbf{x}_i^T$ , assuming  $\mathbb{E}[\mathbf{x}] = \mathbf{0}$ . Then we compute  $\hat{\mathbf{z}}_i := \mathbf{W}^T \mathbf{x}_i$ .

---

<sup>255</sup>It's easier to think of this graphically. Our model asserts that  $\mathbf{x}_i \perp \mathbf{x}_{j \neq i} | \mathbf{z}$ . Independence implies zero correlation, and we cement this by constraining  $\boldsymbol{\Psi}$  to be diagonal. See your related note on LFM, chapter 13 of the DL book.

## Proof: PCA

**Case: L=1.** Goal: Find the best 1-d solution,  $\mathbf{w} \in \mathbb{R}^D$ ,  $z_i \in \mathbb{R}$ ,  $\mathbf{z} \in \mathbb{R}^N$ . Remember that  $\|\mathbf{w}\|_2 = 1$ .

$$J(\mathbf{w}, \mathbf{z}) = \frac{1}{N} \sum_i^N \|\mathbf{x}_i - z_i \mathbf{w}\|^2 = \frac{1}{N} [\mathbf{x}_i^T \mathbf{x}_i - 2z_i \mathbf{w}^T \mathbf{x}_i + z_i^2] \quad (1144)$$

$$\frac{\partial J}{\partial z_i} = 0 \rightarrow z_i = \mathbf{w}^T \mathbf{x}_i \quad (1145)$$

$$J(\mathbf{w}) = \frac{1}{N} \sum_i^N [\mathbf{x}_i^T \mathbf{x}_i - z_i^2] = \text{const} - \frac{1}{N} \sum_i^N z_i^2 \quad (1146)$$

$$\therefore \arg \min_{\mathbf{w}} J(\mathbf{w}) = \arg \max_{\mathbf{w}} \text{Var}[\tilde{\mathbf{z}}] \quad (1147)$$

This shows why PCA finds directions of maximal variance – aka the [analysis view](#) of PCA. Before finding the optimal  $\mathbf{w}$ , don't forget the Lagrange multipliers for constraining unit norm,

$$\tilde{J}(\mathbf{w}) = \mathbf{w}^T \hat{\Sigma} \mathbf{w} + \lambda(\mathbf{w}^T \mathbf{w} - 1) \quad (1148)$$

**Singular Value Decomposition (SVD) (12.2.3).** Any real  $N \times D$  matrix  $\mathbf{X}$  can be decomposed as

$$\mathbf{X} = \underbrace{\mathbf{U}}_{N \times N} \underbrace{\mathbf{S}}_{N \times D} \underbrace{\mathbf{V}^T}_{D \times D} \quad (1149)$$

where the columns of  $\mathbf{U}$  are the left singular vectors, and the columns of  $\mathbf{V}$  are the right singular vectors. **Economy-sized SVD** will shrink  $\mathbf{U}$  to be  $N \times D$  and  $\mathbf{S}$  to be  $D \times D$  (we're assuming  $N > D$ ).

## Markov and Hidden Markov Models (Ch. 17)

Table of Contents Local

Written by Brandon McKinzie

Kevin P. Murphy (2012). Markov and Hidden Markov Models.

*Machine Learning: A Probabilistic Perspective.*

**Hidden Markov Models** (17.3). HMMs model a joint distribution over a sequence of  $T$  observations  $\mathbf{x}_{\langle 1 \dots T \rangle}$  and hidden states  $\mathbf{z}_{\langle 1 \dots T \rangle}$ ,

$$p(\mathbf{z}_{\langle 1 \dots T \rangle}, \mathbf{x}_{\langle 1 \dots T \rangle}) = p(\mathbf{z}_{\langle 1 \dots T \rangle})p(\mathbf{x}_{\langle 1 \dots T \rangle} | \mathbf{z}_{\langle 1 \dots T \rangle}) = \left[ p(z_1) \prod_{t=2}^T p(z_t | z_{t-1}) \right] \left[ \prod_{t=1}^T p(\mathbf{x}_t | z_t) \right] \quad (1150)$$

where each hidden state is discrete:  $z_t \in \{1, \dots, K\}$ , while each observation  $\mathbf{x}_t$  can be discrete or continuous.

**The Forwards Algorithm** (17.4.2). Goal: compute the filtered<sup>256</sup> marginals  $p(z_t | \mathbf{x}_{\langle 1 \dots t \rangle})$ .

1. **Prediction step.** Compute the *one-step-ahead predictive density*  $p(z_t | \mathbf{x}_{\langle 1 \dots t-1 \rangle})$ ,

$$p(z_t=j | \mathbf{x}_{\langle 1 \dots t-1 \rangle}) = \sum_i p(z_t=j, z_{t-1}=i | \mathbf{x}_{\langle 1 \dots t-1 \rangle}) \quad (1151)$$

which will serve as our prior for time  $t$  (since it does not take into account observed data at time  $t$ ).

2. **Update step.** We “update” our beliefs by observing  $\mathbf{x}_t$ ,

$$\alpha_t(j) \triangleq p(z_t=j | \mathbf{x}_{\langle 1 \dots t \rangle}) \quad (1152)$$

$$= \frac{p(z_t=j, \mathbf{x}_t, \mathbf{x}_{\langle 1 \dots t-1 \rangle})}{p(\mathbf{x}_{\langle 1 \dots t \rangle})} \quad (1153)$$

$$= \frac{\cancel{p(\mathbf{x}_{\langle 1 \dots t-1 \rangle})} p(z_t=j | \mathbf{x}_{\langle 1 \dots t-1 \rangle}) p(\mathbf{x}_t | z_t=j, \cancel{\mathbf{x}_{\langle 1 \dots t-1 \rangle}})}{\cancel{p(\mathbf{x}_{\langle 1 \dots t-1 \rangle})} p(\mathbf{x}_t | \mathbf{x}_{\langle 1 \dots t-1 \rangle})} \quad (1154)$$

$$= \frac{1}{Z_t} p(z_t=j | \mathbf{x}_{\langle 1 \dots t-1 \rangle}) p(\mathbf{x}_t | z_t=j) \quad (1155)$$

Notice that we can also use the values of  $Z_t$  to compute the log probability of the evidence:

$$\log p(\mathbf{x}_{\langle 1 \dots T \rangle} | \boldsymbol{\theta}) = \sum_{t=1}^T \log p(\mathbf{x}_t | \mathbf{x}_{\langle 1 \dots t-1 \rangle}) = \sum_{t=1}^T \log Z_t \quad (1156)$$

---

<sup>256</sup>They’re called “filtered” because they use all observations  $\mathbf{x}_{\langle 1 \dots t \rangle}$  instead of just  $\mathbf{x}_t$ , which reduces/filters out the noise more.

### Forwards Algorithm (Algorithm 17.1)

We are given transition matrix  $T_{i,j} = p(z_t = j \mid z_{t-1} = i)$ , evidence vectors  $\psi_t(j) = p(\mathbf{x}_t \mid z_t=j)$ , and initial state distribution  $\pi(j) = p(z_1=j)$ .

1. First, compute the initial  $[\alpha_1, Z_1] = \text{normalize}(\psi_1 \odot \pi)$ .
2. For time  $2 \leq t \leq T$ , compute  $[\alpha_t, Z_t] = \text{normalize}(\psi_t \odot \mathbf{T}^T \alpha_{t-1})$ .
3. Return  $\boldsymbol{\alpha}_{\langle 1 \dots T \rangle}$  and  $\log p(\mathbf{x}_{\langle 1 \dots T \rangle}) = \sum_t \log Z_t$ .

**The Forwards-Backwards Algorithm** (17.4.3). Goal: compute the smoothed marginals,  $p(z_t = j \mid \mathbf{x}_{\langle 1 \dots T \rangle})$ .

$$\gamma_t(j) \triangleq p(z_t = j \mid \mathbf{x}_{\langle 1 \dots T \rangle}) \quad (1157)$$

$$\propto p(z_t = j \mid \mathbf{x}_{\langle 1 \dots t \rangle}) p(\mathbf{x}_{\langle t+1 \dots T \rangle} \mid z_t = j) \quad (1158)$$

$$= \alpha_t(j) \beta_t(j) \quad (1159)$$

$$\text{where } \beta_t(j) \triangleq p(\mathbf{x}_{\langle t+1 \dots T \rangle} \mid z_t=j) \quad (1160)$$

$$= p(\mathbf{x}_{t+1}, \mathbf{x}_{\langle t+2 \dots T \rangle} \mid z_t=j) \quad (1161)$$

$$= \sum_i p(z_{t+1}=i, \mathbf{x}_{t+1}, \mathbf{x}_{\langle t+2 \dots T \rangle} \mid z_t=j) \quad (1162)$$

$$= \sum_i p(z_{t+1}=i, \mathbf{x}_{t+1} \mid z_t=j) p(\mathbf{x}_{\langle t+2 \dots T \rangle} \mid z_{t+1} \neq j, z_{t+1}=i, \mathbf{x}_{t+1}) \quad (1163)$$

$$= \sum_i p(z_{t+1}=i \mid z_t=j) p(\mathbf{x}_{t+1} \mid z_{t+1}=i) p(\mathbf{x}_{\langle t+2 \dots T \rangle} \mid z_{t+1}=i) \quad (1164)$$

$$= \sum_i p(z_{t+1}=i \mid z_t=j) p(\mathbf{x}_{t+1} \mid z_{t+1}=i) \beta_{t+1}(i) \quad (1165)$$

Using the same notation as Algorithm 17.1 above, the matrix-vector form for  $\beta$  is

$$\boldsymbol{\beta}_t = \mathbf{T}(\boldsymbol{\psi}_t \odot \boldsymbol{\beta}_{t+1}) \quad (1166)$$

with base case  $\boldsymbol{\beta}_T = \mathbf{1}$ .

**The Viterbi Algorithm** (17.4.4). Denote the [probability for] the most probable path leading to  $z_t=j$  as

$$\delta_t(j) \triangleq \max_{\mathbf{z}_{\langle 1 \dots t-1 \rangle}} p(\mathbf{z}_{\langle 1 \dots t-1 \rangle}, z_t=j \mid \mathbf{x}_{\langle 1 \dots t \rangle}) \quad (1167)$$

$$= \max_i \delta_{t-1}(i) \cdot T_{i,j} \cdot \psi_t(j) \quad (1168)$$

It is common to work in the log-domain when computing  $\delta$ .

with initialization of  $\delta_1(j) = \pi_j \psi_1(j)$ . We compute this until termination at  $z_T^* = \arg \max_i \delta_T(i)$ . Note the arg max here instead of a max – we keep track of both for all time steps. We do this so we can perform **traceback** to get the full most probable state sequence, starting at  $T$  and ending at  $t=1$ :

$$z_t^* = a_{t+1}(z_{t+1}^*) \quad (1169)$$

where  $a_t(j)$ , the most probable state at time  $t-1$  leading to state  $j$  at time  $t$ , is the same formula as  $\delta_t(j)$  but with an arg max.

## Undirected Graphical Models (Ch. 19)

Table of Contents Local

Written by Brandon McKinzie

Kevin P. Murphy (2012). Undirected Graphical Models.

*Machine Learning: A Probabilistic Perspective.*

**Learning** (19.5). Consider a MRF in log-linear form over  $C$  cliques and its log-likelihood (scaled by  $1/N$ ):

$$p(\mathbf{y} \mid \boldsymbol{\theta}) = \frac{1}{Z(\boldsymbol{\theta})} \exp \left\{ \sum_c \boldsymbol{\theta}_c^T \phi_c(\mathbf{y}) \right\} \quad (1170)$$

$$\ell(\boldsymbol{\theta}) = \frac{1}{N} \sum_i \left[ \sum_c \boldsymbol{\theta}_c^T \phi_c(\mathbf{y}_i) - \log Z(\boldsymbol{\theta}) \right] \quad (1171)$$

We know from chapter 9 that this log-likelihood is concave in  $\boldsymbol{\theta}$ , and that  $\frac{\partial}{\partial \boldsymbol{\theta}_c} \log Z = \mathbb{E} [\phi_c]$ . So the gradient of the LL is

$$\frac{\partial \ell}{\partial \boldsymbol{\theta}_c} = \left[ \frac{1}{N} \sum_i \phi_c(\mathbf{y}_i) \right] - \mathbb{E}_{p(\mathbf{y} \mid \boldsymbol{\theta})} [\phi_c(\mathbf{y})] \quad (1172)$$

Note that the first (“clamped”) term only needs to be computed once for any  $\mathbf{y}_i$ , it is completely independent of any parameters. It’s just evaluating feature functions, which e.g. for CRFs are often all indicator functions.

**CRF Training** (19.6.3). For [linear-chain] CRFs, the equations change slightly (but importantly).

$$\ell(\mathbf{w}) \triangleq \frac{1}{N} \sum_i \left[ \sum_c \mathbf{w}_c^T \phi_c(\mathbf{y}_i, \mathbf{x}_i) - \log Z(\mathbf{w}, \mathbf{x}_i) \right] \quad (1173)$$

$$\frac{\partial \ell}{\partial \mathbf{w}_c} = \frac{1}{N} \sum_i [\phi_c(\mathbf{y}_i, \mathbf{x}_i) - \mathbb{E}_{p(\mathbf{y} \mid \mathbf{x}_i)} [\phi_c(\mathbf{y}_i, \mathbf{x}_i)]] \quad (1174)$$

It’s important to recognize that the gradient of the log partition function now must be computed for *each* instance  $\mathbf{x}_i$ .

# CONVEX OPTIMIZATION

## CONTENTS

10.1 Convex Sets (Ch. 2) . . . . .	370
------------------------------------	-----

## Convex Sets (Ch. 2)

Table of Contents Local

Written by Brandon McKinzie

Boyd and Vandenberghe (2004). Convex Sets.

*Convex Optimization.*

## Lines and line segments

Viewed as a function of  $\theta \in \mathbb{R}$ , we can express the equation for a **line** in  $\mathbb{R}^n$  in the following two ways:

$$y = \theta x_1 + (1 - \theta)x_2 \quad (1175)$$

$$y = x_2 + \theta(x_1 - x_2) \quad (1176)$$

for some  $x_1, x_2 \neq x_1 \in \mathbb{R}^n$ . If we restrict  $\theta \in [0, 1]$ , we have a **line segment**.

## Affine sets

A set  $C \subseteq \mathbb{R}^n$  is **affine** if the line (not just segment) through any two distinct points in  $C$  also lies in  $C$ . More generally, this implies that for any set of points  $\{x_1, \dots, x_k\}$ , with each  $x_i \in C$ , all **affine combinations**,

$$\sum_{i=1}^k \theta_i x_i, \quad \text{where} \quad \sum_i \theta_i = 1 \quad (1177)$$

are in  $C$ , too. Related terminology:

- **affine hull** of any set  $C \subseteq \mathbb{R}^n$ , denoted  $\mathbf{aff}C$ , is the set of all affine combinations of points in  $C$ .
- **affine dimension** of a set  $C$  is the dimension of its affine hull.
- **relative interior** of a set  $C$ , denoted  $\mathbf{relint}C$ , is its interior<sup>257</sup> relative to  $\mathbf{aff}C$ ,

$$\mathbf{relint}C \triangleq \{x \in C \mid B(x, r) \cap \mathbf{aff}C \subseteq C \text{ for some } r > 0\} \quad (1178)$$

There's a lot of neat things to say here, but I only have space to state the results:

- If  $C$  is an affine set and  $x_0 \in C$ , then

$$V = C - x_0 \triangleq \{x - x_0 \mid x \in C\} \quad (1179)$$

is a **subspace**, i.e. closed under sums and scalar multiplication.

---

<sup>257</sup> The **interior of a set**  $C \subseteq \mathbb{R}^n$ , denoted  $\mathbf{int}C$ , is the set of all points interior to  $C$ . A point  $x \in C$  is interior to  $C$  if  $\exists \epsilon > 0$  for which *all points* in the set

$$\{y \in \mathbb{R}^n \mid \|y - x\|_2 \leq \epsilon\}$$

are also in  $C$ .

- The solution set of a system of linear equations,  $C = \{x \mid Ax = b\}$ , is an affine set. The subspace associated with  $C$  is the *nullspace* of  $A$ .

## Convex sets

A set  $C$  is **convex** if it contains all **convex combinations**,

$$\sum_{i=1}^k \theta_i x_i, \quad \text{where } \sum_i \theta_i = 1, \text{ and } \theta_i \geq 0 \quad (1180)$$

Related terminology:

- **convex hull** a set  $C$ , denoted  $\text{conv}C$ , is the set of all convex combinations of points in  $C$ .

## Cones

A set  $C$  is called a **cone** if  $(\forall x \in C)(\theta \geq 0)$  we have  $\theta x \in C$ . A set  $C$  is called a **convex cone** if it contains all **conic combinations**,

$$\sum_{i=1}^k \theta_i x_i, \quad \text{where } \theta_i \geq 0 \quad (1181)$$

of points in  $C$ . Related terminology:

- **conic hull** of a set  $C$  is the set of all conic combinations of points in  $C$ .

# BAYESIAN DATA ANALYSIS

## CONTENTS

11.1	Probability and Inference (Ch. 1) . . . . .	373
11.2	Single-Parameter Models (Ch. 2) . . . . .	375
11.3	Asymptotics and Connections to Non-Bayesian Approaches (Ch. 4) . . . . .	378
11.4	Gaussian Process Models (Ch. 21) . . . . .	381

## Probability and Inference (Ch. 1)

Table of Contents Local

Written by Brandon McKinzie

The process of Bayesian Data Analysis can be divided into the following 3 steps:

1. Setting up a *full probability model*.
2. Conditioning on observed data. Calculating the *posterior distribution* over unobserved quantities, given observed data.
3. Evaluating the fit of the model and implications of the posterior.

**Notation:** In general, we let  $\theta$  denote unobservable vector quantities or population *parameters* of interest, and  $y$  as collected data. This means our *posterior* takes the form  $p(\theta | y)$ , and our *likelihood* takes the form  $p(y | \theta)$ .

### Means and Variances of Conditional Distributions.

$$\mathbb{E}[u] = \mathbb{E}_v[\mathbb{E}_u[u | v]] \quad (1182)$$

$$\text{Var}[u] = \mathbb{E}_v[\text{Var}[u | v]] + \text{Var}[\mathbb{E}_u[u | v]] \quad (1183)$$

#### Proofs

$$\mathbb{E}[u] = \int dv \Pr[v] \int du \Pr[u | v] \quad (1184)$$

$$= \mathbb{E}_v[\mathbb{E}_u[u | v]] \quad (1185)$$

$$\text{Var}[u] = \mathbb{E}[u^2] - \mathbb{E}[u]^2 \quad (1186)$$

$$= \mathbb{E}_v[\mathbb{E}_u[u^2 | v]] - \mathbb{E}_v[\mathbb{E}_u[u | v]]^2 \quad (1187)$$

$$= \mathbb{E}_v[\mathbb{E}_u[u^2 | v] - \mathbb{E}_u[u | v]^2] + \mathbb{E}_v[\mathbb{E}_u[u | v]^2] - \mathbb{E}_v[\mathbb{E}_u[u | v]]^2 \quad (1188)$$

$$= \mathbb{E}_v[\text{Var}[u | v]] + \text{Var}[\mathbb{E}_u[u | v]] \quad (1189)$$

### Transformation of Variables.

$$\Pr_v[v] = |J| \Pr_u[f^{-1}(v)] \quad (1190) \quad J_{i,j} = \frac{\partial u}{\partial v} = \frac{\partial f^{-1}(v)}{\partial v}$$

where  $u$  and  $v$  have the same dimensionality, and  $|J|$  is the determinant of the Jacobian of the transformation  $u = f^{-1}(v)$ . When working with parameters defined on the open unit interval,  $(0, 1)$ , we often use the logistic transformation:

$$\text{logit}(u) = \log\left(\frac{u}{1-u}\right) \quad (1191)$$

$$\text{logit}^{-1}(v) = \log\left(\frac{e^v}{1+e^v}\right) \quad (1192)$$

## Standard Probability Distributions<sup>258</sup>.

Distribution	Notation	Density Function
Beta	Beta( $\alpha, \beta$ )	$p(\theta) = \frac{\Gamma(\alpha+\beta)}{\Gamma(\alpha)\Gamma(\beta)}\theta^{\alpha-1}(1-\theta)^{\beta-1}$
Inverse-gamma	Inv-gamma( $\alpha, \beta$ )	$p(\theta) = \frac{\beta^\alpha}{\Gamma(\alpha)}\theta^{-(\alpha+1)}e^{-\beta/\theta}$
Normal (univariate)	N( $\mu, \sigma^2$ )	$p(\theta) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{1}{2\sigma^2}(\theta - \mu)^2\right)$
Scaled inverse- $\chi^2$	Inv- $\chi^2(\nu, s^2)$	$\theta \sim \text{Inv-gamma}\left(\frac{\nu}{2}, \frac{\nu}{2}s^2\right)$

---

<sup>258</sup>The **gamma function**,  $\Gamma(x)$ , is defined as

$$\Gamma(x) = \int_0^\infty t^{x-1} e^{-t} dt \quad (1193)$$

or simply as  $(x-1)!$  if  $x \in \mathbb{Z}^+$ .

## Single-Parameter Models (Ch. 2)

Table of Contents Local

Written by Brandon McKinzie

Since we'll be referring to the binomial model frequently, below are the main distributions for reference:

$$p(y | \theta) = \text{Bin}(y | n, \theta) = \binom{n}{y} \theta^y (1 - \theta)^{n-y} \quad (1194)$$

$$p(\theta | y) \propto \theta^y (1 - \theta)^{n-y} \quad (1195)$$

$$\theta | y \sim \text{Beta}(y + 1, n - y + 1) \quad (1196)$$

where  $y$  is the number of successes out of  $n$  trials. We assume a uniform prior over the interval  $[0, 1]$ .

**Posterior as compromise between data and prior information.** Intuitively, the prior and posterior distributions over  $\theta$  should have some general relationship showing how the process of observing data  $y$  updates our distribution on  $\theta$ . We can use the identities from the previous chapter to see that

$$\mathbb{E} [\theta] = \mathbb{E}_y [\mathbb{E}_\theta [\theta | y]] \quad (1197)$$

$$\text{Var} [\theta] = \mathbb{E}_y [\text{Var} [\theta | y]] + \text{Var} [\mathbb{E}_\theta [\theta | y]] \quad (1198)$$

where:

- Equation 1197 states the obvious: our prior expectation for  $\theta$  is the expectation, taken over the distribution of possible data, of the posterior expectation.
- Equation 1198 states: *the posterior variance is on average smaller than the prior variance*, by an amount that depends on the variation [in posterior means] over the distribution of possible data. Stated another way: we can reduce our uncertainty with regard to  $\theta$  by larger amounts with models whose [expected] posteriors are strongly informed by the data.

**Informative Priors.** We now discuss some the issues that arise in assigning a prior distribution  $p(\theta)$  that reflects substantive information. Instead of using a uniform prior for our binomial model, let's explore the prior  $\theta \sim \text{Beta}(\alpha, \beta)$ <sup>259</sup>. Now our posterior takes the form

$$\theta | y \propto \text{Beta}(\alpha + y, \beta + n - y) \quad (1199)$$

The property that the posterior follows the same parametric form as the prior is called **conjugacy**.

---

<sup>259</sup> Assume we can select reasonable values for  $\alpha$  and  $\beta$ .

**Conjugacy, families, and sufficient statistics.** Formally, if  $\mathcal{F}$  is a class of sampling distributions  $\{\Pr_i[y | \theta]\}$ , and  $\mathcal{P}$  is a class of prior distributions,  $\{\Pr_j[\theta]\}$ , then the class  $\mathcal{P}$  is **conjugate** for  $\mathcal{F}$  if<sup>260</sup>

$$\forall \Pr[\cdot | \theta] \in \mathcal{F}, \Pr[\cdot] \in \mathcal{P} : \quad \Pr[\theta | y] \in \mathcal{P} \quad (1200)$$

Probability distributions that belong to an **exponential family** have natural conjugate prior distributions. The class  $\mathcal{F}$  is an exponential family if all its members have the form,

$$\Pr[y_i | \theta] = f(y_i)g(\theta)e^{\phi(\theta)^T u(y_i)} \quad (1201)$$

$$\Pr[y | \theta] = \left( \prod_{i=1}^n f(y_i) \right) g(\theta)^n \exp \left( \phi(\theta)^T \sum_{i=1}^n u(y_i) \right) \quad (1202)$$

$$\propto g(\theta)^n \exp \left( \phi(\theta)^T t(y) \right) \quad (1203)$$

where

- $y = (y_1, \dots, y_n)$  denotes  $n$  iid observations.
- $\phi(\theta)$  is called the **natural parameter** of  $\mathcal{F}$ .
- $t(y) = \sum_{i=1}^n u(y_i)$  is said to be a **sufficient statistic** for  $\theta$ , because the likelihood for  $\theta$  depends on the data  $y$  only through the value of  $t(y)$ .

**Normal distribution with known variance**<sup>261</sup>. Consider a single scalar observation  $y$  drawn from  $N(\theta, \sigma^2)$ , where we assume  $\sigma^2$  is known. The family of conjugate prior densities for the Gaussian likelihood as well as our choice of parameterization are, respectively,

$$p(\theta) \propto e^{A\theta^2 + B\theta + C} \quad (1205) \quad \text{Defining our conjugate prior}$$

$$p(\theta) \propto e^{-\frac{1}{2\tau_0^2}(\theta - \mu_0)^2} \quad (1206)$$

By definition, this implies that the posterior should also be normal. Indeed, after some basic arithmetic/substitutions, we find

$$\Pr[\theta | y] \propto \exp \left( -\frac{1}{2\tau_1^2}(\theta - \mu_1)^2 \right) \quad (1207) \quad \text{Writing our posterior precision and mean.}$$

$$\mu_1 = \frac{\frac{1}{\tau_0^2}\mu_0 + \frac{1}{\sigma^2}y}{\frac{1}{\tau_0^2} + \frac{1}{\sigma^2}} \quad \text{and} \quad \frac{1}{\tau_1^2} = \frac{1}{\tau_0^2} + \frac{1}{\sigma^2} \quad (1208)$$

---

<sup>260</sup>In English: A class of prior distributions is conjugate for a class of sampling distributions if, for any pair of sampling distribution and prior distribution [from those two respective classes], the associated posterior distribution is *also* in the same class of prior distributions.

<sup>261</sup>The following will be useful to remember:

$$\int_{-\infty}^{\infty} e^{-ax^2 + bx + c} dx = \frac{\pi}{a} e^{\frac{b^2}{4a} + c} \quad (1204)$$

where we see that the posterior **precision** (inverse of variance) equals the prior prior precision plus the data precision. We can see the posterior mean  $\mu_1$  expressed as a weighted average of the prior mean and the observed value<sup>262</sup>  $y$ , with weights proportional to the precisions.

**Normal distribution with unknown variance.** Now, we assume the mean  $\theta$  is known, and the variance  $\sigma^2$  is unknown. The likelihood for a vector  $y = (y_1, \dots, y_n)$  of  $n$  iid observations is

$$\Pr[y | \sigma^2] \propto (\sigma^2)^{-n/2} \exp\left(-\frac{n}{2\sigma^2}v\right) \quad (1209)$$

$$v := \frac{1}{n} \sum_{i=1}^n (y_i - \theta)^2 \quad (1210)$$

Computing our likelihood for  $n$  IID observations.

where  $v$  is the sufficient statistic. The corresponding conjugate prior density is the inverse-gamma. This and our choice for parameterization (how we define  $\alpha$  and  $\beta$ ) are, respectively,

$$\Pr[\sigma^2] \propto (\sigma^2)^{-(\alpha+1)} e^{-\beta/\sigma^2} \quad (1211)$$

$$\sigma^2 \sim \text{Inv-}\chi^2(\nu_0, \sigma_0^2) = \text{Inv-gamma}(\frac{\nu_0}{2}, \frac{\nu_0}{2}\sigma_0^2) \quad (1212)$$

Defining our conjugate prior.

All that's left is computing our posterior,

$$p(\sigma^2 | y) \propto p(\sigma^2)p(y | \sigma^2) \quad (1213)$$

$$\sigma^2 | y \sim \text{Inv-}\chi^2\left(\nu_0 + 2, \frac{\nu_0\sigma_0^2 + nv}{\nu_0 + n}\right) \quad (1214)$$

**Jeffrey's Invariance Principle.** An approach for defining noninformative prior distributions. Let  $\phi = h(\theta)$ , where the function  $h$  is one-to-one. By transformation of variables,

$$p(\phi) = p(\theta) \left| \frac{d\theta}{d\phi} \right| = p(\theta) |h'(\theta)|^{-1} \quad (1215)$$

Jeffrey's principle is to basically take the above equation as a true equivalence – there is no difference between finding  $p(\theta)$  and applying the equation above [to get  $p(\phi)$ ] and directly finding  $p(\phi)$ .

Let  $J(\theta)$  denote the **Fisher information** for  $\theta$ , defined as

$$J(\theta) = \mathbb{E} \left[ \left( \frac{d \log \Pr[y | \theta]}{d\theta} \right)^2 | \theta \right] = -\mathbb{E} \left[ \frac{d^2 \log \Pr[y | \theta]}{d\theta^2} | \theta \right] \quad (1216)$$

Jeffrey's prior model defines the noninformative prior density as  $\Pr[\theta] \propto [J(\theta)]^{1/2}$ . We can work out that this model is indeed invariant to parameterization<sup>263</sup>.

---

<sup>262</sup>For now, we're considering the single data point case.

<sup>263</sup>Evaluate  $J(\phi)$  at  $\theta = h^{-1}(\phi)$ . You should find that  $J(\phi)^{1/2} = J(\theta)^{1/2} \left| \frac{d\theta}{d\phi} \right|$

## Asymptotics and Connections to Non-Bayesian Approaches (Ch. 4)

Table of Contents Local

Written by Brandon McKinzie

**Normal Approximations to the Posterior Distribution.** If the posterior  $\Pr[\theta | y]$  is unimodal and roughly symmetric, it can be convenient to approximate it by a normal distribution. Here we'll consider a quadratic approximation via the Taylor series expansion up to second-order,

$$\log \Pr[\theta | y] = \log \Pr[\hat{\theta} | y] + \frac{1}{2}(\theta - \hat{\theta})^T \left[ \frac{d^2}{d\theta^2} \log \Pr[\theta | y] \right]_{\theta=\hat{\theta}} (\theta - \hat{\theta}) \quad (1217)$$

where  $\hat{\theta}$  is the posterior mode. The remainder terms of higher order fade in importance relative to the quadratic term when  $\theta$  is close to  $\hat{\theta}$  and  $n$  is large. We'd like to cast this into a normal distribution. First, let

$$I(\theta) \triangleq -\frac{d^2}{d\theta^2} \log \Pr[\theta | y] \quad (1218)$$

which we will refer to as the **observed information**. We can then rewrite our approximation as<sup>264</sup>

$$\Pr[\theta | y] \approx \mathcal{N}(\hat{\theta}, [I(\hat{\theta})]^{-1}) \quad (1221)$$

*Under the normal approximation, the posterior distribution is summarized by its mode,  $\hat{\theta}$ , and the curvature of the posterior density,  $I(\hat{\theta})$ ; that is, asymptotically, these are sufficient statistics.*

---

<sup>264</sup>I also found it helpful to explicitly write the substitution after raising eq 1217 by power of  $e$  (all logs are assumed natural logs)

$$\Pr[\theta | y] = e^{\log \Pr[\hat{\theta} | y] - \frac{1}{2}(\theta - \hat{\theta})^T [I(\hat{\theta})]^{-1}(\theta - \hat{\theta})} \quad (1219)$$

$$= \Pr[\hat{\theta} | y] e^{-\frac{1}{2}(\theta - \hat{\theta})^T [I(\hat{\theta})]^{-1}(\theta - \hat{\theta})} \quad (1220)$$

**Example.** Let  $y_1, \dots, y_n$  be independent observations from  $\mathcal{N}(\mu, \sigma^2)$ . Define  $\theta := (\mu, \log \sigma)$  as the parameters of interest, and assume a uniform prior<sup>265</sup>  $\Pr[\theta]$ . Recall that (equation 3.2 in textbook)

$$\Pr[\theta = (\mu, \log \sigma) | y] \propto \sigma^{-(n+2)} \exp\left(-\frac{1}{2\sigma^2} \sum_{i=1}^n (y_i - \mu)^2\right) \quad (1223) \quad \text{Note that } \sum_{i=1}^n (y_i - \bar{y}) = 0$$

$$= \sigma^{-(n+2)} \exp\left(-\frac{1}{2\sigma^2} \left[n(\bar{y} - \mu)^2 + \sum_{i=1}^n (y_i - \bar{y})^2\right]\right) \quad (1224)$$

$$= \sigma^{-(n+2)} \exp\left(-\frac{1}{2\sigma^2} \left[n(\bar{y} - \mu)^2 + (n-1)s^2\right]\right) \quad (1225)$$

where  $s^2 = \frac{1}{n-1} \sum_{i=1}^n (y_i - \bar{y})^2$  is the sample variance of the  $y_i$ 's. The sufficient statistics are  $\bar{y}$  and  $s^2$ . To construct the approximation, we need the second derivatives of the log posterior density<sup>266</sup>,

$$\log \Pr[\mu, \log \sigma | y] = \text{const} - n \log \sigma - \frac{1}{2\sigma^2} \left(n(\bar{y} - \mu)^2 + (n-1)s^2\right) \quad (1226)$$

in order to compute  $I(\hat{\theta})$ . After computing first derivatives, we find that the posterior mode is

$$\hat{\theta} = (\hat{\mu}, \log \hat{\sigma}) = \left(\bar{y}, \log \left(\sqrt{\frac{n-1}{n}} s\right)\right) \quad (1227)$$

We then compute second derivatives and evaluate at  $\theta = \hat{\theta}$  to obtain  $I(\hat{\theta})$ . Combining all this into the final result:

$$\Pr[\mu, \log \sigma | y] \approx \mathcal{N}(\hat{\theta}, [I(\hat{\theta})]^{-1}) \quad (1228)$$

$$= \mathcal{N}\left(\begin{pmatrix} \bar{y} \\ \log \hat{\sigma} \end{pmatrix}, \begin{pmatrix} \hat{\sigma}^2/n & 0 \\ 0 & 1/(2n) \end{pmatrix}\right) \quad (1229)$$

where  $\hat{\sigma}^2/n$  is the variance along the  $\mu$  dimension, and  $1/(2n)$  is the variance along the  $\log \sigma$  direction. This example was just meant to illustrate, with a simple case, how we work through constructing the approximate normal distribution.

---

<sup>265</sup>Recall from Ch 3.2 that the uniform prior on  $\mu, \log \sigma$  is

$$\Pr[\mu, \sigma^2] \propto (\sigma^2)^{-1} \quad (1222)$$

where we continue to assume  $\mu$  is uniform in  $[0, 1]$  for some reason.

<sup>266</sup>I'm not sure why  $n+2$  has seemingly turned into  $n$ .

**Large-Sample Theory.** Asymptotic normality of the posterior distribution: as more data arrives from the same underlying distribution  $f(y)$ , the posterior distribution of the *parameter vector*  $\theta$  approaches multivariate normality, even if the true distribution of the data is not within the parametric family under consideration.

Suppose the data are modeled by a parametric family,  $\Pr[y | \theta]$ , with a prior distribution  $\Pr[\theta]$ . If the true distribution,  $f(y)$ , is included in the parametric family (i.e. if  $\exists \theta_0 : f(y) = \Pr[y | \theta_0]$ ), then it's also true that **consistency** holds<sup>267</sup>:  $\Pr[\theta | y]$  converges to a point mass at the true parameter value,  $\theta_0$  as  $n \rightarrow \infty$ .

---

<sup>267</sup>So, what if the true  $f(y)$  is *not* included in the parametric family? In that case, there is no longer a true value  $\theta_0$ , but its role in the theoretical result is replaced by a value  $\theta_0$  that makes the model distribution  $\Pr[y | \theta]$  closest to the true distribution  $f(y)$ , in a technical sense involving the **Kullback-Leibler divergence**.

## Gaussian Process Models (Ch. 21)

Table of Contents Local

Written by Brandon McKinzie

Since I like to begin by motivating what we're going to talk about, and since BDA doesn't really do this, I'm going to start with an excerpt from chapter 15 of Kevin Murphy's book:

*In supervised learning, we observe some inputs  $\mathbf{x}_i$  and some outputs  $y_i$ . We assume that  $y_i = f(\mathbf{x}_i)$ , for some unknown function  $f$ , possibly corrupted by noise. The optimal approach is to infer a distribution over functions given the data,  $p(f | \mathbf{X}, \mathbf{y})$ , and then to use this to make predictions given new inputs, i.e., to compute*

$$p(y_* | \mathbf{x}_*, \mathbf{X}, \mathbf{y}) = \int p(y_* | f, \mathbf{x}_*) p(f | \mathbf{X}, \mathbf{y}) df \quad (1230)$$

**Gaussian Processes** or GPs define a prior over functions  $p(f)$  which can be converted into a posterior over functions  $p(f | \mathbf{X}, \mathbf{y})$  once we've seen some data.

**Gaussian Process Regression** (20.1). We write a GP as  $\mu \sim \text{GP}(m, k)$  ( $\mu$  is now taking the place of  $f$  from Kevin Murphy's notation), parameterized in terms of a mean function  $m$  and a covariance function  $k$ . Remember,  $\mu$  is supposed to represent the predictor function for obtaining output predictions given inputs,  $y = \mu(x)$ . Instead of making the usual assumption that there is some "best" function  $\mu^*$  and trying to learn it via fitting a parameterized  $\hat{\mu}(\theta)$ , we are going *full meta*<sup>268</sup> (i.e. full Bayesian) and learning the distribution over predictors.

Apparently, we only need to consider a finite (and arbitrary) set of points  $x_1, \dots, x_n$  to consider when evaluating any given  $\mu$ . The **GP prior** on  $\mu$  is defined as

$$\mu(x_1), \dots, \mu(x_n) \sim \mathcal{N}(\{m(x_1), \dots, m(x_n)\}, K(x_1, \dots, x_n)) \quad (1231)$$

with mean  $m$  and covariance  $K$ <sup>269</sup>. The covariance function  $k$  specifies the covariance between the process at any two points, with  $K$  an  $n \times n$  covariance matrix with  $K_{p,q} = k(x_p, x_q)$ . The covariance function controls the smoothness of realizations from the GP<sup>270</sup> and the degree of shrinkage towards the mean.

<sup>268</sup> What if there are like, a whole space of different predictors, man? Like, what if there is an infinite sea of predictor functions, all with their own unique traits and quirks? Woah.

<sup>269</sup> Don't confuse the notation –  $\mathcal{N}$  uses covariance  $K$  as an argument, while  $GP$  uses covariance function  $k$  as an argument.

<sup>270</sup> In English: How similar we expect different samples of  $\mu$  to look as a function of  $x$ . The reason this was weird to think about at first is because I'm used to thinking about covariance/smoothness over  $x$  rather than sampled functions of  $x$ . Meta.

A common choice the squared exponential,

$$k(x, x') = \tau^2 \exp\left(-\frac{|x - x'|^2}{2\ell^2}\right) \quad (1232)$$

where  $\tau$  controls the magnitude and  $\ell$  the smoothness of the function.

# GAUSSIAN PROCESSES FOR MACHINE LEARNING

## CONTENTS

12.1 Regression (Ch. 2) . . . . .	384
-----------------------------------	-----

## Regression (Ch. 2)

Table of Contents Local

Written by Brandon McKinzie

Rasmussen and Williams (2006). Regression. *Gaussian Processes for Machine Learning*.

**Weight-space view** (2.1). We review the standard probabilistic view of linear regression.

$$f(\mathbf{x}) = \mathbf{x}^T \mathbf{w} \quad (1233)$$

$$y = f(\mathbf{x}) + \varepsilon \quad \text{where } \varepsilon \sim \mathcal{N}(0, \sigma_n^2) \quad (1234)$$

$$\begin{aligned} [\text{likelihood}] \quad & \mathbf{y} | X, \mathbf{w} \sim \mathcal{N}(X^T \mathbf{w}, \sigma_n^2 I) \\ [\text{prior}] \quad & \mathbf{w} \sim \mathcal{N}(\mathbf{0}, \Sigma_p) \end{aligned} \quad (1235) \quad X \in \mathbb{R}^{d \times n}$$

$$\begin{aligned} [\text{posterior}] \quad p(\mathbf{w} | \mathbf{y}, X) &= \frac{p(\mathbf{y} | \mathbf{w}, X)p(\mathbf{w})}{p(\mathbf{y} | X)} \\ &\sim \mathcal{N}\left(\frac{1}{\sigma_n^2} A^{-1} X \mathbf{y}, A^{-1}\right) \end{aligned} \quad (1237) \quad (1238)$$

where  $A = \sigma_n^{-2} X X^T + \Sigma_p^{-1}$ , and we often set  $\Sigma_p = I$ . When analyzing the contour plots for the likelihood, remember that it is *not* a probability distribution, but rather it's interpreted as a function of  $\mathbf{w}$ , i.e.  $\text{likelihood}(\mathbf{w}) := \mathcal{N}(\mathbf{y}; X^T \mathbf{w}, \sigma_n^2 I)$ .

Note that we often use Bayesian techniques without realizing it. For example, what does the following remind you of?

$$\ln p(\mathbf{w}) \propto \frac{1}{2} \mathbf{w}^T \Sigma_p \mathbf{w} \quad (1239)$$

It's the l2 penalty from ridge regression (where typically  $\Sigma_p = I$ ). We can also project the inputs to a higher-dimensional space, often referred to as the *feature space*, by passing them through feature functions  $\phi(x)$ . As we'll see later (ch 5), GPs actually tell us how to define the basis functions  $h_i(x)$  which define the value of  $\phi_i(x)$ , the  $i$ th element of the feature vector. The author then proceeds to give an overview of the kernel trick.

**Function-space view** (2.2). Instead of inference over parameters  $\mathbf{w}$ , we can equivalently consider inference in function space with a **Gaussian process** (GP), formally defined as

A Gaussian process is a collection of random variables, any finite number of which have a joint Gaussian distribution.

A GP is completely specified by its mean function and covariance function. Define mean function  $m(\mathbf{x})$  and covariance function  $k(\mathbf{x}, \mathbf{x}')$  of a real [Gaussian] process  $f(\mathbf{x})$  as

$$m(\mathbf{x}) = \mathbb{E}_f [f(\mathbf{x})] \quad (1240)$$

$$k(\mathbf{x}, \mathbf{x}') = \mathbb{E}_f [(f(\mathbf{x}) - m(\mathbf{x})) (f(\mathbf{x}') - m(\mathbf{x}'))] \quad (1241)$$

$$f(\mathbf{x}) \sim \mathcal{GP}(m(\mathbf{x}), k(\mathbf{x}, \mathbf{x}')) \quad (1242)$$

Note that the expectations are over the *random variable*  $f(\mathbf{x})$  for any *given* (non-random)  $\mathbf{x}$ . In other words, the expectation is over the space of possible functions, each evaluated at point  $\mathbf{x}$ . Concretely, this is often an expectation over the parameters  $\mathbf{w}$ , as is true for our linear regression example. We can now write our Bayesian linear regression model (with feature functions) as a GP.

$$\mathbb{E}_{\mathbf{w} \sim \mathcal{N}(\mathbf{0}, \Sigma_p)} [f(\mathbf{x}; \mathbf{w})] = \phi(\mathbf{x})^T \mathbb{E}[\mathbf{w}] = 0 \quad (1243)$$

$$\mathbb{E}_{\mathbf{w} \sim \mathcal{N}(\mathbf{0}, \Sigma_p)} [f(\mathbf{x}) f(\mathbf{x}')] = \phi(\mathbf{x})^T \mathbb{E}[\mathbf{w} \mathbf{w}^T] \phi(\mathbf{x}') = \phi(\mathbf{x})^T \Sigma_p \phi(\mathbf{x}') \quad (1244)$$

A common covariance function is the **squared exponential** (a.k.a. the RBF kernel),

$$k(\mathbf{x}_p, \mathbf{x}_q) = \text{Cov}[f(\mathbf{x}_p), f(\mathbf{x}_q)] = \exp(-\frac{1}{2} |\mathbf{x}_p - \mathbf{x}_q|^2) \quad (1245)$$

where it's important to recognize that, whereas we've usually seen this is the RBF kernel for the purposes of kernel methods on inputs, we are now using it to specify the *covariance of the outputs*.

Ok, so how do we sample some functions and plot them? Below is an overview for our current linear regression running example.

1. Choose a number of input points  $X_*$ . For our linear regression example, we could set this to `np.arange(-5, 5.1, 0.1)` to get evenly spaced  $x$  in  $[-5, 5]$  in intervals of 0.1.
2. Write out the covariance matrix defined by  $K_{p,q} = k(x_p, x_q)$  using our squared exponential covariance function, for all pairs of inputs.
3. We can now generate samples of function  $f$ , represented as a random vector with size equal to the number of inputs  $|X_*|$ , by sampling from the **GP prior**

$$\mathbf{f}_* \sim \mathcal{N}(\mathbf{0}, K(X_*, X_*)) \quad (1246)$$

So far, we've only dealt with the GP *prior*. What do we do when we get labeled training observations? How do we make predictions on unlabeled test data? Well, for the simple case where our observations are noise free<sup>271</sup>, that is we know  $\{(\mathbf{x}_i, f_i) \mid i = 1, \dots, n\}$ , the joint

---

<sup>271</sup>For example, noise-free linear regression would mean we model  $y=f(x)$ , implicitly defining  $\varepsilon = 0$ .

GP prior over the train set inputs  $X$  and test set inputs  $X_*$  is defined exactly how we did it earlier (zero mean, elementwise evaluation of  $k$ ). In other words, our GP prior models the train outputs  $\mathbf{f}$  and test outputs  $\mathbf{f}_*$  as random vectors sampled via

$$\begin{bmatrix} \mathbf{f} \\ \mathbf{f}_* \end{bmatrix} \sim \mathcal{N} \left( \mathbf{0}, \begin{bmatrix} K(X, X) & K(X, X_*) \\ K(X_*, X) & K(X_*, X_*) \end{bmatrix} \right) \quad (1247)$$

You may be wondering: *why are we talking about sampling from the prior on inputs? We already know the outputs!*, and you'd be correct. The way we obtain our posterior is by restricting our joint prior to only those functions that agree with the observed training data  $X, \mathbf{f}$ , which we can do by simply conditioning on them. Our posterior for sampling test outputs given test inputs  $X_*$  is thus

$$\begin{aligned} \mathbf{f}_* | X_*, X, \mathbf{f} \sim \mathcal{N} \left( K(X_*, X)K(X, X)^{-1}\mathbf{f}, \right. \\ \left. K(X_*, X_*) - K(X_*, X)K(X, X)^{-1}K(X, X_*) \right) \end{aligned} \quad (1248)$$

# BLOGS

## CONTENTS

13.1	Conv Nets: A Modular Perspective . . . . .	388
13.2	Understanding Convolutions . . . . .	389
13.3	Deep Reinforcement Learning . . . . .	391
13.4	Deep Learning for Chatbots (WildML) . . . . .	393
13.5	Attentional Interfaces – Neural Perspective . . . . .	395

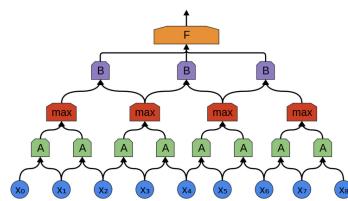
## Conv Nets: A Modular Perspective

Table of Contents Local

Written by Brandon McKinzie

From this post on Colah's Blog.

The title is inspired by the following figure. Colah mentions how groups of neurons, like  $A$ , that appear in multiple places are sometimes called **modules**, and networks that use them are sometimes called modular neural networks. You can feed the output of one convolutional layer into another. With each layer, the network can detect higher-level, more abstract features.



- Function of the  $A$  neurons: compute certain *features*.
- Max pooling layers: kind of “zoom out”. They allow later convolutional layers to work on larger sections of the data. They also make us invariant to some very small transformations of the data.

## Understanding Convolutions

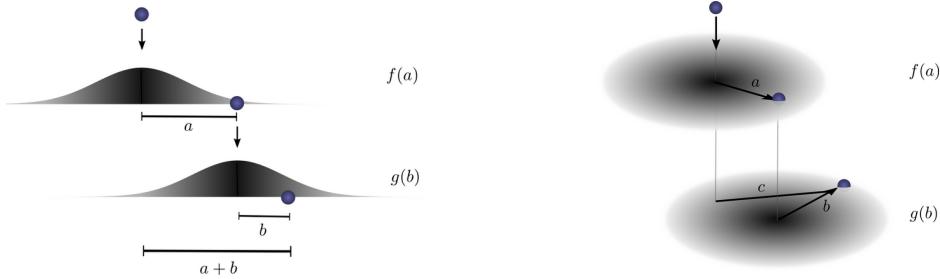
[Table of Contents](#)   [Local](#)

*Written by Brandon McKinzie*

From Colah's Blog.

**Ball-Dropping Example.** The posed problem:

Imagine we drop a ball from some height onto the ground, where it only has one dimension of motion. How likely is it that a ball will go a distance  $c$  if you drop it and then drop it again from above the point at which it landed?



From basic probability, we know the result is a sum over possible outcomes, constrained by  $a + b = c$ . It turns out this is actually the definition of the convolution of  $f$  and  $g$ .

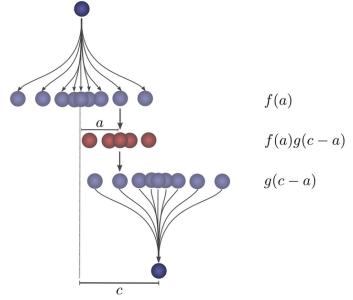
$$\Pr(a + b = c) = \sum_{a+b=c} f(a) \cdot g(b) \quad (1249)$$

$$(f * g)(c) = \sum_{a+b=c} f(a) \cdot g(b) \quad (1250)$$

$$= \sum_a f(a) \cdot g(c - a) \quad (1251)$$

**Visualizing Convolutions.** Keeping the same example in the back of our heads, consider a few interesting facts.

- **Flipping directions.** If  $f(x)$  yields the probability of landing a distance  $x$  away from where it was dropped, what about the probability that it was dropped a distance  $x$  from where it *landed*? It is  $f(-x)$ .
- Above is a visualization of one term in the summation of  $(f * g)(c)$ . It is meant to show how we can move the bottom around to think about evaluating the convolution for different  $c$  values.



$$\begin{matrix} 0 & 0 & 0 & 0 & 0 \\ 0 & 1/9 & 1/9 & 1/9 & 0 \\ 0 & 1/9 & 1/9 & 1/9 & 0 \\ 0 & 1/9 & 1/9 & 1/9 & 0 \\ 0 & 0 & 0 & 0 & 0 \end{matrix}$$

$$\begin{matrix} 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 \\ 0 & -1 & 1 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 \end{matrix}$$

We can relate these ideas to image recognition. Below are two common kernels used to convolve images with.

On the left is a kernel for *blurring* images, accomplished by taking simple averages. On the right is a kernel for *edge detection*, accomplished by taking the difference between two pixels, which will be largest at edges, and essentially zero for similar pixels.

## Deep Reinforcement Learning

[Table of Contents](#)   [Local](#)

*Written by Brandon McKinzie*

Link to tutorial – Part I of “Demystifying deep reinforcement learning.”

**Reinforcement Learning.** Vulnerable to the *credit assignment problem* - i.e. unsure which of the preceding actions was responsible for getting some reward and to what extent. Also need to address the famous *explore-exploit dilemma* when deciding what strategies to use.

**Markov Decision Process.** Most common method for representing a reinforcement problem. MDPs consist of states, actions, and rewards. Total reward is sum of current (includes previous) and *discounted* future rewards:

$$R_t = r_t \gamma (r_{t+1} + \gamma (r_{t+2} + \dots)) = r_t + \gamma R_{t+1} \quad (1252)$$

**Q - learning.** Define function  $Q(s, a)$  to be best possible score at end of game after performing action  $a$  in state  $s$ ; the “quality” of an action from a given state. The recursive definition of Q (for one transition) is given below in the *Bellman equation*.

$$Q(s, a) = r + \gamma \max_{a'} Q(s', a')$$

and updates are computed with a learning rate  $\alpha$  as

$$Q(s_t, a_t) = (1 - \alpha) \cdot Q(s_{t-1}, a_{t-1}) + \alpha \cdot (r + \gamma \max_{a'} Q(s'_{t+1}, a'_{t+1}))$$

**Deep Q Network.** Deep learning can take deal with issues related to prohibitively large state spaces. The implementation chosen by DeepMind was to represent the Q-function with a neural network, with the states (pixels) as the input and Q-values as output, where the number of output neurons is the number of possible actions from the input state. We can optimize with simple squared loss:

$$L = \frac{1}{2} [\underbrace{r + \max_{a'} Q(s', a')}_{\text{target}} - \underbrace{Q(s, a)}_{\text{prediction}}]^2$$

and our algorithm from some state  $s$  becomes

1. **First forward pass** from  $s$  to get all predicted Q-values for each possible action. Choose action corresponding to max output, leading to next  $s'$ .

2. **Second forward pass** from  $s'$  and again compute  $\max_{a'} Q(s', a')$ .
3. **Set target output** for each action  $a'$  from  $s'$ . For the action corresponding to max (from step 2) set its target as  $r + \gamma \max_{a'} Q(s', a')$ , and for all other actions set target to same as originally returned from step 1, making the error 0 for those outputs. (Interpret as update to our guess for the best Q-value, and keep the others the same.)
4. **Update weights** using backprop.

**Experience Replay.** This the most important trick for helping convergence of Q-values when approximating with non-linear functions. During gameplay all the experience  $< s, a, r, s' >$  are stored in a replay memory. When training the network, random minibatches from the replay memory are used instead of the most recent transition.

**Exploration.** One could say that initializing the Q-values randomly and then picking the max is essentially a form of exploitation. However, this type of exploration is *greedy*, which can be tamed/fixed with  **$\epsilon$ -greedy exploration**. This incorporates a degree of randomness when choosing next action at *all* time-steps, determined by probability  $\epsilon$  that we choose the next action randomly. For example, DeepMind decreases  $\epsilon$  over time from 1 to 0.1.

### Deep Q-Learning Algorithm.

```

initialize replay memory D
initialize action-value function Q with random weights
observe initial state s
repeat
    select an action a
    with probability  $\epsilon$  select a random action
    otherwise select  $a = \text{argmax}_{a'} Q(s, a')$ 
    carry out action a
    observe reward r and new state  $s'$ 
    store experience  $< s, a, r, s' >$  in replay memory D

    sample random transitions  $< ss, aa, rr, ss' >$  from replay memory D
    calculate target for each minibatch transition
        if  $ss'$  is terminal state then  $tt = rr$ 
        otherwise  $tt = rr + \gamma \max_{a'} Q(ss', aa')$ 
    train the Q network using  $(tt - Q(ss, aa))^2$  as loss

     $s = s'$ 
until terminated

```

## Deep Learning for Chatbots (WildML)

Table of Contents   Local

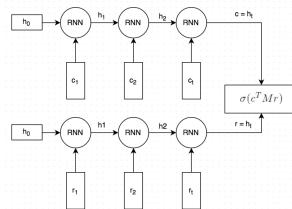
*Written by Brandon McKinzie*

### Overview.

- **Model.** Implementing a retrieval-based model. Input: conversation/context  $c$ . Output: response  $r$ .
- **Data.** Ubuntu Dialog Corpus (UDC). 1 million examples of form (context, utterance, label). The label can be 1 (utterance was actual response to the context) or a 0 (utterance chosen randomly). Using NLTK, the data has been . . .
  - **Tokenized.** dividing strings into lists of substrings.
  - **Stemmed. IDK**
  - **Lemmatized. IDK**

The test/validation set consists (context, ground-truth utterance, [9 distractors (incorrect utterances)]). The distractors are picked at random<sup>272</sup>

### Dual-Encoder LSTM.



1. **Inputs.** Both the context and the response text are split by words, and each word is embedded into a vector and fed into the same RNN.
2. **Prediction.** Multiply the [vector representation ("meaning")]  $c$  with param matrix  $M$  to predict some response  $r'$ .
3. **Evaluation.** Measure similarity of predicted  $r'$  to actual  $r$  via simple dot product. Feed this into sigmoid to obtain a probability [of  $r'$  being the correct response]. Use (binary) cross-entropy for loss function:

$$L = -y \cdot \ln(y') - (1 - y) \cdot \ln(1 - y') \quad (1253)$$

where  $y'$  is the predicted probability that  $r'$  is correct response  $r$ , and  $y \in \{0, 1\}$  is the true label for the context-response pair ( $c, r$ ).

---

<sup>272</sup>Better example/approach: Google's Smart Reply uses clustering techniques to come up with a set of possible responses.

**Data Pre-Processing.** Courtesy of WildML, we are given 3 files after preprocessing: train.tfrecords, validation.tfrecords, and test.tfrecords, which use TensorFlow's 'Example' format. Each Example consists of . . .

- context: Sequence of word ids.
- context\_len: length of the aforementioned sequence.
- utterance: seq of word ids representing utterance (response).
- utterance\_len.
- label: only in training data. 0 or 1.
- distractor\_[N]: Only in test/validation. N ranges from 0 to 8. Seq of word ids reppin the distractor utterance.
- distractor\_[N]\_len.

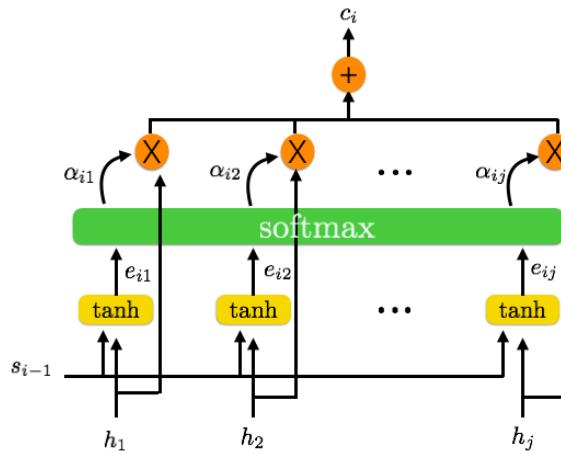
## Attentional Interfaces – Neural Perspective

Table of Contents Local

Written by Brandon McKinzie

[\[Link to article\]](#)

**Attention Mechanism.** Below is a close-up view/diagram of an attention layer. Technically, it only corresponds to a single time step  $i$ ; we are using the previous decoder state  $s_{i-1}$  to compute the  $i$ th context vector  $c_i$  which will be fed as an input to the decoder for step  $i$ .



For convenience, I'll rewrite the familiar equations for computing quantities at some step  $i$ .

$$\text{decoder state} \quad s_i = f(s_{i-1}, y_{i-1}, c_i) \quad (1254)$$

$$\text{context vect} \quad c_i = \sum_{j=1}^{T_x} \alpha_{ij} h_j \quad (1255)$$

$$\begin{aligned} \text{annotation weights} \quad \alpha_{ij} &= \frac{\exp(e_{ij})}{\sum_{k=1}^{T_x} \exp(e_{ik})} \\ e_{ij} &= a(s_{i-1}, h_j) \end{aligned} \quad (1256) \quad (1257)$$

Now we can see just how simple this really is. Recall that Bahdanau *et al.*, 2015 use the wording: “ $e_{ij}$  is an alignment model which scores how well the inputs around position  $j$  and the output at position  $i$  match.” But we can see an example implementation of an alignment model above: the  $\tanh$  function (that’s it).

# APPENDIX

## CONTENTS

14.1	Common Distributions and Models . . . . .	397
14.2	Math . . . . .	399
14.3	Matrix Cookbook . . . . .	408
14.4	Main Tasks in NLP . . . . .	409
14.5	Misc. Topics . . . . .	412
14.5.1	BLEU Score . . . . .	412
14.5.2	Connectionist Temporal Classification (CTC) . . . . .	413
14.5.3	Perplexity . . . . .	415
14.5.4	Byte Pair Encoding . . . . .	417
14.5.5	Grammars . . . . .	417
14.5.6	Bloom Filter . . . . .	418
14.5.7	Distributed Training . . . . .	418
14.5.8	Traditional Language Modeling . . . . .	419

## Common Distributions and Models

Table of Contents   Local

*Written by Brandon McKinzie*

### Continuous Distributions.

Distributions with support  $\theta > 0$ :

$$(\forall n \in \mathbb{Z}^+) \Gamma(n) = (n-1)!$$

Distribution	Density Function	Notation
Chi-Square	$p(\theta) = \frac{2^{-\nu/2}}{\Gamma(\nu/2)} \theta^{\nu/2-1} e^{-\theta/2}$	$\theta \sim \chi_\nu^2$
Gamma	$p(\theta) = \frac{\beta^\alpha}{\Gamma(\alpha)} \theta^{\alpha-1} e^{-\beta\theta}$	$\theta \sim \text{Gamma}(\alpha, \beta)$
Inverse-gamma	$p(\theta) = \frac{\beta^\alpha}{\Gamma(\alpha)} \theta^{-\alpha-1} e^{-\beta/\theta}$	$\theta \sim \text{Inv-gamma}(\alpha, \beta)$
Inverse-chi-square	$p(\theta) = \frac{2^{-\nu/2}}{\Gamma(\nu/2)} \theta^{-\nu/2-1} e^{-\frac{1}{2\theta}}$	$\theta \sim \text{Inv-}\chi_\nu^2$

$$\Gamma(z) = \int_0^\infty x^{z-1} e^{-x} dx$$

Distributions with support  $\theta \in [0, 1]$ :

Distribution	Density Function	Notation
Beta	$p(\theta) = \frac{\Gamma(\alpha+\beta)}{\Gamma(\alpha)+\Gamma(\beta)} \theta^{\alpha-1} (1-\theta)^{\beta-1}$	$\theta \sim \text{Beta}(\alpha, \beta)$
Dirichlet	$p(\theta) = \frac{\Gamma(\sum_k \alpha_k)}{\prod_k \Gamma(\alpha_k)} \prod_k \theta_k^{\alpha_k-1}, \quad \sum_k \theta_k = 1$	$\theta \sim \text{Dirichlet}(\alpha_1, \dots, \alpha_K)$

### Discrete Distributions.

Distribution	Density Function	Notation
Bernoulli	$p(x; \theta) = \theta^{\mathbb{1}_{x=1}} (1-\theta)^{\mathbb{1}_{x=0}}$	$X \sim \text{Ber}(\theta)$
Binomial	$p(x; n) = \binom{n}{x} p^x (1-p)^{n-x}$	$x \sim \text{Bin}(n, p)$
Multinomial	$p(x_1, \dots, x_k; n) = \frac{n!}{\prod_i x_i!} \prod_i p_i^{x_i}$	

**Logistic Regression.** Perhaps the simplest *linear method*<sup>273</sup> for classification is **logistic regression**. Let  $K$  be the number of classes that  $y$  can take on. The model is defined as

$$\Pr[y = k \mid \mathbf{x}] = \frac{\exp(\boldsymbol{\theta}_k^T \mathbf{x})}{1 + \sum_{\ell=1}^{K-1} \exp(\boldsymbol{\theta}_\ell^T \mathbf{x})}, \quad \text{for } 1 \leq k \leq K-1 \quad (1258)$$

$$\Pr[y = K \mid \mathbf{x}] = \frac{1}{1 + \sum_{\ell=1}^{K-1} \exp(\boldsymbol{\theta}_\ell^T \mathbf{x})} \quad (1259)$$

and we often denote  $\Pr[y = k \mid \mathbf{x}]$  under the entire set of parameters  $\boldsymbol{\theta}$  simply as  $p_k(\mathbf{x}; \boldsymbol{\theta})$  or just  $p_k(\mathbf{x})$ . The decision boundaries are the set of points in the domain of  $\mathbf{x}$  for which some  $p_k(\mathbf{x}) = p_{j \neq k}(\mathbf{x})$ . Equivalently, the model can be specified by  $K-1$  log-odds or logit

---

<sup>273</sup>We say a classification method is *linear* if its **decision boundary** is linear.

transformations of the form

$$\log \left( \frac{p_i(\mathbf{x})}{p_K(\mathbf{x})} \right) = \boldsymbol{\theta}_i^T \mathbf{x} \quad \text{for } 1 \leq i \leq K - 1 \quad (1260)$$

Also note that the parameter vectors are orthogonal to the K-1 decision boundaries. For any  $x, x'$  on the decision boundary defined as the set of points  $\{x : p_a(x) = p_b(x)\}$ , we know that the vector  $x - x'$  is parallel to the decision boundary (by definition), and can derive

$$\frac{p_a(x)}{p_b(x)} = 1 = \exp(\boldsymbol{\theta}_a^T x - \boldsymbol{\theta}_b^T x) \implies \boldsymbol{\theta}_a^T x = \boldsymbol{\theta}_b^T x \quad (1261)$$

$$\therefore \boldsymbol{\theta}_a^T (x - x') = \boldsymbol{\theta}_b^T (x - x') = 0 \quad (1262)$$

and thus  $\boldsymbol{\theta}_a$  and  $\boldsymbol{\theta}_b$  are both perpendicular to the decision boundary where  $p_a(x) = p_b(x)$ .

# Math

Table of Contents   Local

Written by Brandon McKinzie

Fancy math definitions/concepts for fancy authors who require fancy terminology.

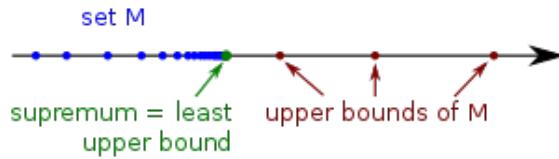
- **Support.** Sloppy definition you'll see in most places: The *set-theoretic support* of a real-valued function  $f : X \mapsto \mathbb{R}$  is defined as

$$\text{supp}(f) \triangleq \{x \in X : f(x) \neq 0\}$$

Note that Wikipedia is surprisingly sloppy with how it defines and/or uses support in various articles. After some digging, I finally found the formal definition for probability and measure theory:

*If  $X : \Omega \mapsto \mathbb{R}$  is a random variable on  $(\Omega, \mathcal{F}, P)$ , then the **support** of  $X$  is the smallest closed set  $R_X \subset \mathbb{R}$  such that  $P(X \in R_X) = 1$ .*

- **Infimum and Supremum.** The fancy-math way of saying minimum and maximum. Yes, I recognize that these are important in certain (usually rather abstract) settings, but often in ML it is used when sup means exactly the same thing as max, but the authors want to look sophisticated. Here I'll give the formal definition for sup. You have a partially ordered set<sup>274</sup>  $P$ , and are for the moment checking out some subset  $S \subseteq P$ . Someone asks you, "hey, give me an *upper bound* of  $S$ ." You just gotta find *some*  $b \in P$  (the larger/global set) that is greater than or equal to every element in  $S$ . The person then comes back and says "ok no, I need the *supremum* of  $S$ ." Now you need to find the *smallest* value out of all the possible upper bounds.



Hopefully it is clear why this is only relevant in real-valued cases where the “edges” aren’t well-defined.

- **Probability Measure.** Informal definition: a probability distribution<sup>275</sup>. Formal definition: a function  $\mu : \alpha \mapsto \mathbb{R}[0, 1]$  from events to scalar values, where  $\mu(\alpha) = 1$  if  $\alpha = \Omega$  (the full space) and  $\mu(\emptyset) = 0$ . Also  $\mu$  must satisfy the countable additivity property:  $\mu(\cup_i \alpha_i) = \sum_i \mu(\alpha_i)$  for pairwise disjoint sets  $\{\alpha\}_i$ .

<sup>274</sup>A partially ordered set  $(P, \leq)$  is a set of elements such that element  $i$  is less than or equal to element  $i + i$ .

<sup>275</sup>See this great answer detailing how the difference between “measure” and “distribution” is basically just context.

**Linear Algebra.** Feeling like I need a quick recap from my adv. linalg course and an area where I can ramble my interpretations. In what follows, let  $V$  and  $W$  be vector spaces over some field  $F$ .

### Linear Transformation

A function  $T : V \mapsto W$  is called a *linear transformation* from  $V$  to  $W$  if  $\forall x, y \in V$  and  $\forall c \in F$ :

- $T(x + y) = T(x) + T(y)$ .
- $T(cx) = cT(x)$ .

Now suppose that  $V$  and  $W$  have ordered bases  $\beta = \{v_1, \dots, v_n\}$  and  $\gamma = \{w_1, \dots, w_m\}$ , respectively. Then for each basis vector  $v_j$ , there exist unique scalars  $a_{ij} \in F$  such that

$$T(v_j) = \sum_{i=1}^m a_{ij} w_i \quad (1263)$$

Remember that each  $v_j$  and  $w_i$  are members of a *vector space* (they are not scalars). And also be careful to not associate the representation of any vector space element with its coordinate vector relative to a specific ordered basis, which *itself is a different linear transformation from some  $V \mapsto F^n$* . Keep it abstract.

### Matrix Representation of a Linear Transformation

We call the  $m \times n$  matrix  $A$  defined by  $A_{ij} = a_{ij}$  the *matrix representation of  $T$*  in the ordered bases  $\beta$  and  $\gamma$  and write  $A = [T]_{\beta}^{\gamma}$ . If  $V = W$  and  $\beta = \gamma$ , then  $A = [T]_{\beta}$ .

Given this definition, I think it's wise to interpret matrix representations by the column-vector point of view. Each column vector  $[T(v_j)]_{\gamma}$ , read as “the coordinate vector of  $T(v_j)$  relative to ordered basis  $\gamma$ ,” tells you how each basis vector  $v_j$  in domain  $V$  gets mapped to a [coordinate] vector in output space  $W$  [relative to a given ordered basis  $\gamma$ ]. For some reason, my brain has always had a hard time with the fact that the matrix row indices  $i$  correspond to output space, while the column indices  $j$  represent the input space. The easiest way (I think) to help ground this the right way is to remember that  $L_A(x) \triangleq Ax$ , i.e. the operator view. At the same time, notice how the effect of  $Ax$  is to take successive linear combinations over each element of  $x$ .

I just realized another reason why the interpretation felt backwards to my brain: when we are taught matrix multiplication, we do the computations  $Ax$  in our heads “row by row” along  $A$ , taking inner products of the row with  $x$ , so I've been taught to think of the rows as the main “units” of the matrix. I'm not sure how to fix this, but notice that the matrix is basically just a blueprint/roadmap/whatever you want to call it for taking coordinate vectors in one basis to coordinate vectors in another basis. It's really important to remember the coefficients of  $A$  are intimately tied to the input/output bases.

**AHA.** I've been thinking about this all wrong. For the longest time, I've been trying to force an interpretation of matrix multiplication that “feels” like scalar multiplication. I realize now that this is going about it all the wrong way. *Matrix multiplication need only be considered from*

*the lens of a linear transformation.* After all, that's exactly the purpose of matrices anyway! It's so glaringly obvious from the paragraphs above, but I guess I never took them seriously enough. Matrices are simply convenient ways for us to write down linear transformations on vectors in a given [ordered] basis. The jth column of the matrix defines how the original jth basis vector is transformed. **AHA** (again). Now I see why I missed this crucial connection – *everything above focuses on the formal definition of input ordered basis  $\beta$  to output ordered basis  $\gamma$ , but 99 percent of the time in real life we have either  $\beta \subset \gamma$  or  $\gamma \subset \beta$  (we usually are mapping from  $\mathbb{R}^n$  to  $\mathbb{R}^m$ ).* For example, let  $\mathbf{A} \in M^{m \times n}$  and  $\mathbf{x} \in \mathbb{R}^n$ ; the following is always true:

$$\mathbf{Ax} = L_A(\mathbf{x}) \quad (1264)$$

$$= \begin{bmatrix} L_A(\hat{\mathbf{e}}_1) & L_A(\hat{\mathbf{e}}_2) & \cdots & L_A(\hat{\mathbf{e}}_n) \end{bmatrix} \mathbf{x} \quad (1265)$$

$$= \sum_{i=1}^n L_A(\hat{\mathbf{e}}_i) x_i \quad (1266)$$

This viewpoint is painfully obvious to me now, but I guess I hadn't thought deeply enough about the implications of the definition of a linear transformation, and I *definitely* took the **representation** of a matrix way too seriously, rather than focusing on its **sole purpose**: provide a convenient way to write down linear transformations. For example, the above is actually a direct consequence of the definition of a L.T. itself:

$$L_A(\mathbf{x}) = L_A(x_1 \hat{\mathbf{e}}_1 + \cdots + x_n \hat{\mathbf{e}}_n) \quad (1267)$$

$$= x_1 L_A(\hat{\mathbf{e}}_1) + \cdots + x_n L_A(\hat{\mathbf{e}}_n) \quad (1268)$$

Time to really nail in the understanding. I also remember getting screwed up trying to think about *ok, so how do I conceptualize of the ith element of  $\mathbf{x}$  after the transformation? It's just a bunch of summed up goobley-gook!*. On one hand, yes that's true, but focus on the following before/after representations of  $\mathbf{x}$  to make your life easier:

$$\mathbf{x} \triangleq \sum_i^n x_i \hat{\mathbf{e}}_i \quad \xrightarrow{Ax} \quad L_A(\mathbf{x}) \triangleq \sum_i^n x_i L_A(\hat{\mathbf{e}}_i) \quad (1269)$$

**Matrix-Matrix Multiplication.** Continuing with the viewpoint that a matrix is nothing more than a convenient way to represent a linear transformation, recognize that any matrix multiplication  $AB$  represents a linear transformation itself, defined as  $T := T_A \circ T_B$ , the composition of  $A$  and  $B$ .

**Matrix Multiplication and Neural Networks.** Let's use the previous interpretations in the context of neural networks. A basic feedforward network with one hidden layer will compute outputs  $\mathbf{o}$  given inputs  $\mathbf{x}$ , each of which are vectors of possibly different dimension:

$$\mathbf{o}(\mathbf{x}) = \mathbf{W}^{(o)}\phi(\mathbf{x}) \quad (1270)$$

$$\phi(\mathbf{x}) = \mathbf{W}^{(h)}\mathbf{x} \quad (1271)$$

where  $\mathbf{W}^{(o)}$  and  $\mathbf{W}^{(h)}$  are the output and hidden parameter matrices, respectively. We already know that we can interpret each columns of these matrices as how the input basis vectors get mapped to the hidden or output space. However, since we usually think of the parameter matrices as representing the weighted edges of a network, we often think in terms of individual units. For example, the  $i$ th unit of the hidden layer vector  $\mathbf{h}$  is given by  $h_i = \sum_j^{n_{in}} W_{ij}x_j = \langle \mathbf{W}_{i,:}, \mathbf{x} \rangle$ . One interesting interpretation is that the  $i$ th element of  $\mathbf{h}$  is a **projection**<sup>276</sup> of the input  $\mathbf{x}$  onto the  $i$ th row of  $\mathbf{W}$ . This is of course true for any linear transformation; we can always think of the elements of the transformed vector as the result of projections of the original vector along a particular direction.

**Determinants.**  $\det A$  is the volume of space that a unit [hyper] cube is mapped to. Starting with the simplest non-trivial case, let  $A \in M^{2 \times 2}$ , and define  $A$  s.t. it simply scales the basis vectors (zero rotation). In other words,  $A_{i,j \neq i} := 0$ . In this case,  $\det A = a_{11}a_{22}$ , which is the area enclosed by the new scaled basis vectors. Skipping straight to the general case of [necessarily square] matrix  $A \in M^{n \times n}$  using Einstein summation notation and the Levi-Cevita symbol<sup>277</sup>:

$$\det A \triangleq \varepsilon_{i_1, \dots, i_n} a_{1,i_1} \dots a_{n,i_n} = \varepsilon_{i_1, \dots, i_n} \prod_{j=1}^n a_{j,i_j} \quad (1273)$$

Consider that if  $\det A = 0$ , then  $T_A$  “squishes” the volume of space in such a way that we essentially lose one or more dimensions. Notice how it only takes *one* lost dimension, since the volume of any region in  $\mathbb{R}^n$  is zero unless there is *some* amount in all dimensions (e.g. a cube with zero width has zero volume, regardless of its height/depth). It's also interesting to consider the relationships here with the invertible matrix theorem (defined a few paragraphs below). Having the intuition that determinants can be thought of as a change-in-volume makes it much more obvious why the equivalence statements of the invertible matrix theorem are indeed equivalent.

<sup>276</sup>This is informally worded. See the footnotes in the dot products section to understand why the element is technically just the result of a transformation (a projection would require re-mapping the scalar back to the space that  $\mathbf{W}_{i,:}$  lives (input space)).

<sup>277</sup>Recall that

$$\varepsilon_{i_1, \dots, i_n} \triangleq \begin{cases} +1 & \text{if } (i_1, \dots, i_n) \text{ is even perm of } (1, \dots, n) \\ -1 & \text{if } (i_1, \dots, i_n) \text{ is odd perm of } (1, \dots, n) \\ 0 & \text{otherwise} \end{cases} \quad (1272)$$

Note that this implies equal-to-zero if any of the indices are equal.

**Dot Products and Projections.** First, recall that a projection is *defined* to be a **linear** transformation  $P$  that is idempotent ( $P^n = P$  for  $n \geq 1$ ). Also, note that what you generally think of as a projection is technically an *orthogonal projection*<sup>278</sup>.

Here we'll show the intimate link between the dot product and [orthogonal] projection. Let  $P_{\mathbf{u}}$  define the [orthogonal] projection onto some **unit vector**  $\mathbf{u} \in \mathbb{R}^n$  (more generally, we could project onto a subspace instead of a single vector<sup>279</sup>). We can re-cast this as a linear transformation  $T_{\mathbf{u}} : \mathbb{R}^n \mapsto \mathbb{R}$  (technically not a *projection*, which would require re-mapping the output scalar back to  $\mathbb{R}^n$ ). We interpret the scalar output of  $T_{\mathbf{u}}(\mathbf{x})$  as the coordinate along the line spanned by  $\{\mathbf{u}\}$  that input vector  $\mathbf{x}$  gets mapped to. But wait, didn't we just talk a bunch about how to represent/conceptualize of the matrix representation of a transformation? Yes, we did. Well then, what would the matrix representation of  $T_{\mathbf{u}}$  look like? Don't forget that we've defined  $\|\mathbf{u}\| = 1$ .

$$[T_{\mathbf{u}}]_{\mathbb{R}^n}^{\mathbb{R}} = \begin{bmatrix} u_1 & \cdots & u_n \end{bmatrix} \quad (1274)$$

$$T_{\mathbf{u}}(\mathbf{x}) = \sum_i^n u_i x_i \quad (1275)$$

$$\longrightarrow = \mathbf{x} \cdot \mathbf{u} \quad (1276)$$

Furthermore, since linear transformations satisfy  $T(c\mathbf{x}) = cT(\mathbf{x})$  by definition, the final result is true even when  $\mathbf{u}$  is not a unit vector.

**Invertibility and Isomorphisms.** A function is invertible IFF it is both one-to-one and onto (i.e. bijective). Recall that  $\text{rank}(T)$  is the dimensionality of the range of  $T$ , which is the subspace of  $W$  that  $T$  maps to.

*Let  $T : V \mapsto W$  be a linear transformation, where  $V$  and  $W$  are finite-dimensional spaces of equal dimension. Then  $T$  is invertible IFF  $\text{rank}(T) = \dim(V)$ .*

For any invertible functions  $T$  and  $U$ :

- $(TU)^{-1} = U^{-1}T^{-1}$ . One easy way to show this is

$$(TU)U^{-1}T^{-1}(x) = T(UU^{-1})T^{-1}(x) = TT^{-1}(x) = x \quad (1277)$$

- $(T^{-1})^{-1} = T$ . The inverse of  $T$  is itself invertible.

---

<sup>278</sup>The more general definition uses wording “projection of vector  $\mathbf{x}$  *along*  $\mathbf{k}$  onto  $\mathbf{m}$ ”, where the distinction is shown in italics. An orthogonal projection implicitly defines  $\mathbf{k}$  as its null space; for any  $\alpha \in \mathbb{R}$ , an orthogonal projection satisfies  $P(\alpha\mathbf{k}) = \mathbf{0}$

<sup>279</sup>And technically, you don't project onto a *vector*, but rather you project onto a *line*, which is itself technically a subspace of 1 dimension. Yada yada yada.

## Inverse of a partitioned matrix

Consider a general partitioned matrix,

$$\mathbf{M} = \begin{pmatrix} \mathbf{E} & \mathbf{F} \\ \mathbf{G} & \mathbf{H} \end{pmatrix} \quad (1278)$$

where  $\mathbf{E}$  and  $\mathbf{H}$  are invertible. Then

$$\mathbf{M}^{-1} = \begin{pmatrix} (\mathbf{M}/\mathbf{H})^{-1} & -(\mathbf{M}/\mathbf{H})^{-1} \mathbf{F} \mathbf{H}^{-1} \\ \mathbf{E}^{-1} + \mathbf{E}^{-1} \mathbf{F} (\mathbf{M}/\mathbf{E})^{-1} \mathbf{G} \mathbf{E}^{-1} & (\mathbf{M}/\mathbf{E})^{-1} \end{pmatrix} \quad (1279)$$

$$\text{where } \mathbf{M}/\mathbf{H} \triangleq \mathbf{E} - \mathbf{F} \mathbf{H}^{-1} \mathbf{G} \quad (1280)$$

$$\mathbf{M}/\mathbf{E} \triangleq \mathbf{H} - \mathbf{G} \mathbf{E}^{-1} \mathbf{F} \quad (1281)$$

where  $\mathbf{M}/\mathbf{H}$  denotes the **Schur complement** of  $\mathbf{M}$  w.r.t.  $\mathbf{H}$ .

## Invertible Matrix Theorem

Let  $\mathbf{A}$  be a square  $n \times n$  matrix over some field  $K$ . The following statements are equivalent (I'll group statements that are nearly identical, too):

- The ones I consider useful:
    1.  $\mathbf{A}$  is invertible.
    2.  $\mathbf{A}$  is row-equivalent (and thus column-equivalent) to  $\mathbf{I}_n$ .
    3.  $\det \mathbf{A} \neq 0$ .
    4.  $\text{rank}(\mathbf{A}) = n$ .
    5. The columns of  $\mathbf{A}$  are linearly independent. They span  $K^n$ .  $\text{Col}(\mathbf{A}) = K^n$ .
    6. The transformation  $T(\mathbf{x}) = \mathbf{A}\mathbf{x}$  is a bijection from  $K^n$  to  $K^n$ .
  - The rest:
    1.  $\mathbf{A}$  has  $n$  pivot positions.
    2.  $\mathbf{A}$  can be expressed as a finite product of elementary matrices.
- 

## Understanding correlation vs independence.

- Two events  $A$  and  $B$  are **independent** iff  $P(A \cap B) = P(A)P(B)$ .
- Although  $\text{Indep}(X, Y) \implies \text{cov}(X, Y)=0$ , the converse is *not* true. It's useful to see that statement more explicitly:

$$[\mathbb{E}[XY] = \mathbb{E}[X]\mathbb{E}[Y]] \iff [P(X, Y) = P(X)P(Y)] \quad (1282)$$

### Example: Uncorrelated $\not\Rightarrow$ Independent

To get an intuition for this, I immediately try formalizing the possible cases where this is true. It seems that symmetry is always present in such cases, and they do seem like edge cases. The simplest and by far most common example is the case where we have x,y coordinates  $\{(-1, 0), (0, -1), (0, 1), (1, 0)\}$ .

It's obvious that  $X * Y$  equals zero for all of these points, and also that both X and Y are symmetric about the origin, meaning that  $\mathbb{E}[XY] = 0 = \mathbb{E}[X]\mathbb{E}[Y]$ . In other words, they are **uncorrelated**. The **key insight** comes from understanding *why* this is so: Regardless of whether one variable is either *positive* or *negative* the other is zero. I really want to nail this insight down, because it made me realize I was thinking about correlation wrong – I was thinking about it more as independence in the first place, and so looking back it's no wonder I was confused about the difference. You simply cannot think about correlation from the perspective of a single instance. For example, when I first read this, I thought “well if I know X is 1, then I know automatically that Y is zero”, and although that is technically true, *that is not what correlation is about*. Rather, correlation is about *trends* of multiple instances. A more correct thought would be “Regardless of whether X is positive or negative, Y is zero, therefore positive values of X are neither positively nor negatively correlated with values of Y.”

Now that we've got the thornier part (for me at least) out of the way, recognize that although X and Y are uncorrelated, they are *not* independent. This should be fairly obvious, since given either X or Y, we can say what the other's value is with higher certainty than otherwise.

**Least-Squares Regression.** Note how I emphasized *least-squares*, since in this section we measure how good an estimator is based on least-squares loss. Recall that least-squares arises naturally as a result of modeling  $Y \sim \mathcal{N}(f(X), \varepsilon^2)$  and conducting MLE on the log probability of the data.

- **Linear Least Squares Estimate (LLSE).** The LLSE of response  $Y$  to  $X$ , which we'll denote as  $L[Y | X]$ , is defined as

$$L[Y | X] \triangleq \arg \min_{\hat{Y}} \mathbb{E}_{(X,Y) \sim \mathcal{D}} [(Y - \hat{Y}(X))^2] \quad (1283)$$

$$\text{where } \hat{Y}(X) := a + bX \quad (1284)$$

where finding the optimal linear function  $\hat{Y}(X)$  amounts to finding the optimal coefficients  $(a, b)$  over the dataset  $\mathcal{D}$ . Unfortunately, it seems that the main way of “deriving” the result is actually to just proven, *given* the result, that it does indeed minimize the MSE. So, with that said, we begin with the result:

$$L[Y | X] = \mathbb{E}[Y] + \frac{\text{cov}(X, Y)}{\text{Var}(X)} (X - \mathbb{E}[X]) \quad (1285)$$

### Proof: Eq 1285 Minimizes the MSE

Formally, let  $\mathcal{L}(X) \triangleq \{aX + b \mid a, b \in \mathbb{R}\}$ . Prove that  $\forall \hat{Y} \in \mathcal{L}(X)$ ,

$$\mathbb{E} [(Y - L[Y \mid X])^2] \leq \mathbb{E} [(Y - \hat{Y})^2] \quad (1286)$$

1. Expand the general form of

$$\mathbb{E} [(Y - aX - b)^2] = \mathbb{E} [((Y - L[Y \mid X]) + (L[Y \mid X] - aX - b))^2] \quad (1287)$$

$$\begin{aligned} &= \mathbb{E} [(Y - L[Y \mid X])^2] \\ &\quad + 2\mathbb{E} [(Y - L[Y \mid X])(L[Y \mid X] - aX - b)] \\ &\quad + \mathbb{E} [(L[Y \mid X] - aX - b)^2] \end{aligned} \quad (1288)$$

Our next goal is to evaluate the term in red (spoiler alert: it is zero).

2. First, it is easy to show that  $\mathbb{E} [Y - L[Y \mid X]] = 0$  by simple substitution/arithmetic. We can also show that<sup>a</sup>  $\forall aX + b \in \mathcal{L}(X)$ ,

$$\mathbb{E} [(Y - L[Y \mid X])(aX + b)] = 0$$

as well.

3. Since  $L[Y \mid X] \in \mathcal{L}(X)$ , it is *also* true that  $\forall \hat{Y} \in \mathcal{L}(X)$ , we know  $(L[Y \mid X] - \hat{Y}) \in \mathcal{L}(X)$ , too. Therefore, the red term from step 1 equates to zero.
4. We now know that our formula from step 1 can be written

$$\mathbb{E} [(Y - aX - b)^2] = \mathbb{E} [(Y - L[Y \mid X])^2] + \mathbb{E} [(L[Y \mid X] - aX - b)^2] \quad (1289)$$

Clearly, this minimized when  $aX + b = L[Y \mid X]$ .

---

<sup>a</sup>Also via simple substitution and using  $\text{cov}(x, y) = \mathbb{E}[xy] - \mathbb{E}[x]\mathbb{E}[y]$

**TODO:** Figure out how this formulation is equivalent to the typical multivariate expression:

$$\hat{\mathbf{y}} = \hat{\mathbf{w}}_{OLS}^T \mathbf{x} \quad (1290)$$

$$\hat{\mathbf{w}}_{OLS} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y} \quad (1291)$$

### Questions.

- **Q:** In general, how can one tell if a matrix  $\mathbf{A}$  has an eigenvalue decomposition? [insert more conceptual matrix-related questions here . . . ]
- **Q:** Let  $\mathbf{A}$  be real-symmetric. What can we say about  $\mathbf{A}$ ?
  - Proof that eigendecomposition  $\mathbf{A} = Q\Lambda Q^T$  exists: Wow this is apparently quite hard to prove according to many online sources. Guess I don't feel so bad now that it wasn't (and still isn't) obvious.
  - Eigendecomposition not unique. This is apparently because two or more eigenvectors may have same eigenvalue.

This is the principal axis theorem: if  $A$  symmetric, then orthonorm basis of e-vects exists.

## Stuff I Forget.

- Existence of eigenvalues/eigenvectors. Let  $\mathbf{A} \in \mathbb{R}^{n \times n}$ .
  - $\lambda$  is an eigenvalue of  $A$  iff it satisfies  $\det(\lambda\mathbf{I} - \mathbf{A}) = 0$ . Why? Because it is an equivalent statement as requiring that  $(\lambda\mathbf{I} - \mathbf{A})\mathbf{x} = 0$  has a nonzero solution for  $\mathbf{x}$ .
  - The following statements are equivalent:
    - \*  $\mathbf{A}$  is diagonalizable.
    - \*  $\mathbf{A}$  has  $n$  linearly independent eigenvectors.
  - The **eigenspace** of  $\mathbf{A}$  corresponding to  $\lambda$  is the solution space of the homogeneous system  $(\lambda\mathbf{I} - \mathbf{A})\mathbf{x} = 0$ .
  - $\mathbf{A}$  has at most  $n$  distinct eigenvalues.
- Diagonalizability notes from 5.2 of advanced linear alg. book (261). Recall that  $\mathbf{A}$  is defined to be diagonalizable if and only if there exists an ordered basis  $\beta$  for the space consisting of eigenvectors of  $\mathbf{A}$ .
  - If the standard way of finding eigenvalues leads to  $k$  distinct  $\lambda_i$ , then the corresponding set of  $k$  eigenvectors  $v_i$  are guaranteed to be linearly independent (but might not span the full space).
  - If  $\mathbf{A}$  has  $n$  linearly independent eigenvectors, then  $\mathbf{A}$  is diagonalizable.
  - The characteristic polynomial of any diagonalizable linear operator splits (can be factored into product of linear factors). The **algebraic multiplicity** of an eigenvalue  $\lambda$  is the largest positive integer  $k$  for which  $(t - \lambda)^k$  is a factor of  $f(t)$ .
- Expectation of a random vector. Defined as

$$\mathbb{E}[\mathbf{x}] = \begin{bmatrix} \mathbb{E}[x_1] \\ \vdots \\ \mathbb{E}[x_d] \end{bmatrix} \quad (1292)$$

You can work out that it separates like that (which is not intuitive/immediately obvious imo) by considering e.g. the case where  $d = 2$ . You'll end up finding that

$$\mathbb{E}[\mathbf{x}] = \sum_{x_1} \sum_{x_2} \mathbf{x} p(\mathbf{x} = \mathbf{x}) \quad (1293)$$

$$= \begin{bmatrix} \mathbb{E}[x_1] \\ \mathbb{E}_{x_1} [\mathbb{E}_{x_2 \sim p(x_2|x_1)} [x_2 | x_1]] \end{bmatrix} \quad (1294)$$

and since we know from CS70 that  $\mathbb{E}[\mathbb{E}[Y | X]] = \mathbb{E}[Y]$ , we get the desired result.

Most info here comes from chapter 5 of your "Elementary Linear Algebra" textbook (around pg305)

Recall that a linear operator is a special case of a linear map where the input space is the same as the output space.

# Matrix Cookbook

Table of Contents   Local

*Written by Brandon McKinzie*

Compiling a list of most useful equations from the matrix cookbook.

$$\frac{\partial}{\partial \mathbf{X}} \|\mathbf{X}\|_F^2 = \frac{\partial}{\partial \mathbf{X}} \text{Tr}(\mathbf{X} \mathbf{X}^T) = 2\mathbf{X} \quad (1295)$$

$$[\text{chain rule}] \quad \frac{\partial g(\mathbf{U})}{\partial X_{ij}} = \text{Tr} \left[ \left( \frac{\partial g(\mathbf{U})}{\partial \mathbf{U}} \right)^T \frac{\partial \mathbf{U}}{\partial X_{ij}} \right] \quad (1296)$$

# Main Tasks in NLP

Table of Contents   Local

*Written by Brandon McKinzie*

Going to start compiling a list of the main tasks in NLP (alphabetized). Note that NLP-Progress, a site dedicated to this purpose, is a much more detailed. I'm going for short-and-sweet here.

## Constituency Parsing.

- **Task:** Generate parse tree of a sentence. Nodes are typically labeled by parts of speech and/or chunks.
  - A constituency parse tree breaks a text into sub-phrases, or constituents. Non-terminals in the tree are types of phrases, the terminals are the words in the sentence.

## Coreference Resolution.

- **Task:** clustering mentions in text that refer to the same underlying real world entities.
- **SOTA:** End-to-end Neural Coreference Resolution.

## Dependency Parsing.

- **Task:** Given a sentence, generate its dependency tree (DT). A DT is a labeled directed tree whose nodes are individual words, and whose edges are directed arcs labeled with dependency types.
    - A dependency parser analyzes the grammatical structure of a sentence, establishing relationships between “head” words and words which modify those heads.
  - **SOTA:** Deep Biaffine Attention for Neural Dependency Parsing.
- Related Tasks:**
- Constituency parsing. See this great wiki explanation of dependency vs constituency.

## Information Extraction.

- **Task:** Given a (typically long) portion of raw text, recover information about pre-specified relations, entities, events, etc.

## Language Modeling.

- **Task:** Learning the probability distribution over text sequences. Often used for predicting the next word in a sequence, given the K previous words.
- **SOTA:** ELMo.

## Machine Translation.

## Semantic Parsing.

- **Task:** Translating natural language into a formal meaning representation *on which a machine can act*.

## Semantic Role Labeling.

- **Task:** Given a sentence, extract the predicates<sup>280</sup> and their respective arguments.
- **Historical Approaches.**
  - **CCG Semantic Parsing.** Zettlemoyer & Collins 2005, 2007.
  - Seq2seq. Dong & Lapata, 2016.

## Sentiment Analysis.

- **Task:** Determining whether a piece of text is positive, negative, or neutral.
- **SOTA:** Biattentive classification network (BCN) from Learned in Translation: Contextualized Word Vectors (the CoVe paper) with ELMo embeddings.

## Summarization.

## Textual Entailment.

- **Task:** Given a premise, determine whether a proposed hypothesis is true.
- **SOTA:** Enhanced Sequential Inference (ESIM) model from Enhanced LSTM for Natural Language Inference with ELMo embeddings.

---

<sup>280</sup>The predicate of a sentence mostly corresponds to the main verb and any auxiliaries that accompany the main verb; whereas the arguments of that predicate (e.g. the subject and object noun phrases) are outside the predicate.

## **Topic Modeling.**

**Question Answering.** Also called **machine reading comprehension**.

- **Task:** Given a paragraph of text, “read” it and then answer questions pertaining to the text.
- **Dataset:** The main benchmark dataset is the *Stanford Question Answering Dataset* (SQuAD). Each sample has the form (*{question, paragraph}*, *answer*), where the answer is a subsequence found somewhere in the paragraph (i.e. this is an *extractive* task, not abstractive).
- **SOTA:** Improved versions of the Bidirectional Attention Flow (BiDAF) model, with ELMo embeddings.
- **Related tasks:**

## **Word Sense Disambiguation.**

- **Task:** Associating words in context with their most suitable entry in a pre-defined sense inventory.

## Misc. Topics

Table of Contents   Local

*Written by Brandon McKinzie*

### 14.5.1 BLEU SCORE

BiLingual Evaluation Understudy. For scoring machine-generated translations when we have access to one or more reference translations.

- Unigram precision: Really naive and basically useless version:

$$P = \frac{\text{num pred words that also appear somewhere in ref words}}{\text{total num pred words}} \quad (1297)$$

It's important to emphasize how ridiculous this really is. It literally means that we walk along each word in the prediction and ask "is this word somewhere in any of the reference translations?" and if that answer is "yes", we +1 to the numerator. Period.

- Modified unigram precision: actually considering how many times we've mentioned a given word  $w$  when incorporating it into the precision calculation. Now, when walking along a sentence, we add to the aforementioned question, "...and have I seen it less than N times already?" where  $N = [\max(\text{count}(\text{sent}, w)) \text{ for } \text{sent} \in \text{refs}]$ . This means our numerator can now be at most  $N$  for any given word.
- Generalize to n-grams. Below is the formula for Blue score on n-grams only:

$$p_n(\hat{y}) = \frac{\sum_{\text{ngrams} \in \hat{y}} \text{Count}_{clip}(\text{ngram}, \text{refs})}{\sum_{\text{ngrams} \in \hat{y}} \text{Count}(\text{ngram})} \quad (1298)$$

- Combined Blue score.

$$= \text{BP} \cdot \exp \left\{ \frac{1}{4} \sum_{n=1}^4 \log p_n \right\} \quad (1299)$$

$$\text{BP} = \begin{cases} 1 & \text{len(pred)} > \text{len(ref)} \\ \exp\{1 - \text{len(pred)}/\text{len(ref)}\} & \text{otherwise} \end{cases} \quad (1300)$$

where BP is the **brevity penalty**.

---

### 14.5.2 CONNECTIONIST TEMPORAL CLASSIFICATION (CTC)

---

Approach for mapping input sequences  $X = \{x_1, \dots, x_T\}$  to label sequences  $Y = \{y_1, \dots, y_U\}$ , where the lengths may vary ( $T \neq U$ ). First, we need to define the meaning of an **alignment** between input sequence  $X$  and label sequence  $Y$ . Most generally, an alignment is a composition of one or more functions, and accept input  $X$  and ultimately map to output  $Y$ . Take, for example, the label sequence  $Y = \{h, e, l, l, o\}$  and some input sequence (e.g. raw audio)  $X = \{x_1, \dots, x_{12}\}$ . CTC places the following constraints on the [first function of the] alignment sequence:

1. It must be the same length as the input sequence  $X$ .
2. It has the same vocabulary as  $Y$ , plus an additional token  $\epsilon$  to denote blanks.
3. At position  $i$ , it either (a) repeats the aligned token at  $i - 1$ , (b) assigns the empty token  $\epsilon$ , or (c) assigns the next letter of the label sequence.

For our example, we could have an aligned sequence  $A = \{h, h, e, \epsilon, \epsilon, l, l, l, \epsilon, l, l, o\}$ . Then we apply the following two steps (can interpret as functions) to map from  $A$  to  $Y$ :

1. Merge any repeated [contiguous] characters.
2. Remove any  $\epsilon$  tokens.

When you hear someone say “the CTC loss,” they usually mean “MLE using a CTC posterior.” In other words, there is no “CTC loss” function, but rather there is the standard maximum likelihood objective, but we use a particular form for the posterior  $p(Y | X)$  over possible output labels  $Y$  given raw input sequence  $X$ :

$$p(Y | X) = \sum_{\mathcal{A} \in \mathcal{A}_{X,Y}} \prod_{t=1}^T p_t(a_t | X) \quad (1301)$$

where  $\mathcal{A}$  is one of the valid alignments from  $\mathcal{A}_{X,Y}$ , the full set of valid alignments from  $X$  to a given output sequence  $Y$ . The per-timestep probabilities  $p(a_t | X)$  can be given by, for example, an RNN.

## Number of Valid Alignments

Given  $X$  of length  $T$  and  $Y$  of length  $U \leq T$  (and no repeating letters), how many valid alignments exist?

The differences between alignments fall under two categories:

1. Indices where we transition from one label to the next.
2. Indices where we insert the blank token,  $\epsilon$ .

Stated even simpler, the alignments differ first and foremost by *which elements of  $X$  are “extra” tokens*, where I’m using “extra” to mean either blank or repeat token. Given a set of  $T$  tokens, there are  $\binom{T}{T-U}$  different ways to assign  $T - U$  of them as “extra.” The tricky part is that we can’t just randomly decide to repeat or insert a blank, since a sequence of one or more blanks is *always* followed by a transition to next letter, by definition. And remember, we have defined  $Y$  to have no repeated [contiguous] labels.

Apparently, the answer is  $\binom{T+U}{T-U}$  total valid alignments.

Computing forward probabilities  $\alpha_t(s)$ , defined as the probability of arriving at [prefix of] augmented label sequence  $\ell'_{(1\dots s)}$  given unmerged alignments up to step  $t$ . There are two cases to consider.

1. (1.1) The augmented label at step  $s$ ,  $\ell'_s$  is the blank token  $\epsilon$ . Remember,  $\epsilon$  occurs at every other position in the augmented sequence  $\ell'$ . At the previous RNN output (time  $t - 1$ ), we could’ve emitted either a blank token  $\epsilon$  or the previous token in the augmented label sequence,  $\ell'_{s-1}$ . In other words,

$$y_{\ell'_s=\epsilon}^t (\alpha_{t-1}(s) + \alpha_{t-1}(s-1)) \quad (1302)$$

2. (1.2) The augmented label at step  $s$ ,  $\ell'_s$  is the same augmented label as at step  $s - 2$ . This occurs when the [not augmented] label sequence has repeated labels next to each other.

$$y_{\ell'_s=\ell'_{s-2}}^t (\alpha_{t-1}(s) + \alpha_{t-1}(s-1)) \quad (1303)$$

In this situation,  $\alpha_{t-1}(s)$  corresponds to us just emitting the same token as we did at  $t - 1$  or emitting a blank token  $\epsilon$ , and  $\alpha_{t-1}(s - 1)$  corresponds to a transition to/from  $\epsilon$  and a label.

3. (2) The augmented label at step  $s - 1$ ,  $\ell'_{s-1}$  is the blank token  $\epsilon$  between unique characters. In addition to the two  $\alpha_{t-1}$  terms from before, we now also must consider the possibility that our RNN emitted  $\ell'_{s-2}$  at the previous time ( $t - 1$ ) and then emitted  $\ell'_s$  immediately after at time  $t$ .

---

### 14.5.3 PERPLEXITY

---

Per Wikipedia:

*In information theory, perplexity is a measurement of how well a probability distribution or probability model predicts a sample.*

The perplexity,  $\mathcal{P}$ , of discrete probability distribution  $p$  over word sequences  $W = \{w_1, \dots, w_T\} = \mathbf{w}_{(1\dots T)}$  of length  $T$  is defined as:

$$\mathcal{P}(p) = 2^{-\mathbb{E}_{\mathbf{x} \sim p}[\lg p(\mathbf{x})]} = 2^{H(p)} \quad (1304)$$

$$= 2^{-\sum_{\mathbf{w}_{(1\dots T)}} p(\mathbf{w}_{(1\dots T)}) \lg p(\mathbf{w}_{(1\dots T)})} \quad [\text{theory}] \quad (1305)$$

$$\approx 2^{-\frac{1}{N} \sum_{\mathbf{w}_{(1\dots T)} \in \mathcal{T}} \lg q(\mathbf{w}_{(1\dots T)})} \quad [\text{empirical}] \quad (1306)$$

where  $H$  is entropy (in bits). It's important to note that, in practice, we are never able to use the theoretical version since we don't know  $p$  exactly (we are usually trying to estimate it) – instead of  $H(p)$  we thus usually think in terms of  $H(p, q)$ , the *cross entropy*<sup>281</sup>. The empirical definition is when we have  $N$  samples in some test set  $\mathcal{T}$ , and a model  $q$  that we want to approximate the true distribution  $p$ .

In NLP, it is more common to want the *per-word* perplexity of a language model. We typically do this by flattening out a sequence of words in some test set containing  $M$  words total and simply compute

$$\mathcal{P} = 2^{-\frac{1}{M} \lg p(\{w_1, \dots, w_M\})} \quad (1307)$$

$$= \frac{1}{p(\{w_1, \dots, w_M\})^{\frac{1}{M}}} \quad (1308)$$

In other words, NLP nearly always defines perplexity as the **inverse probability of the test set**, normalized by number of words. So, why is this valid? We are implicitly assuming that language sources are **ergodic**:

#### Ergodic

*A random process is **ergodic** if its (asymptotic) time average is the same as its expectation value over all possible states (w.r.t the specified probability distribution).*

*Informally, this means that the system eventually reaches all states, and such that the probability of observing it in state  $s$  is  $p(s)$ , where  $p$  is the true generating distribution.*

---

<sup>281</sup>Recall the relationship between entropy  $H(p)$  and cross entropy  $H(p, q)$ :

$$H(p, q) = H(p) + D_{KL}(p||q)$$

In the per-word NLP case, this means we can assume that

$$\lim_{m \rightarrow \infty} \mathbb{E} [\lg p(\{w_1, \dots, w_m\})] = \lim_{m \rightarrow \infty} \frac{1}{m} \lg p(\{w_1, \dots, w_m\}) \quad (1309)$$

where the sequence on the RHS is any sample from  $p^{282}$ .

**Intuition.** Ok, now that we've got definitions out of the way, what does it actually mean? First consider some limiting cases. If the distribution  $p$  is uniform over  $N$  possible outcomes, then  $\mathcal{P}(p) = 2^{\lg N} = N$ . Since the uniform distribution has the highest possible entropy,  $N$  is also the largest possible value for perplexity of a discrete distribution  $p$  over  $N$  possible outcomes.

Consider the interpretation of the cross entropy loss as the negative log-likelihood:

$$NLL(p_{data}) = -\frac{1}{M} \sum_{i=1}^M \log p(w=w_i) = \mathbb{E}_{w_i \sim p_{data}} \left[ \log \frac{1}{p(w)} \right] \quad (1310)$$

we see that  $NLL$  (and thus  $\mathcal{P}=\exp(NLL)$ ) decreases as our model assigns higher probabilities to samples drawn from  $p_{data}$ . *Better models of  $p_{data}$  are less surprised by samples from  $p_{data}$ .* If we use the typical interpretation of entropy as the number of bits needed (on average) to represent a sample from  $p$ , then the perplexity can be interpreted as the total number of possible results (on average) when drawing a sample from  $p$ .

In the case of language modeling, this represents the total number of reasonable next-word predictions for  $w_{t+1}$  given some context  $w_1, \dots, w_t$ . As our model assigns higher probabilities to the true samples in  $p_{data}$ , the number of bits required to specify each word, on average, becomes smaller. Therefore, you can roughly think of per-word perplexity as telling you the number of possible choices, on average, your model considers uniformly at random at a given step. For example,  $\mathcal{P} = 42$  could be interpreted loosely as “to predict the next word out of some vocabulary  $V$ , my model can narrow it down on average to about 42 choices, and chooses uniformly at random from that subset”, where typically  $|V| \gg 42$ .

---

<sup>282</sup>Something feels off here. I'm synthesizing what I'm reading from wikipedia and this source from berkeley but I can't fix the sloppy tone of the wording.

---

#### 14.5.4 BYTE PAIR ENCODING

---

---

#### 14.5.5 GRAMMARS

---

In formal language theory, a **formal grammar** is a set of production rules for strings in a formal language. The rules describe how to form strings from the language's alphabet that are valid according to the language's syntax. A grammar does not describe the meaning of the strings or what can be done with them in whatever context—only their form.

- **Regular Grammar:** no rule has more than one nonterminal in its right-hand side, and each of these nonterminals is at the same end of the right-hand side. Every regular grammar corresponds directly to a nondeterministic finite automaton.
- A context-free grammar (**CFG**) is a formal grammar that consists of:
  - **Terminal symbols:** characters that appear in the strings generated by the grammar.
  - **Nonterminal symbols:** placeholders for patterns of terminal symbols that can be generated by the nonterminal symbols.
  - **Productions:** rules for replacing (or rewriting) nonterminal symbols (on the LHS) in a string with other nonterminal or terminal symbols (on the RHS), *which can be applied regardless of context*.
  - **Start symbol:** a special nonterminal symbol that appears in the initial string generated by the grammar.

To generate a string of terminal symbols from a CFG, we:

1. Begin with a string consisting of the start symbol;
2. Apply one of the productions with the start symbol on the left hand size, replacing the start symbol with the right hand side of the production;
3. Repeat the process of selecting nonterminal symbols in the string, and replacing them with the right hand side of some corresponding production, until all nonterminals have been replaced by terminal symbols.

- A **Probabilistic CFG** extends CFGs the same way HMMs extend regular grammars, by defining the set P of probabilities on production rules.

---

#### 14.5.6 BLOOM FILTER

---

Data structure for querying whether a data point is a member of some set. It returns either “no” or “maybe”. It is implemented as a bit vector. Each member of the set is passed through  $k$  hash functions. Each hash function maps an element to an integer index. For each member, we set the  $k$  output indices of the hash functions to 1 in our bit vector. To answer if some data point  $x$  is in the set, we pass  $x$  through the  $k$  hash functions, which gives us  $k$  indices. If all  $k$  indices have their bit value set to 1, the answer is “maybe”, otherwise (if any bit value is 0) the answer is “no”.

---

#### 14.5.7 DISTRIBUTED TRAINING

---

##### Asynchronous SGD.

$$W_{i+1} = W_i - \frac{\alpha}{N_x} \sum_{j=1}^{N_x} \frac{\partial L(\mathbf{x}^{(j)})}{\partial W_i} \quad (1311)$$

$$\text{[SyncSGD]} \quad W_{i+1} = W_i - \lambda \sum_{j=1}^{N_w} \sum_{k=1}^{N_x(j)} \alpha \frac{\partial L(\mathbf{x}^{(k)})}{\partial W_i} \quad (1312)$$

where  $N_x$  is the number of data samples.

In asynchronous SGD, we just apply the gradient updates to a global version of the parameters whenever they are available. In practice, this can result in **stale gradients**, which happens when a worker takes a long time to compute some gradient step, while the master version of the parameters has been updated many times. This results in the master computing an update like  $W_{t+1} = W_t - \lambda \Delta W_{t-D}$  for larger-than-desired values of  $D$  (delay in num updates).

---

#### 14.5.8 TRADITIONAL LANGUAGE MODELING

---

Some quick terminology/definitions since my brain forgets things.

- **Backoff:** when we want to estimate e.g.  $p(w_1, w_2, w_3)$  but we've never seen the sequence  $w_1, w_2, w_3$ . We can instead *backoff* to bigrams if we *have* seen e.g. the sequences  $w_1, w_2$  and  $w_2, w_3$ . More generally, if we have many N-gram LMs with different values of N, we backoff to the highest order LM that contains the  $N$ -gram we are querying for.
  - Example failure mode of backoff: we typically have to backoff for unusual sequences of words (by definition). The lower order N gram model could drastically overestimate the backoff probabilities<sup>283</sup>.
- **Interpolation:** Using a combination of N-gram LMs with different values of  $N$  as an attempt to utilize the strengths of each and mitigate their weaknesses.
- **Kneser-Ney:** (link)

$$P_{KN}(w_t | w_{t-1}) = \frac{\max(c(w_{t-1}, w_t), 0)}{c(w_{t-1})} + \lambda \frac{|\{w_{t-1} : c(w_{t-1}, w_t) > 0\}|}{|\{w_{t'-1} : c(w_{t'-1}, w'_t) > 0\}|} \quad (1313)$$

Click this link to see a really good overview of the terms above and more.

---

<sup>283</sup>Good example is how a unigram model assigns a decent probability for "York" but a human bigram could tell you that it is nearly certain that "New" preceded it. The backoff model would tend to overestimate  $p(\text{York})$  since it has no contextual information