

## CONTENTS

<b>1 Neural Machine Translation</b>	<b>5</b>
1.1 Neural Conversation Model . . . . .	6
1.2 NMT By Jointly Learning to Align & Translate . . . . .	8
1.2.1 Detailed Model Architecture . . . . .	9
1.3 Using Large Vocabularies for NMT . . . . .	11
1.4 Candidate Sampling – TensorFlow . . . . .	14
1.5 Attention Terminology . . . . .	16
1.6 Effective Approaches to Attention-Based NMT . . . . .	18
<b>2 Conditional Random Fields</b>	<b>20</b>
2.1 Conditional Random Fields . . . . .	21
2.2 An Introduction to Conditional Random Fields . . . . .	24
2.2.1 Inference (Sec. 4) . . . . .	28
2.2.2 Parameter Estimation (Sec. 5) . . . . .	31
2.2.3 Related Work and Future Directions (Sec. 6) . . . . .	34
2.3 Hidden-Unit Conditional Random Fields . . . . .	35
2.3.1 Detailed Derivations . . . . .	37
2.4 Pre-training of Hidden-Unit CRFs . . . . .	42
2.5 Neural Conditional Random Fields . . . . .	44
2.6 Bidirectional LSTM-CRF Models for Sequence Tagging . . . . .	46
2.7 On Herding and the Perceptron Cycling Theorem . . . . .	47
<b>3 Structured Prediction</b>	<b>49</b>
3.1 Joint Event Extraction via RNNs . . . . .	50
3.2 Event Extraction via Bidi-LSTM Tensor NNs . . . . .	52
3.3 Reasoning with Neural Tensor Networks . . . . .	54
3.4 Language to Logical Form with Neural Attention . . . . .	55
3.5 Seq2SQL: Generating Structured Queries from NL using RL . . . . .	57
3.6 SLING: A Framework for Frame Semantic Parsing . . . . .	60
3.7 DeepWalk: Online Learning of Social Representations . . . . .	62
3.8 Review of Relational Machine Learning for Knowledge Graphs . . . . .	64
3.9 Fast Top-K Search in Knowledge Graphs . . . . .	67
3.10 Dynamic Recurrent Acyclic Graphical Neural Networks (DRAGNN) . . . . .	69
3.10.1 More Detail: Arc-Standard Transition System . . . . .	72
3.11 Joint Extraction of Events and Entities within a Document Context . . . . .	73
3.12 Globally Normalized Transition-Based Neural Networks . . . . .	76
3.13 Structured Attention Networks . . . . .	79
3.14 Relation Extraction: A Survey . . . . .	81
3.15 Neural Relation Extraction with Selective Attention over Instances . . . . .	84

<b>4 Language Modeling</b>	<b>86</b>
4.1 Improving Language Understanding by Generative Pre-Training . . . . .	87
4.2 Exploring the Limits of Language Modeling . . . . .	88
4.3 Noise Contrastive Estimation . . . . .	90
4.3.1 Self-Normalized NCE . . . . .	92
4.4 Improving Neural Language Models with a Continuous Cache . . . . .	94
4.5 Context Dependent RNN Language Model . . . . .	95
4.6 Strategies for Training Large Vocabulary Neural Language Models . . . . .	96
4.7 Large Memory Layers with Product Keys . . . . .	98
4.8 XLNet . . . . .	100
4.9 Transformer-XL . . . . .	102
4.10 Efficient Softmax Approximation for GPUs . . . . .	103
4.11 Adaptive Input Representations for Neural Language Modeling . . . . .	104
4.12 Fine-Tuning Language Models from Human Preferences . . . . .	105
4.13 Cross-lingual Language Model Pretraining . . . . .	106
4.14 Learning and Evaluating General Linguistic Intelligence . . . . .	108
4.15 Scaling Laws for Neural Language Models . . . . .	111
4.16 Language Models are Few-Shot Learners . . . . .	113
4.17 Generating Long Sequences with Sparse Transformers . . . . .	118
4.18 The Curious Case of Neural Text Degeneration . . . . .	120
<b>5 Embeddings and Transfer Learning</b>	<b>121</b>
5.1 Simple Baseline for Sentence Embeddings . . . . .	122
5.2 Deep Sentence Embedding Using LSTMs . . . . .	124
5.3 Supervised Universal Sentence Representations (InferSent) . . . . .	127
5.4 Dist. Rep. of Sentences from Unlabeled Data (FastSent) . . . . .	128
5.5 Attention Is All You Need . . . . .	130
5.6 Poincaré Embeddings for Learning Hierarchical Representations . . . . .	134
5.7 Enriching Word Vectors with Subword Information (FastText) . . . . .	136
5.8 Deep Contextualized Word Representations . . . . .	138
5.9 BERT . . . . .	140
5.10 Wasserstein is all you need . . . . .	142
5.11 On the Dimensionality of Word Embedding . . . . .	144
5.12 Universal Language Model Fine-Tuning for Text Classification . . . . .	145
5.13 Hopfield Networks is All You Need . . . . .	147
5.13.1 Review: Hopfield Networks . . . . .	150
<b>6 Vision and Language</b>	<b>152</b>
6.1 Show, Ask, Attend, and Answer . . . . .	153
6.2 Neural Module Networks . . . . .	154
6.3 Learning to Compose Neural Networks for QA . . . . .	156
6.4 End-to-End Module Networks for VQA . . . . .	158
6.5 Fast Multi-language LSTM-based Online Handwriting Recognition . . . . .	160

6.6	Multi-Language Online Handwriting Recognition . . . . .	161
6.7	Learning Two-Branch Neural Networks for Image-Text Matching Tasks . . . . .	163
6.8	Learning by Abstraction: The Neural State Machine . . . . .	165
6.9	Tags and Attributes based Visual Commonsense Reasoning Baselines . . . . .	168
6.10	From Recognition to Cognition: Visual Commonsense Reasoning . . . . .	170
6.11	GQA . . . . .	172
6.12	A negative case analysis of visual grounding methods for VQA . . . . .	175
6.13	Bottom-Up and Top-Down Attention for Image Captioning and Visual Question Answering . . . . .	177
6.14	ViLBERT . . . . .	179
6.14.1	12-in-1: Multi-Task Vision and Language Representation Learning . . . . .	180
6.15	Language-Conditioned Graph Networks for Relational Reasoning . . . . .	181
6.16	Towards VQA Models That Can Read . . . . .	183
6.17	Accuracy vs. Complexity: A Trade-off in VQA Models . . . . .	184
6.18	Iterative Answer Prediction with Multimodal Transformers for TextVQA . . . . .	186
6.19	CVPR 2020 . . . . .	189
6.19.1	Unified Visual-Semantic Embeddings: Bridging Vision and Language with Structured Meaning Representations . . . . .	189
6.19.2	Counterfactual Samples Synthesizing for Robust Visual Question Answering . . . . .	190
6.19.3	What Makes Training Multi-modal Classification Networks Hard? . . . . .	191
<b>7</b>	<b>Generative Models</b>	<b>192</b>
7.1	WaveNet . . . . .	193
7.2	Latent Dirichlet Allocation . . . . .	197
7.3	Generative Adversarial Nets . . . . .	200
7.4	Modular Generative Adversarial Networks . . . . .	203
7.5	Transfer Learning from Speaker Verification to TTS . . . . .	205
7.6	Tacotron 2 . . . . .	206
7.7	Glow . . . . .	208
7.8	WaveGlow . . . . .	209
<b>8</b>	<b>Meta Learning</b>	<b>210</b>
8.1	Online Meta-Learning . . . . .	211
<b>9</b>	<b>Optimization</b>	<b>213</b>
9.1	Non-Convex Optimization for Machine Learning . . . . .	214
9.1.1	Non-Convex Projected Gradient Descent (3) . . . . .	217
9.2	Co-sampling: Training Robust Networks for Extremely Noisy Supervision . . . . .	218
9.3	Connectionist Temporal Classification . . . . .	219
9.3.1	Sequence Modeling With CTC . . . . .	221
9.4	Large-Scale Study of Curiosity Driven Learning . . . . .	224
9.5	The Marginal Value of Adaptive Gradient Methods in Machine Learning . . . . .	225
9.6	Deep Double Descent . . . . .	226
9.7	A Tutorial on Bayesian Optimization of Expensive Cost Functions . . . . .	227

9.8	Practical Bayesian Optimization of Machine Learning Algorithms . . . . .	229
9.8.1	Review of Gaussian Processes . . . . .	231
9.9	Dragonfly . . . . .	233
<b>10</b>	<b>Miscellaneous</b>	<b>235</b>
10.1	Neural Style . . . . .	236
10.2	TextRank . . . . .	238
10.2.1	Keyword Extraction . . . . .	240
10.2.2	Sentence Extraction . . . . .	241
10.3	Survey of Text Clustering Algorithms . . . . .	242
10.3.1	Distance-based Clustering Algorithms . . . . .	245
10.3.2	Probabilistic Document Clustering and Topic Models . . . . .	246
10.3.3	Online Clustering with Text Streams . . . . .	248
10.4	Clustering Massive Text Streams . . . . .	250
10.5	Hierarchical Attention Networks . . . . .	252
10.6	Neural Architecture Search with Reinforcement Learning . . . . .	255
10.7	Neural Ordinary Differential Equations . . . . .	257
10.8	A Framework for Intelligence and Cortical Function . . . . .	259
10.9	A Theoretically Grounded Application of Dropout in Recurrent Neural Networks . . . . .	260
10.10	Protection Against Reconstruction and Its Applications in Private Federated Learning . . . . .	261
10.11	Product quantization for nearest neighbor search . . . . .	263
10.12	Did the Model Understand the Question? . . . . .	264
10.13	Solving Rubik's Cube with a Robot Hand . . . . .	265
10.14	Decoupled Weight Decay Regularization . . . . .	267
10.15	SmartChoices . . . . .	268
10.16	The Contextual Bandits Problem . . . . .	269
10.17	The principles of adaptation in organisms and machines . . . . .	272

# NEURAL MACHINE TRANSLATION

## CONTENTS

1.1	Neural Conversation Model . . . . .	6
1.2	NMT By Jointly Learning to Align & Translate . . . . .	8
1.2.1	Detailed Model Architecture . . . . .	9
1.3	Using Large Vocabularies for NMT . . . . .	11
1.4	Candidate Sampling – TensorFlow . . . . .	14
1.5	Attention Terminology . . . . .	16
1.6	Effective Approaches to Attention-Based NMT . . . . .	18

## Neural Conversation Model

[Table of Contents](#)   [Local](#)

*Written by Brandon McKinzie*

**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.

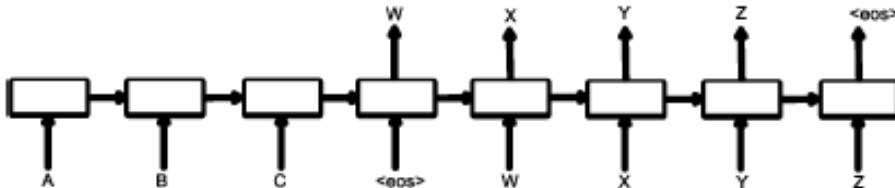
**Related Work.** The authors' approach is based on the following (linked and saved) papers on seq2seq:

- 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>1</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>1</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 &amp; 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>2</sup> and translate.

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

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

$$s_i = f(s_{i-1}, y_{i-1}, c_i) \quad (2)$$

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:

$$c_i = \sum_{j=1}^{T_x} \alpha_{ij} h_j \quad (3)$$

$$\alpha_{ij} = \frac{\exp(e_{ij})}{\sum_{k=1}^{T_x} \exp(e_{ik})} \quad (4)$$

$$e_{ij} = a(s_{i-1}, h_j) \quad (5)$$

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>2</sup>By “align” the authors are referring to aligning the source-search to the relevant parts for prediction.

---

### 1.2.1 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:**

$$z_t = \sigma(W_z y_{t-1} + U_z [s_{t-1}, c_t]) \quad [\text{update gate}] \quad (6)$$

$$r_t = \sigma(W_r y_{t-1} + U_r [s_{t-1}, c_t]) \quad [\text{reset gate}] \quad (7)$$

(8)

3. **Proposal state:**

$$\tilde{s}_t = \tanh(W y_{t-1} + U[r_t \circ s_{t-1}, c_t]) \quad (9)$$

4. **Hidden state:**

$$s_t = (1 - z_t) \circ s_{t-1} + z_t \circ \tilde{s}_t \quad (10)$$

**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:

$$a_j = a(s_{t-1}, h_j) = v_a^T \tanh(W_a s_{t-1} + U_a h_j) \quad (11)$$

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

$$a_j \leftarrow \frac{e^{a_j}}{\sum_{k=0}^{L_{enc}-1} e^{a_k}} \quad (12)$$

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

$$c = \sum_{j=1}^{L_{enc}} a_j h_j \quad (13)$$

**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  $t$ ]:

$$\Pr(y | s, c) \propto e^{y^T W_o u} \quad (14)$$

$$u = [\max\{\tilde{u}_{2j-1}, \tilde{u}_{2j}\}]_{j=1,\dots,\ell}^T \quad (15)$$

$$\tilde{u} = U_o[s_{t-1}, c] + V_o y_{t-1} \quad (16)$$

**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.

## 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:

$$x = (x_1, \dots, x_T) \quad [\text{source sentence}] \quad (17)$$

$$h = (h_1, \dots, h_T) \quad [\text{encoder state seq}] \quad (18)$$

$$h_t = f(x_t, h_{t-1}) \quad (19)$$

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:

$$y = (y_1, \dots, y'_T) \quad [y_i \in \mathbb{Z}] \quad (20)$$

$$\Pr[y_t | y_{<t}, x] \propto e^{q(y_{t-1}, z_t, c_t)} \quad (21)$$

$$z_t = g(y_{t-1}, z_{t-1}, c_t) \quad [\text{decoder hidden?}] \quad (22)$$

$$c_t = r(z_{t-1}, h_1, \dots, h_T) \quad [\text{decoder inp?}] \quad (23)$$

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,

$$\theta^* = \arg \max_{\theta} \sum_{n=1}^N \sum_{t=1}^{T_n} \log (\Pr[y_t^n | y_{<t}^n, x^n]) \quad (24)$$

**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  $c_t$ .**

$$c_t = \sum_{i=1}^T \alpha_i h_i \quad (25)$$

$$\alpha_t = \frac{e^{a(h_t, z_{t-1})}}{\sum_k e^{a(h_k, z_{t-1})}} \quad (26)$$

$a$  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 21 as<sup>3</sup>,

$$\Pr[y_t | y_{<t}, x] = \frac{1}{Z} e^{\mathbf{w}_t^T \phi(y_{t-1}, z_t, c_t) + b_t} \quad (27)$$

$$Z = \sum_{k: y_k \in V} e^{\mathbf{w}_k^T \phi(y_{t-1}, z_t, c_t) + b_k} \quad (28)$$

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>4</sup>, written in terms of the energy  $\mathcal{E}$ .

$$\nabla \log (\Pr[y_t | y_{<t}, x]) = \nabla \mathcal{E}(y_t) - \sum_{k: y_k \in V} \Pr[y_k | y_{<t}, x] \nabla \mathcal{E}(y_k) \quad (29)$$

$$\mathcal{E}(y_j) = \mathbf{w}_j^T \phi(y_{t-1}, z_t, c_t) + b_j \quad (30)$$

---

<sup>3</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>4</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  $V'$ .

$$\mathbb{E}_P [\nabla \mathcal{E}(y)] \approx \sum_{k: y_k \in V'} \frac{\omega_k}{\sum_{k': y_{k'} \in V'} \omega_{k'}} \nabla \mathcal{E}(y_k) \quad (31)$$

$$\omega_k = e^{\mathcal{E}(y_k) - \log Q(y_k)} \quad (32)$$

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

$$\omega_k = \frac{e^{\mathcal{E}(y_k)}}{Q(y_k)} \quad \text{thus} \quad p(y_k | y_{<t}, x) = \omega_k \frac{Q(y_k)}{Z} \quad (33)$$

$$\rightarrow e^{\mathcal{E}(y_k)} = Z \cdot p(y_k | y_{<t}, x) = Q(y_k) \cdot \omega_k \quad (34)$$

### 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.

$$\mathbb{E}_P [\nabla \mathcal{E}(y)] = \sum_{k: y_k \in V} \frac{\omega_k \cdot Q(y_k)}{\sum_{k': y_{k'} \in V} \omega_{k'} \cdot Q(y_{k'})} \nabla \mathcal{E}(y_k) \quad (35)$$

$$\mathbb{E}_P [\nabla \mathcal{E}(y)] = \sum_{k: y_k \in V'} \frac{\omega_k}{\sum_{k': y_{k'} \in V'} \omega_{k'}} \nabla \mathcal{E}(y_k) \quad (36)$$

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:

$$\Pr(y_i^{true} = y | C_i, x_i) = \frac{\Pr(y_i^{true} = y | x_i) \cdot \Pr(C_i | y_i^{true} = y, x_i)}{\Pr(C_i | x_i)} \quad (37)$$

$$= \frac{\Pr(y | x_i)}{Q(y | x_i)} \cdot \frac{1}{K(x_i, C_i)} \quad (38)$$

where they've just defined

$$K(x_i, C_i) \triangleq \frac{\Pr(C_i | x_i)}{\prod_{y' \in C_i} Q(y' | x_i) \prod_{y' \in (L - C_i)} (1 - Q(y' | x_i))} \quad (39)$$

## 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

$$F(x, y) = \log(\Pr(y | x)) + K(x) \quad (40)$$

which is the general form because

$$\text{Softmax}(\log(\Pr(y | x)) + K(x)) = \frac{e^{\log(\Pr(y | x)) + K(x)}}{\sum_{y'} e^{\log(\Pr(y' | x)) + K(x)}} \quad (41)$$

$$= \Pr(y | x) \quad (42)$$

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 40 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 38. 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

$$\text{Softmax Input: } F(x, y) - \log(Q(y | x)) \quad (43)$$

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

$$e_i = v^T \tanh(\text{FC}(s_{i-1}) + \text{FC}(h)) \quad (44)$$

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>5</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>5</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:

$$\text{score}(\mathbf{h}_t, \bar{\mathbf{h}}_s) = \begin{cases} \mathbf{h}_t^T \bar{\mathbf{h}}_s & \text{dot} \\ \mathbf{h}_t^T \mathbf{W}_a \bar{\mathbf{h}}_s & \text{general} \\ \mathbf{v}_a^T \tanh(\mathbf{W}_a[\mathbf{h}_t; \bar{\mathbf{h}}_s]) & \text{concat} \end{cases} \quad (45)$$

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.

$$\mathbf{a}_t(s) = \text{align}(\mathbf{h}_t, \bar{\mathbf{h}}_s) \quad (46)$$

$$= \frac{\exp(\text{score}(\mathbf{h}_t, \bar{\mathbf{h}}_s))}{\sum_{s'} \exp(\text{score}(\mathbf{h}_t, \bar{\mathbf{h}}_{s'}))} \quad (47)$$

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.
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:

$$\tilde{\mathbf{h}}_t = \tanh(\mathbf{W}_c [\mathbf{c}_t, \mathbf{h}_t]) \quad (48)$$

$$p(y_t | y_{<t}, x) = \text{softmax}(\mathbf{W}_s \tilde{\mathbf{h}}_t) \quad (49)$$

## 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:

$$\tilde{\mathbf{h}}_t = \tanh(\mathbf{W}_c[\mathbf{c}_t; \mathbf{h}_t]) \quad (50)$$

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

$$p(y_t | y_{<t}, x) = \text{softmax}(\mathbf{W}_s \tilde{\mathbf{h}}_t) \quad (51)$$

**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:

$$\text{score}(\mathbf{h}_t, \bar{\mathbf{h}}_s) = \begin{cases} \mathbf{h}_t^T \bar{\mathbf{h}}_s & \text{dot} \\ \mathbf{h}_t^T \mathbf{W}_a \bar{\mathbf{h}}_s & \text{general} \\ \mathbf{v}_a^T \tanh(\mathbf{W}_a[\mathbf{h}_t; \bar{\mathbf{h}}_s]) & \text{concat} \end{cases} \quad (52)$$

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

$$\mathbf{a}_t(s) = \text{align}(\mathbf{h}_t, \bar{\mathbf{h}}_s) \quad (53)$$

$$= \frac{\exp(\text{score}(\mathbf{h}_t, \bar{\mathbf{h}}_s))}{\sum_{s'} \exp(\text{score}(\mathbf{h}_t, \bar{\mathbf{h}}_{s'}))} \quad (54)$$

3.  $\mathbf{c}_t$ : The weighted average over all source (encoder) hidden states<sup>6</sup>:

$$\mathbf{c}_t = \sum_{i=1}^{L_{enc}} \mathbf{a}_t(i) \bar{\mathbf{h}}_i \quad (55)$$

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

$$\tilde{\mathbf{h}}_t = \tanh(\mathbf{W}_c[\mathbf{c}_t; \mathbf{h}_t]) \quad (56)$$

**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>6</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.*

# CONDITIONAL RANDOM FIELDS

## CONTENTS

2.1	Conditional Random Fields . . . . .	21
2.2	An Introduction to Conditional Random Fields . . . . .	24
2.2.1	Inference (Sec. 4) . . . . .	28
2.2.2	Parameter Estimation (Sec. 5) . . . . .	31
2.2.3	Related Work and Future Directions (Sec. 6) . . . . .	34
2.3	Hidden-Unit Conditional Random Fields . . . . .	35
2.3.1	Detailed Derivations . . . . .	37
2.4	Pre-training of Hidden-Unit CRFs . . . . .	42
2.5	Neural Conditional Random Fields . . . . .	44
2.6	Bidirectional LSTM-CRF Models for Sequence Tagging . . . . .	46
2.7	On Herding and the Perceptron Cycling Theorem . . . . .	47

## 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>7</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:*

$$\Pr[Y_v | X, Y_w, w \neq v] = \Pr[Y_v | X, Y_w, w \sim v] \quad (57)$$

*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>8</sup> are the edges and vertices. Take note that  $\mathbf{X}$  is not a member of the vertices

---

<sup>7</sup>We assume all components  $\mathbf{Y}_i$  can only take on values in some finite label set  $\mathcal{Y}$ .

<sup>8</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:

$$p_\theta(\mathbf{y} \mid \mathbf{x}) \propto \exp \left( \sum_{e \in E, k} \lambda_k f_k(e, \mathbf{y}|_e, \mathbf{x}) + \sum_{v \in V, k} \mu_k g_k(v, \mathbf{y}|_v, \mathbf{x}) \right) \quad (58)$$

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  $\mathbf{M}_t(\mathbf{x})$

$$\mathbf{M}_t(y', y, \mid \mathbf{x}) = \exp(\Lambda_t(y', y \mid x)) \quad (59)$$

$$\Lambda_t(y', y \mid x) = \sum_k \lambda_k f_k(y', y, \mathbf{x}) + \sum_k \mu_k g_k(y, \mathbf{x}) \quad (60)$$

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:

$$Z_\theta(\mathbf{x}) = \left[ \prod_{t=1}^{T+1} \mathbf{M}_t(\mathbf{x}) \right]_{y_0, y_{T+1}} \quad (61)$$

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:

$$p_\theta(\mathbf{y} \mid \mathbf{x}) = \frac{\prod_{t=1}^{T+1} \mathbf{M}_t(y_{t-1}, y_t \mid \mathbf{x})}{\left[ \prod_{t=1}^{T+1} \mathbf{M}_t(\mathbf{x}) \right]_{y_0, y_{T+1}}} \quad (62)$$

**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>9</sup>. Their recurrence relations are

$$\alpha_t(\mathbf{x}) = \alpha_{t-1}(\mathbf{x}) \mathbf{M}_t(\mathbf{x}) \quad (63)$$

$$\beta_t(\mathbf{x})^T = \mathbf{M}_{t+1}(\mathbf{x}) \beta_{t+1}(\mathbf{x}) \quad (64)$$

---

<sup>9</sup>Remember that  $y_0$  and  $y_{T+1}$  are their own fixed symbolic constants representing a fixed start/stop state.

## 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

$$p(y) = \frac{1}{Z} \prod_{a=1}^A \psi_a(\mathbf{y}_a) \quad (65)$$

$$Z = \sum_{\mathbf{y}} \prod_{a=1}^A \psi_a(\mathbf{y}_a) \quad (66)$$

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>10</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,

$$p(y, \mathbf{x}) = p(y) \prod_{k=1}^K p(x_k | y) \quad (67)$$

$$p(y | \mathbf{x}) = \frac{1}{Z(\mathbf{x})} \exp \left( \sum_{k=1}^K \theta_k f_k(y, \mathbf{x}) \right) \quad (68)$$

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>10</sup>i.e. a particular set of factors.

An example generative model for sequence prediction is the **HMM**. Recall that an HMM defines

$$p(\mathbf{y}, \mathbf{x}) = \prod_{t=1}^T p(y_t | y_{t-1})p(x_t | y_t) \quad (69)$$

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

$$p(y, \mathbf{x}) = \frac{\exp(\sum_k \theta_k f_k(y, \mathbf{x}))}{\sum_{\tilde{y}, \tilde{\mathbf{x}}} \exp(\sum_k \theta_k f_k(\tilde{y}, \tilde{\mathbf{x}}))} \quad (70)$$

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.

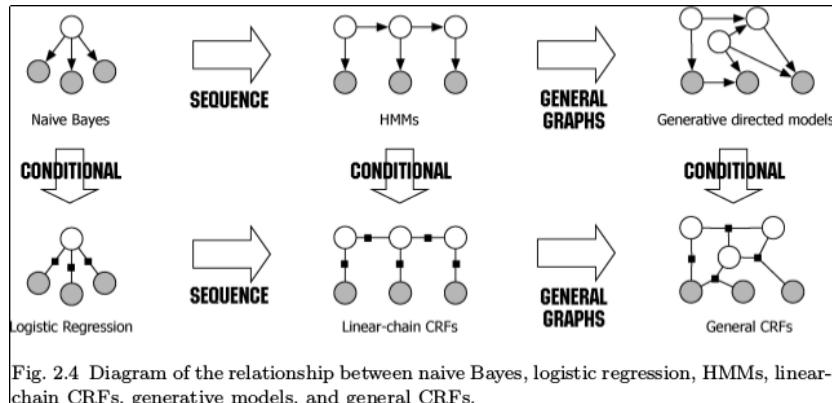


Fig. 2.4 Diagram of the relationship between naive Bayes, logistic regression, HMMs, linear-chain CRFs, generative models, and general CRFs.

**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:

$$p(\mathbf{y}, \mathbf{x}) = \frac{1}{Z} \prod_{t=1}^T \exp \left( \sum_{i,j \in S} \theta_{i,j} \mathbb{1}_{\{y_t=i, y_{t-1}=j\}} + \sum_{i \in S, o \in O} \mu_{o,i} \mathbb{1}_{\{y_t=i, x_t=o\}} \right) \quad (71) \quad \text{HMM joint distribution}$$

$$= \frac{1}{Z} \prod_{t=1}^T \exp \left( \sum_{k=1}^K \theta_k f_k(y_t, y_{t-1}, x_t) \right) \quad (72)$$

and the latter provides the more compact notation<sup>11</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:

$$p(\mathbf{y} \mid \mathbf{x}) = \frac{1}{Z(\mathbf{x})} \prod_{t=1}^T \exp \left( \sum_k \theta_k f_k(y_t, y_{t-1}, \mathbf{x}_t) \right) \quad (73)$$

$$Z(\mathbf{x}) = \sum_{\mathbf{y}} \prod_{t=1}^T \exp \left( \sum_k \theta_k f_k(y_t, y_{t-1}, \mathbf{x}_t) \right) \quad (74)$$

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:

$$\psi_t(y_t, y_{t-1}, x_t) = \exp \left( \sum_k \theta_k f_k(y_t, y_{t-1}, x_t) \right) \quad (75)$$

**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

$$p(\mathbf{y} \mid \mathbf{x}) = \frac{1}{Z(\mathbf{x})} \prod_{a=1}^A \psi_a(\mathbf{y}_a, \mathbf{x}_a) \quad (76)$$

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

$$p(\mathbf{y} \mid \mathbf{x}) = \frac{1}{Z(\mathbf{x})} \prod_{\psi_a \in F} \exp \left( \sum_{k=1}^{K(a)} \theta_{a,k} f_{a,k}(\mathbf{y}_a, \mathbf{x}_a) \right) \quad (77)$$

---

<sup>11</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>12</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

$$p(\mathbf{y} | \mathbf{x}) = \frac{1}{Z(\mathbf{x})} \prod_{C_p \in \mathcal{C}} \prod_{\psi_c \in C_p} \psi_c(\mathbf{x}_c, \mathbf{y}_c; \boldsymbol{\theta}_p) \quad (78)$$

$$= \frac{1}{Z(\mathbf{x})} \prod_{C_p \in \mathcal{C}} \prod_{c \in C_p} \exp \left\{ \sum_{k=1}^{K(p)} \theta_{p,k} f_{p,k}(\mathbf{x}_c, \mathbf{y}_c) \right\} \quad (79)$$

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:

$$f_{p,k}(\mathbf{y}_c, \mathbf{x}_c) = \mathbf{1}_{\{\mathbf{y}_c = \tilde{\mathbf{y}}_c\}} q_{p,k}(\mathbf{x}_c) \quad (80)$$

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.

$$[\text{edge-obs}] \quad f(y_t, y_{t-1}, \mathbf{x}_t) = q_m(\mathbf{x}_t) \mathbf{1}_{y_t=y, y_{t-1}=y'} \quad (\forall y, y' \in \mathcal{Y}, \forall m) \quad (81)$$

$$[\text{node-obs}] \quad f(y_t, y_{t-1}, \mathbf{x}_t) = \mathbf{1}_{y_t=y, y_{t-1}=y'} \quad (\forall y, y' \in \mathcal{Y}) \quad (82)$$

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>13</sup>.

- **Feature Induction.** The model begins with a number of base features, and the training procedure adds conjunctions of those features.

---

<sup>12</sup>Note how, for CRFs, the actual parameters  $\theta$  are tightly coupled with the feature functions  $f$ .

<sup>13</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.

---

### 2.2.1 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:

$$p(\mathbf{x}) = \sum_{\mathbf{y}} p(\mathbf{x}, \mathbf{y}) \quad (83) \quad \text{We define the } \psi_t \text{ as } p(y_t | y_{t-1})p(x_t | y_t)$$

$$= \sum_{\mathbf{y}} \prod_{t=1}^T \psi_t(y_t, y_{t-1}, x_t) \quad (84)$$

$$= \sum_{y_1} p(y_1) p(x_1 | y_1) \sum_{y_2} p(y_2 | y_1) p(x_2 | y_2) \sum_{y_3} \cdots \sum_{y_T} p(y_T | y_{T-1}) p(x_T | y_T) \quad (85)$$

$$= \sum_{y_1} \psi_1(y_1, x_1) \sum_{y_2} \cdots \sum_{y_T} \psi_T(y_T, y_{T-1}, x_T) \quad (86)$$

$$= \sum_T \sum_{T-1} \psi_T(y_T, y_{T-1}, x_T) \sum_{y_{T-2}} \cdots \sum_{y_1} \psi_1(y_1, x_1) \quad (87)$$

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  $M$ :

$$\alpha_t(j) \triangleq p(\mathbf{x}_{\langle 1 \dots t \rangle}, y_t = j) \quad (88)$$

$$= \sum_{\mathbf{y}_{\langle 1 \dots t-1 \rangle}} p(\mathbf{x}_{\langle 1 \dots t \rangle}, \mathbf{y}_{\langle 1 \dots t-1 \rangle}, y_t = j) \quad (89)$$

$$= \sum_{\mathbf{y}_{\langle 1 \dots t-1 \rangle}} \psi_t(j, y_{t-1}, x_t) \prod_{t'=1}^{t'-1} \psi_{t'}(y_{t'}, y_{t-1}, x_{t'}) \quad (90)$$

$$= \sum_{i \in S} \psi_t(j, i, x_t) \sum_{\mathbf{y}_{\langle 1 \dots t-2 \rangle}} \psi_{t-1}(y_{t-1}, y_{t-2}, x_{t-1}) \prod_{t'=1}^{t-2} \psi_{t'}(y_{t'}, y_{t-1}, x_{t'}) \quad (91)$$

$$= \sum_{i \in S} \psi_t(j, i, x_t) \alpha_{t-1}(i) \quad (92)$$

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

$$p(\mathbf{x}) = \sum_{j \in S} \alpha_T(j) \quad (93)$$

Notice how in the step from equation 88 to 89, 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  $t$ ).

$$\beta_t(i) \triangleq p(\mathbf{x}_{\langle t+1 \dots T \rangle} \mid y_t = i) \quad (94)$$

$$= \sum_{\mathbf{y}_{\langle t+1 \dots T \rangle}} \psi_{t+1}(y_{t+1}, i, x_{t+1}) \prod_{t'=t+2}^T \psi_{t'}(y_{t'}, y_{t'-1}, x_{t'}) \quad (95)$$

$$= \sum_{y_{t+1}} \psi_{t+1}(y_{t+1}, i, x_{t+1}) \beta_{t+1}(y_{t+1}) \quad (96)$$

We initialize  $\beta_T(i) = 1$ .

Similar to how we obtained equation 93, we can rewrite  $p(\mathbf{x})$  in terms of the  $\beta$ :

$$p(\mathbf{x}) = \beta_0(y_0) \triangleq \sum_{y_1} \psi_1(y_1, y_0, x_1) \beta_1(y_1) \quad (97)$$

We can then combine the definition for  $\alpha$  and  $\beta$  to compute marginals of the form  $p(y_{t-1}, y_t \mid \mathbf{x})$ :

$$p(y_{t-1}, y_t \mid \mathbf{x}) = \frac{1}{p(\mathbf{x})} \alpha_{t-1}(y_{t-1}) \psi_t(y_t, y_{t-1}, x_t) \beta_t(y_t) \quad (98)$$

$$\text{where } p(\mathbf{x}) = \sum_{j \in S} \alpha_T(j) = \beta_0(y_0) \quad (99)$$

In summary, the forward-backward algorithm consists of the following steps:

1. Compute  $\alpha_t$  for all  $t$  using equation 93.
2. Compute  $\beta_t$  for all  $t$  using equation 97.
3. Return the marginal distributions computed from equation 98.

**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 87 with maximization. The analog of  $\alpha$  for viterbi are defined as:

$$\delta_t(j) = \max_{\mathbf{y}_{(1 \dots t-1)}} \psi_t(j, y_{t-1}, x_t) \prod_{t'=1}^{t-1} \psi_{t'}(y_{t'}, y_{t'-1}, x_{t'}) \quad (100)$$

$$= \max_{i \in S} \psi_t(j, i, x_t) \delta_{t-1}(i) \quad (101)$$

and the maximizing assignment is computed by a backwards recursion,

$$y_T^* = \arg \max_{i \in S} \delta_T(i) \quad (102)$$

$$y_t^* = \arg \max_{i \in S} \psi_t(y_{t+1}^*, i, x_{t+1}) \delta_t(i) \text{ for } t < T \quad (103)$$

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

$$p(\mathbf{y} \mid \mathbf{x}) = \frac{1}{Z(\mathbf{x})} \prod_{t=1}^T \psi_t(y_t, y_{t-1}, x_t) \quad (104)$$

$$\text{where } \psi_t(y_t, y_{t-1}, x_t) = \exp \left\{ \sum_k \theta_k f_k(y_t, y_{t-1}, x_t) \right\} \quad (105)$$

and the results for the forward-backward algorithm become

$$p(y_{t-1}, y_t \mid \mathbf{x}) = \frac{1}{Z(\mathbf{x})} \alpha_{t-1}(y_{t-1}) \psi_t(y_t, y_{t-1}, x_t) \beta_t(y_t) \quad (106)$$

$$p(y_t \mid \mathbf{x}) = \frac{1}{Z(\mathbf{x})} \alpha_t(y_t) \beta_t(y_t) \quad (107)$$

$$\text{where } Z(\mathbf{x}) = \sum_{j \in S} \alpha_T(j) = \beta_0(y_0) \quad (108)$$

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

$$\alpha_t(j) = \sum_{\mathbf{y}_{(1 \dots t-1)}} \exp \left\{ \sum_k \theta_k f_k(j, y_{t-1}, x_t) \right\} \prod_{t'=1}^{t-1} \exp \left\{ \sum_k \theta_k f_k(y_{t'}, y_{t'-1}, x_{t'}) \right\} \quad (109)$$

$$\beta_t(i) = \sum_{\mathbf{y}_{(t+1 \dots T)}} \exp \left\{ \sum_k \theta_k f_k(y_{t+1}, i, x_{t+1}) \right\} \prod_{t'=t+2}^T \exp \left\{ \sum_k \theta_k f_k(y_{t'}, y_{t'-1}, x_{t'}) \right\} \quad (110)$$

$$\delta_t(j) = \max_{\mathbf{y}_{(1 \dots t-1)}} \exp \left\{ \sum_k \theta_k f_k(j, y_{t-1}, x_t) \right\} \prod_{t'=1}^{t-1} \exp \left\{ \sum_k \theta_k f_k(y_{t'}, y_{t'-1}, x_{t'}) \right\} \quad (111)$$

**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

$$p(\mathbf{y}) = Z^{-1} \prod_{a \in F} \psi_a(\mathbf{y}_a) \quad (112)$$

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.

### 2.2.2 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:

$$\ell(\boldsymbol{\theta}) = \sum_{i=1}^N \log p(\mathbf{y}^{(i)} | \mathbf{x}^{(i)}; \boldsymbol{\theta}) - \frac{1}{2\sigma^2} \sum_{k=1}^K \theta_k^2 \quad (113)$$

$$= \sum_{i=1}^N \sum_{t=1}^T \sum_{k=1}^K \theta_k f_k(y_t^{(i)}, y_{t-1}^{(i)}, \mathbf{x}_t^{(i)}) - \sum_{i=1}^N \log Z(\mathbf{x}^{(i)}) - \frac{1}{2\sigma^2} \sum_{k=1}^K \theta_k^2 \quad (114)$$

with regularization parameter  $1/2\sigma^2$ . The partial derivatives are

$$\frac{\partial \ell}{\partial \theta_k} = \sum_{i,t} f_k(y_t^{(i)}, y_{t-1}^{(i)}, \mathbf{x}_t^{(i)}) - \sum_{i,t,y,y'} f_k(y, y', \mathbf{x}_t^{(i)}) p(y, y' | \mathbf{x}^{(i)}) - \frac{\theta_k}{\sigma^2} \quad (115)$$

and the derivation for the partial derivative of  $\log(Z(x))$  is

$$\frac{\partial \log Z(x)}{\partial \theta_k} = \frac{1}{Z(x)} \frac{\partial}{\partial \theta_k} \left[ \sum_{\mathbf{y}_{\langle 1 \dots T \rangle}} \prod_{t=1}^T \exp \left\{ \sum_k \theta_k f_k(y_t, y_{t-1}, x_t) \right\} \right] \quad (116)$$

$$= \frac{1}{Z(x)} \sum_{\mathbf{y}_{\langle 1 \dots T \rangle}} \frac{\partial}{\partial \theta_k} \exp \left\{ \sum_t \sum_k \theta_k f_k(y_t, y_{t-1}, x_t) \right\} \quad (117)$$

$$= \frac{1}{Z(x)} \sum_{\mathbf{y}_{\langle 1 \dots T \rangle}} \sum_t f_k(y_t, y_{t-1}, x_t) \exp \left\{ \sum_t \sum_k \theta_k f_k(y_t, y_{t-1}, x_t) \right\} \quad (118)$$

$$= \sum_t \sum_{y_t} \sum_{y_{t-1}} f_k(y_t, y_{t-1}, x_t) \left[ \sum_{\mathbf{y}_{\langle 1 \dots t-2 \rangle}} \sum_{\mathbf{y}_{\langle t+1 \dots T \rangle}} \frac{1}{Z(x)} \exp \left\{ \sum_t \sum_k \theta_k f_k(y_t, y_{t-1}, x_t) \right\} \right] \quad (119)$$

$$= \sum_t \sum_y \sum_{y'} f_k(y, y', x_t) p(y_t = y, y_{t-1} = y' | x) \quad (120)$$

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*.

$$\frac{\partial \ell}{\partial \theta_k} = \mathbb{E}_{\mathbf{x}, \mathbf{y} \sim \tilde{p}(\mathbf{y}, \mathbf{x})} \left[ \sum_t f_k(y_t, y_{t-1}, x_t) \right] - \mathbb{E}_{\mathbf{x}, \mathbf{y} \sim \hat{p}(\mathbf{y}, \mathbf{x})} \left[ \sum_t f_k(y_t, y_{t-1}, x_t) \right] \quad (121)$$

### 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:

$$\alpha_t(j) = \sum_{i \in S} \exp \left\{ \sum_k \theta_k f_k(j, i, x_t) \right\} \alpha_{t-1}(i) \quad (122)$$

$$\beta_t(j) = \sum_{i \in S} \exp \left\{ \sum_k \theta_k f_k(i, j, x_{t+1}) \right\} \beta_{t+1}(i) \quad (123)$$

3. Compute the marginals:

$$p(y_t, y_{t-1}, | \mathbf{x}) = \frac{1}{Z(\mathbf{x})} \alpha_{t-1}(y_{t-1}) \psi_t(y_t, y_{t-1}, x_t) \beta_t(y_t) \quad (124)$$

$$p(y_t | \mathbf{x}) = \frac{1}{Z(\mathbf{x})} \alpha_t(y_t) \beta_t(y_t) \quad (125)$$

$$\text{where } Z(\mathbf{x}) = \sum_{j \in S} \alpha_T(j) = \beta_0(y_0) \quad (126)$$

**Gradients.** For each parameter  $\theta_k$ , compute:

$$\frac{\partial \ell}{\partial \theta_k} = \sum_{i,t} f_k(y_t^{(i)}, y_{t-1}^{(i)}, \mathbf{x}_t^{(i)}) - \sum_{i,t,y,y'} f_k(y, y', \mathbf{x}_t^{(i)}) p(y, y' | \mathbf{x}^{(i)}) - \frac{\theta_k}{\sigma^2} \quad (127)$$

**CRF with latent variables.** Now we have additional latent variables  $\mathbf{z}$ :

$$p(\mathbf{y}, \mathbf{z} | \mathbf{x}) = \frac{1}{Z(\mathbf{x})} \prod_t \psi_t(y_t, y_{t-1}, \mathbf{z}_t, \mathbf{x}_t) \quad (128)$$

Since we observe  $\mathbf{y}$  during training, what if we instead treat this as a CRF with both  $\mathbf{x}$  and  $\mathbf{y}$  observed?

$$p(\mathbf{z} | \mathbf{y}, \mathbf{x}) = \frac{1}{Z(\mathbf{y}, \mathbf{x})} \prod_t \psi_t(y_t, y_{t-1}, \mathbf{z}_t, \mathbf{x}_t) \quad (129)$$

$$Z(\mathbf{y}, \mathbf{x}) = \sum_{\mathbf{z}} \prod_t \psi_t(y_t, y_{t-1}, \mathbf{z}_t, \mathbf{x}_t) \quad (130)$$

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

$$p(\mathbf{y} \mid \mathbf{x}) = \frac{1}{Z(\mathbf{x})} \sum_{\mathbf{z}} \prod_t \psi_t(y_t, y_{t-1}, \mathbf{z}_t, \mathbf{x}_t) = \frac{Z(\mathbf{y}, \mathbf{x})}{Z(\mathbf{x})} \quad (131)$$

Finally, we can write the gradient as<sup>14</sup>

$$\frac{\partial \ell}{\partial \theta_k} = \sum_{\mathbf{z}} p(\mathbf{z} \mid \mathbf{y}, \mathbf{x}) \frac{\partial}{\partial \theta_k} [\log p(\mathbf{y}, \mathbf{z} \mid \mathbf{x})] \quad (132)$$

$$= \sum_t \sum_{\mathbf{z}_t} \left[ p(\mathbf{z}_t \mid \mathbf{y}, \mathbf{x}) f_k(y_t, y_{t-1}, \mathbf{z}_t, \mathbf{x}_t) - \sum_{y, y'} p(\mathbf{z}_t, y, y' \mid \mathbf{x}_t) f_k(y_t, y_{t-1}, \mathbf{z}_t, \mathbf{x}_t) \right] \quad (133)$$

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:

$$\frac{\partial \ell_i}{\partial \theta_k} = \sum_t f_k(y_t^{(i)}, y_{t-1}^{(i)}, \mathbf{x}_t^{(i)}) - \sum_{t, y, y'} f_k(y, y', \mathbf{x}_t^{(i)}) p(y, y' \mid \mathbf{x}^{(i)}) - \frac{\theta_k}{N\sigma^2} \quad (134)$$

which corresponds to parameter update (remember: we are using the LL, not the negative LL):

$$\theta^{(m)} = \theta^{(m-1)} + \alpha_m \nabla \ell_i(\theta^{(m-1)}) \quad (135)$$

where  $m$  denotes this is the  $m$ th update of the training process.

---

<sup>14</sup>This uses the trick

$$\frac{df}{d\theta} = f(\theta) \frac{d \log f}{d\theta}$$

---

### 2.2.3 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

$$p_{MEMM}(\mathbf{y} \mid \mathbf{x}) = \prod_{t=1}^T p(y_t \mid y_{t-1}, \mathbf{x}) \quad (136)$$

$$p(y_t \mid y_{t-1}, \mathbf{x}) = \frac{\exp \left\{ \sum_{k=1}^K \theta_k f_k(y_t, y_{t-1}, \mathbf{x}_t) \right\}}{Z_t(y_{t-1}, \mathbf{x})} \quad (137)$$

$$Z_t(y_{t-1}, \mathbf{x}) = \sum_{y'} \exp \left\{ \sum_{k=1}^K \theta_k f_k(y_t, y_{t-1}, \mathbf{x}_t) \right\} \quad (138)$$

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:

$$y_{ML}^* \leftarrow \arg \max_y p(y \mid x; \hat{\theta}) \quad (139)$$

$$y_B^* \leftarrow \arg \max_y \mathbb{E}_{\theta \sim p(\theta|x)}(p(y \mid x; \theta)) \quad (140)$$

Unfortunately, computing the exact integral (the expectation) is usually intractable, and we must resort to approximate methods like MCMC.

## 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>15</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

$$p(\mathbf{y} \mid \mathbf{x}) = \frac{1}{Z(\mathbf{x})} \sum_{\mathbf{z}} \exp(E(\mathbf{x}, \mathbf{z}, \mathbf{y})) \quad (141)$$

$$\begin{aligned} E(\mathbf{x}, \mathbf{z}, \mathbf{y}) = & \sum_{t=2}^T [\mathbf{y}_{t-1}^T \mathbf{A} \mathbf{y}_t] + \sum_{t=1}^T [\mathbf{x}_t^T \mathbf{W} \mathbf{z}_t + \mathbf{z}_t^T \mathbf{V} \mathbf{y}_t + \mathbf{b}^T \mathbf{z}_t + \mathbf{c}^T \mathbf{y}_t] \\ & + \mathbf{y}_1^T \boldsymbol{\pi} + \mathbf{y}_T^T \boldsymbol{\tau} \end{aligned} \quad (142)$$

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>15</sup>See the introduction in my CRF notes for recap of label-bias.

shown below:

$$p(\mathbf{y} \mid \mathbf{x}) = \frac{\exp\{\mathbf{y}_1^T \boldsymbol{\pi} + \mathbf{y}_T^T \boldsymbol{\tau}\}}{Z(\mathbf{x})} \prod_{t=1}^T \left[ \exp\{\mathbf{c}^T \mathbf{y}_t + \mathbf{y}_{t-1}^T \mathbf{A} \mathbf{y}_t\} \right. \\ \left. \prod_{h=1}^H \sum_{z_h \in \{0,1\}} \exp\{z_h b_h + z_h \mathbf{w}_h^T \mathbf{x}_t + z_h \mathbf{v}_h^T \mathbf{y}_t\} \right] \quad (143)$$

$$= \frac{\exp\{\mathbf{y}_1^T \boldsymbol{\pi} + \mathbf{y}_T^T \boldsymbol{\tau}\}}{Z(\mathbf{x})} \prod_{t=1}^T \left[ \exp\{\mathbf{c}^T \mathbf{y}_t + \mathbf{y}_{t-1}^T \mathbf{A} \mathbf{y}_t\} \right. \\ \left. \prod_{h=1}^H \left( 1 + \exp\{b_h + \mathbf{w}_h^T \mathbf{x}_t + \mathbf{v}_h^T \mathbf{y}_t\} \right) \right] \quad (144)$$

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  $\mathbf{z}$ :

$$p(y_t, y_{t-1} \mid \mathbf{x}) \propto \alpha_{t-1}(y_{t-1}) \left[ \sum_{\mathbf{z}_t} \Psi_t(\mathbf{x}_t, \mathbf{z}_t, y_t, y_{t-1}) \right] \beta_t(y_t) \quad (145)$$

$$p(y_t \mid \mathbf{x}) \propto \alpha_t(y_t) \beta_t(y_t) \quad (146)$$

$$(147)$$

$$\alpha_t(j) = \sum_{i \in \mathcal{Y}} \sum_{\mathbf{z}_t} \Psi_t(\mathbf{x}_t, \mathbf{z}_t, j, i) \alpha_{t-1}(i) \quad (148)$$

$$\beta_t(j) = \sum_{i \in \mathcal{Y}} \sum_{\mathbf{z}_{t+1}} \Psi_{t+1}(\mathbf{x}_{t+1}, \mathbf{z}_{t+1}, i, j) \beta_{t+1}(i) \quad (149)$$

**Training.** The conditional log likelihood for a single example  $(\mathbf{x}, \mathbf{y})$  is (bias and initial-state terms omitted)

$$\mathcal{L} = \log p(\mathbf{y} \mid \mathbf{x}) \quad (150)$$

$$= \sum_{t=1}^T \log \left( \sum_{\mathbf{z}_t} \Psi_t(\mathbf{x}_t, \mathbf{z}_t, y_{t-1}, y_t) \right) - \log Z(\mathbf{x}) \quad (151)$$

$$\text{where } \Psi_t := \exp\{y_{t-1}^T \mathbf{A} \mathbf{y}_t + \mathbf{x}_t^T \mathbf{W} \mathbf{z}_t + \mathbf{z}_t^T \mathbf{V} \mathbf{y}_t\} \quad (152)$$

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>16</sup> parameters  $v \in \Upsilon$  is given by

$$\frac{\partial \mathcal{L}}{\partial v} = \sum_{t=1}^T \left[ \sum_{k \in \mathcal{Y}} \left( (\mathbb{1}_{y_t=k} - p(y_t=k \mid \mathbf{x})) \sum_{h=1}^H \sigma(o_{hk}(\mathbf{x}_t)) \frac{\partial o_{hk}(\mathbf{x}_t)}{\partial v} \right) \right] \quad (153)$$

$$\text{where } o_{hk}(\mathbf{x}_t) = b_h + c_k + V_{hk} + \mathbf{w}_h^T \mathbf{x}_t \quad (154)$$

Unfortunately, the negative CLL is **non-convex**, and so we are only guaranteed to converge to a *local* maximum of the CLL.

---

<sup>16</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  $\boldsymbol{\tau}$ .

---

### 2.3.1 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>17</sup>

$$p(\mathbf{y} \mid \mathbf{x}) = \frac{\exp\{\mathbf{y}_1^T \boldsymbol{\pi} + \mathbf{y}_T^T \boldsymbol{\tau}\}}{Z(\mathbf{x})} \prod_{t=1}^T \left[ \exp\{\mathbf{c}^T \mathbf{y}_t + \mathbf{y}_{t-1}^T \mathbf{A} \mathbf{y}_t\} \prod_{h=1}^H (1 + \exp\{b_h + \mathbf{w}_h^T \mathbf{x}_t + \mathbf{v}_h^T \mathbf{y}_t\}) \right] \quad (155)$$

$$NLL = - \sum_{i=1}^N \log p(\mathbf{y}^{(i)} \mid \mathbf{x}^{(i)}) \quad (156)$$

$$= - \sum_{i=1}^N \left[ \sum_{t=1}^T \log \left( \sum_{\mathbf{z}_t} \psi_t(\mathbf{x}_t, \mathbf{z}_t, \mathbf{y}_{t-1}, \mathbf{y}_t) \right) - \log Z(\mathbf{x}^{(i)}) \right] \quad (157)$$

The above formula for  $p(\mathbf{y} \mid \mathbf{x})$  implies something that will be very useful:

$$\sum_{\mathbf{z}_t} \psi_t(\mathbf{x}_t, \mathbf{z}_t, \mathbf{y}_{t-1}, \mathbf{y}_t) = \exp\{\mathbf{c}^T \mathbf{y}_t + \mathbf{y}_{t-1}^T \mathbf{A} \mathbf{y}_t\} \prod_{h=1}^H (1 + \exp\{b_h + \mathbf{w}_h^T \mathbf{x}_t + \mathbf{v}_h^T \mathbf{y}_t\}) \quad (158)$$

Using the generalization of the product rule for derivatives over N products, we can derive that

$$\frac{\partial}{\partial v} \prod_h (1 + \exp\{o(h)\}) = \left[ \prod_h (1 + \exp\{o(h)\}) \right] \left[ \sum_h \sigma(o(h)) \frac{\partial o(h)}{\partial v} \right] \quad (159)$$

Which means the derivatives of  $\sum_z \psi$  for the data-dependent params  $v_{dat}$  and transition params  $v_{tr}$ , are:

$$\frac{\partial \sum_{\mathbf{z}_t} \psi_t}{\partial v_{dat}} = \left[ \sum_h \sigma(o(h, y_t)) \frac{\partial o(h, y_t)}{\partial v_{dat}} \right] \sum_{\mathbf{z}_t} \psi_t \quad (160)$$

$$\frac{\partial \sum_{\mathbf{z}_t} \psi_t}{\partial v_{tr}} = \left[ \frac{\partial}{\partial v_{tr}} \mathbf{c}^T \mathbf{y}_t + \mathbf{y}_{t-1}^T \mathbf{A} \mathbf{y}_t \right] \sum_{\mathbf{z}_t} \psi_t \quad (161)$$

which also conveniently means that

$$\frac{\partial}{\partial v_{dat}} \log \left( \sum_{\mathbf{z}_t} \psi_t \right) = \sum_h \sigma(o(h, y_t)) \frac{\partial o(h, y_t)}{\partial v_{dat}} \quad (162)$$

$$\frac{\partial}{\partial v_{tr}} \log \left( \sum_{\mathbf{z}_t} \psi_t \right) = \frac{\partial}{\partial v_{tr}} \mathbf{c}^T \mathbf{y}_t + \mathbf{y}_{t-1}^T \mathbf{A} \mathbf{y}_t \quad (163)$$

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>17</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  $v$ :

$$\frac{\partial NLL}{\partial v} = - \sum_{i=1}^N \left[ \sum_{t=1}^T \frac{\partial}{\partial v} \log \left( \sum_{\mathbf{z}_t} \psi_t(\mathbf{x}_t^{(i)}, \mathbf{z}_t, \mathbf{y}_{t-1}^{(i)}, \mathbf{y}_t^{(i)}) \right) - \frac{\partial}{\partial v} \log Z(\mathbf{x}^{(i)}) \right] \quad (164)$$

$$\frac{\partial \log Z(\mathbf{x}^{(i)})}{\partial v} = \frac{1}{Z(\mathbf{x}^{(i)})} \frac{\partial}{\partial v} \sum_{\mathbf{y}_{\langle 1 \dots T \rangle}} \tilde{p}(\mathbf{y}_1, \dots, \mathbf{y}_T \mid \mathbf{x}^{(i)}) \quad (165)$$

$$\frac{\partial}{\partial v} \tilde{p}(\mathbf{y}_1, \dots, \mathbf{y}_T \mid \mathbf{x}^{(i)}) = \frac{\partial}{\partial v} \prod_t \sum_{\mathbf{z}_t} \psi_t \quad (166)$$

$$= \left[ \prod_t \sum_{\mathbf{z}_t} \psi_t \right] \left[ \sum_t \frac{\frac{\partial}{\partial v} \sum_{\mathbf{z}_t} \psi_t}{\sum_{\mathbf{z}_t} \psi_t} \right] \quad (167)$$

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.

$$\frac{\partial NLL}{\partial v_{dat}} = - \sum_{i=1}^N \left[ \sum_{t=1}^T \frac{\partial}{\partial v_{dat}} \log \left( \sum_{\mathbf{z}_t} \psi_t \right) - \frac{\partial}{\partial v} \log Z(\mathbf{x}^{(i)}) \right] \quad (168)$$

$$= - \sum_{i=1}^N \left[ \sum_t \sum_h \sigma(o(h, y_t)) \frac{\partial o(h, y_t)}{\partial v_{dat}} - \frac{1}{Z(\mathbf{x}^{(i)})} \sum_{\mathbf{y}_{\langle 1 \dots T \rangle}} \left[ \prod_t \sum_{\mathbf{z}_t} \psi_t \right] \left[ \sum_t \frac{\frac{\partial}{\partial v} \sum_{\mathbf{z}_t} \psi_t}{\sum_{\mathbf{z}_t} \psi_t} \right] \right] \quad (169)$$

$$= - \sum_{i=1}^N \left[ \sum_t \sum_h \sigma(o(h, y_t)) \frac{\partial o(h, y_t)}{\partial v_{dat}} - \frac{1}{Z(\mathbf{x}^{(i)})} \sum_{\mathbf{y}_{\langle 1 \dots T \rangle}} \left[ \prod_t \sum_{\mathbf{z}_t} \psi_t \right] \left[ \sum_t \sum_h \sigma(o(h, y_t)) \frac{\partial o(h, y_t)}{\partial v_{dat}} \right] \right] \quad (170)$$

$$= - \sum_{i=1}^N \left[ \sum_t \sum_h \sigma(o(h, y_t)) \frac{\partial o(h, y_t)}{\partial v_{dat}} - \sum_t \sum_y \sum_{y'} \left[ \sum_h \sigma(o(h, y_t)) \frac{\partial o(h, y_t)}{\partial v_{dat}} \right] \xi_{t,y,y'} \right] \quad (171)$$

$$= - \sum_{i=1}^N \left[ \sum_t \sum_h \sigma(o(h, y_t)) \frac{\partial o(h, y_t)}{\partial v_{dat}} - \sum_t \sum_y \left[ \sum_h \sigma(o(h, y_t)) \frac{\partial o(h, y_t)}{\partial v_{dat}} \right] \gamma_{t,y} \right] \quad (172)$$

$$= - \sum_{i=1}^N \left[ \sum_t \left( \sum_h \sigma(o(h, y_t)) \frac{\partial o(h, y_t)}{\partial v_{dat}} - \sum_y \left[ \sum_h \sigma(o(h, y_t)) \frac{\partial o(h, y_t)}{\partial v_{dat}} \right] \gamma_{t,y} \right) \right] \quad (173)$$

$$= - \sum_{i=1}^N \left[ \sum_t \sum_y \left( (\mathbb{1}_{y_t=y} - \gamma_{t,y}) \sum_h \sigma(o(h, y_t)) \frac{\partial o(h, y_t)}{\partial v_{dat}} \right) \right] \quad (174)$$

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

$$\frac{\partial NLL}{\partial v_{tr}} = - \sum_{i=1}^N \left[ \sum_{t=1}^T \frac{\partial}{\partial v_{tr}} \log \left( \sum_{\mathbf{z}_t} \psi_t \right) - \frac{\partial}{\partial v_{tr}} \log Z(\mathbf{x}^{(i)}) \right] \quad (175)$$

$$= - \sum_i^N \left[ \sum_t \frac{\partial}{\partial v_{tr}} [\mathbf{c}^T \mathbf{y}_t + \mathbf{y}_{t-1} \mathbf{A} \mathbf{y}_t] - \frac{\partial}{\partial v_{tr}} \log Z(\mathbf{x}^{(i)}) \right] \quad (176)$$

$$= - \sum_{i=1}^N \left[ \sum_t \frac{\partial}{\partial v_{tr}} [\mathbf{c}^T \mathbf{y}_t + \mathbf{y}_{t-1} \mathbf{A} \mathbf{y}_t] - \frac{1}{Z(\mathbf{x}^{(i)})} \sum_{\mathbf{y}_{(1 \dots T)}} \left[ \prod_t \sum_{\mathbf{z}_t} \psi_t \right] \left[ \sum_t \frac{\partial}{\partial v_{tr}} [\mathbf{c}^T \mathbf{y}_t + \mathbf{y}_{t-1} \mathbf{A} \mathbf{y}_t] \right] \right] \quad (177)$$

$$= - \sum_{i=1}^N \left[ \sum_t \sum_y \left( (\mathbb{1}_{y_t=y} - \gamma_{t,y}) \frac{\partial}{\partial v_{tr}} [\mathbf{c}^T \mathbf{y}_t + \mathbf{y}_{t-1} \mathbf{A} \mathbf{y}_t] \right) \right] \quad (178)$$

### Boundary parameters

$$\frac{\partial NLL}{\partial \pi_\ell} = - \sum_{i=1}^N [\mathbb{1}_{y_1=\ell} - \gamma_{1,\ell}] \quad (179)$$

$$\frac{\partial NLL}{\partial \tau_\ell} = - \sum_{i=1}^N [\mathbb{1}_{y_T=\ell} - \gamma_{T,\ell}] \quad (180)$$

The results of each of the boxes above are summarized below, for the case of  $N = 1$  to save space.

$$\frac{\partial NLL}{\partial v_{dat}} = - \sum_t \sum_y \left( (\mathbb{1}_{y_t=y} - \gamma_{t,y}) \sum_h \sigma(o(h,y)) \frac{\partial o(h,y)}{\partial v_{dat}} \right) \quad (181)$$

$$\frac{\partial NLL}{\partial v_{tr}} = - \sum_t \sum_y \left( (\mathbb{1}_{y_t=y} - \gamma_{t,y}) \frac{\partial}{\partial v_{tr}} [\mathbf{c}^T \mathbf{y}_t + \mathbf{y}_{t-1} \mathbf{A} \mathbf{y}_t] \right) \quad (182)$$

$$\frac{\partial NLL}{\partial \pi_\ell} = - \sum_{i=1}^N [\mathbb{1}_{y_1=\ell} - \gamma_{1,\ell}] \quad (183)$$

$$\frac{\partial NLL}{\partial \tau_\ell} = - \sum_{i=1}^N [\mathbb{1}_{y_T=\ell} - \gamma_{T,\ell}] \quad (184)$$

Now I'll further go through and show how the equations simplify for each type of data-

dependent parameter.

$$\frac{\partial NLL}{\partial W_{c,h}} = - \sum_t \sum_y \left( (\mathbb{1}_{y_t=y} - \gamma_{t,y}) \sum_{h'} \sigma(o(h', y)) \frac{\partial}{\partial W_{c,h}} (b_{h'} + c_y + V_{h',y} + \mathbf{w}_{h'}^T \mathbf{x}_t) \right) \quad (185)$$

$$= - \sum_t \sum_y \left( (\mathbb{1}_{y_t=y} - \gamma_{t,y}) \sum_{h'} \sigma(o(h, y)) \mathbb{1}_{h=h'} \mathbb{1}_{c \in \mathbf{x}_t} \right) \quad (186)$$

$$= - \sum_t \sum_y (\mathbb{1}_{y_t=y} - \gamma_{t,y}) \sigma(o(h, y)) \mathbb{1}_{c \in \mathbf{x}_t} \quad (187)$$

$$\frac{\partial NLL}{\partial V_{h,y}} = - \sum_t (\mathbb{1}_{y_t=y} - \gamma_{t,y}) \sigma(o(h, y)) \quad (188)$$

$$\frac{\partial NLL}{\partial b_h} = - \sum_t \sum_y (\mathbb{1}_{y_t=y} - \gamma_{t,y}) \sigma(o(h, y)) \quad (189)$$

$$\frac{\partial NLL}{\partial c_y} = - \sum_t (\mathbb{1}_{y_t=y} - \gamma_{t,y}) \sum_h \sigma(o(h, y)) \quad (190)$$

$$(191)$$

**Alternative Approach.** The above was a bit more cumbersome than needed. I'll now derive it in an easier way.

$$p(\mathbf{y} \mid \mathbf{x}) = \frac{\exp\{\mathbf{y}_1^T \boldsymbol{\pi} + \mathbf{y}_T^T \boldsymbol{\tau}\}}{Z(\mathbf{x})} \prod_{t=1}^T \left[ \exp\{\mathbf{c}^T \mathbf{y}_t + \mathbf{y}_{t-1}^T \mathbf{A} \mathbf{y}_t\} \prod_{h=1}^H (1 + \exp\{b_h + \mathbf{w}_h^T \mathbf{x}_t + \mathbf{v}_h^T \mathbf{y}_t\}) \right] \quad (192)$$

$$= \frac{\exp\{I+T\}}{Z(\mathbf{x})} \prod_{t=1}^T \prod_{h=1}^H (1 + \exp\{b_h + \mathbf{w}_h^T \mathbf{x}_t + \mathbf{v}_h^T \mathbf{y}_t\}) \quad (193)$$

$$NLL = - \sum_{i=1}^N \log p(\mathbf{y}^{(i)} \mid \mathbf{x}^{(i)}) \quad (194)$$

$$= - \sum_{i=1}^N \left[ I + T + \sum_{t=1}^T \sum_{h=1}^H [\log(1 + \exp\{b_h + \mathbf{w}_h^T \mathbf{x}_t + \mathbf{v}_h^T \mathbf{y}_t\})] - \log Z(\mathbf{x}^{(i)}) \right] \quad (195)$$

Now, focusing on the regular log-likelihood for a single example, we have

$$\frac{\partial \mathcal{L}_i}{\partial v} = \frac{\partial}{\partial v} \log p(\mathbf{y}^{(i)}, \mathbf{x}^{(i)}) \quad (196)$$

$$= \frac{\partial}{\partial v} \left[ I + T + \sum_{t,h} \log(1 + \exp\{b_h + \mathbf{w}_h^T \mathbf{x}_t + \mathbf{v}_h^T \mathbf{y}_t\}) - \log Z(\mathbf{x}^{(i)}) \right] \quad (197)$$

$$\frac{\partial \log Z(\mathbf{x}^{(i)})}{\partial v} = \frac{1}{Z(\mathbf{x}^{(i)})} \sum_{\mathbf{y}_{\langle 1 \dots T \rangle}} \frac{\partial}{\partial v} \tilde{p}(\mathbf{y} \mid \mathbf{x}^{(i)}) \quad (198)$$

$$= \frac{1}{Z(\mathbf{x}^{(i)})} \sum_{\mathbf{y}_{\langle 1 \dots T \rangle}} \tilde{p}(\mathbf{y} \mid \mathbf{x}^{(i)}) \frac{\partial}{\partial v} \log \tilde{p}(\mathbf{y} \mid \mathbf{x}^{(i)}) \quad (199)$$

as our base formula for partial derivatives.

### Transition parameters

$$\frac{\partial \mathcal{L}_i}{\partial A_{i,j}} = \frac{\partial}{\partial A_{i,j}} \left[ I + T + \sum_{t,h} \log \left( 1 + \exp \left\{ b_h + \mathbf{w}_h^T \mathbf{x}_t + \mathbf{v}_h^T \mathbf{y}_t \right\} \right) - \log Z(\mathbf{x}^{(i)}) \right] \quad (200)$$

$$= \frac{\partial}{\partial A_{i,j}} [I + T] - \sum_{\mathbf{y}_{\langle 1 \dots T \rangle}} p(\mathbf{y} \mid \mathbf{x}^{(i)}) \frac{\partial}{\partial A_{i,j}} \left[ y_1^T \pi + y_T^T \tau + \sum_t y_{t-1} A y_t \right] \quad (201)$$

$$= \sum_t \mathbb{1}_{y_{t-1}^{(i)}=i} \mathbb{1}_{y_t^{(i)}=j} - \sum_{\mathbf{y}_{\langle 1 \dots T \rangle}} \frac{1}{Z(\mathbf{x}^{(i)})} \tilde{p}(\mathbf{y} \mid \mathbf{x}^{(i)}) \sum_{t=1}^T \mathbb{1}_{y_{t-1}=i} \mathbb{1}_{y_t=j} \quad (202)$$

$$= \sum_t \mathbb{1}_{y_{t-1}^{(i)}=i} \mathbb{1}_{y_t^{(i)}=j} - \sum_t \sum_{y_t} \sum_{y_{t-1}} \mathbb{1}_{y_{t-1}=i} \mathbb{1}_{y_t=j} \sum_{\mathbf{y}_{\langle 1 \dots t-2 \rangle}} \sum_{\mathbf{y}_{\langle t+1 \dots T \rangle}} \frac{1}{Z(\mathbf{x}^{(i)})} \tilde{p}(\mathbf{y} \mid \mathbf{x}^{(i)}) \quad (203)$$

## 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

$$p_{\theta, \gamma}(\mathbf{y}, \mathbf{z} \mid \mathbf{x}) = \frac{\exp(\boldsymbol{\theta}^T \Phi(\mathbf{x}, \mathbf{z}) + \boldsymbol{\gamma}^T \Psi(\mathbf{z}, \mathbf{y}))}{\sum_{\mathbf{z}', \mathbf{y}' \in \mathcal{Y}(\mathbf{x}, \mathbf{z}')} \exp(\boldsymbol{\theta}^T \Phi(\mathbf{x}, \mathbf{z}') + \boldsymbol{\gamma}^T \Psi(\mathbf{z}', \mathbf{y}'))} \quad (204)$$

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:

$$(\theta_1, \gamma_1) \approx \arg \max_{\theta, \gamma} \sum_{i=1}^M \log p_{\theta, \gamma}(C(u^{(i)}) \mid u^{(i)}) \quad (205)$$

2. Train a final model on the labeled data  $\{(x^{(i)}, y^{(i)})\}_{i=1}^N$ , using  $\theta_1$  as an initialization point:

$$(\theta_2, \gamma_2) \approx \arg \max_{\theta, \gamma} \sum_{i=1}^N \log p_{\theta, \gamma}(y^{(i)} \mid x^{(i)}) \quad (206)$$

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:

$$x \longrightarrow A^T x \quad (207)$$

$$y \longrightarrow B^T y \quad (208)$$

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:

$$D_{i,j} = \sum_{l=1}^n \mathbb{1}_{x_i^{(l)}=1} \mathbb{1}_{y_j^{(l)}=1} \quad (209)$$

$$u_i = \sum_{l=1}^n \mathbb{1}_{x_i^{(l)}=1} \quad (210)$$

$$v_i = \sum_{l=1}^n \mathbb{1}_{y_i^{(l)}=1} \quad (211)$$

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.

## 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

$$p(\mathbf{y} \mid \mathbf{x}) \propto \prod_{c \in C} e^{-E_c(\mathbf{x}, \mathbf{y}_c, \mathbf{w})} = \prod_{c \in C} e^{\langle \mathbf{w}_c^{\mathbf{y}_c}, \Phi_c(\mathbf{x}, \mathbf{w}_{NN}) \rangle} \quad (212)$$

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.

$$p(\mathbf{y} \mid \mathbf{x}) = \frac{1}{Z(\mathbf{x})} \prod_t^T \exp \{-E_t(\mathbf{x}, \mathbf{y}_t, \mathbf{y}_{t-1}, \mathbf{w})\} \quad (213)$$

$$= \frac{1}{Z(\mathbf{x})} \prod_t^T \exp \{-E_{loc}(\mathbf{x}, t, y_t, \mathbf{w}) - E_{trans}(\mathbf{x}, t, y_{t-1}, y_t, \mathbf{w})\} \quad (214)$$

$$(215)$$

where the authors are using a blanket  $\mathbf{w}$  to denote all model parameters<sup>18</sup>.

---

<sup>18</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.

$$\frac{\partial L(\mathbf{w})}{\partial \mathbf{w}} = \frac{1}{n} \sum_{i=1}^n \frac{\partial L_i(\mathbf{w})}{\partial \mathbf{w}} \quad (216)$$

$$\frac{\partial L_i(\mathbf{w})}{\partial \mathbf{w}} = \frac{\partial L_i(\mathbf{w})}{\partial \mathbf{E}(\mathbf{x}^{(i)})} \frac{\partial \mathbf{E}(\mathbf{x}^{(i)})}{\partial \mathbf{w}} \quad (217)$$

$$[\mathbf{E}(\mathbf{x}^{(i)})]_t = E_{loc}(\mathbf{x}, t, y_t, \mathbf{w}) + E_{trans}(\mathbf{x}, t, y_{t-1}, y_t, \mathbf{w}) \quad (218)$$

$$(219)$$

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

$$s(\mathbf{x}, \mathbf{y}, \tilde{\theta}) = \sum_t^T (A_{y_{t-1}, y_t} + [f_\theta]_{y_t, t}) \quad (220)$$

**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.

## 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  $\hat{y}_i = \text{sgn}(\mathbf{w}^T \mathbf{x}_i)$ ,

$$\mathbf{w} \leftarrow \mathbf{w} + \mathbf{x}_i(y_i - \hat{y}_i) \quad (221)$$

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>19</sup>,

$$\left\| \frac{1}{N} \sum_i^N \mathbf{x}_i y_i - \frac{1}{N} \sum_i^N \mathbf{x}_i \hat{y}_i \right\| \sim \mathcal{O}\left(\frac{1}{N}\right) \quad (222)$$

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>19</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  $\mathcal{X}$ ),

$$E(\mathbf{x}) = -\mathbf{w}^T \phi(\mathbf{x}) \quad (223)$$

$$\mathbf{w}_{t+1} = \mathbf{w}_t + \bar{\phi} - \phi(\mathbf{x}_t^*) \quad (224)$$

$$\text{where } \bar{\phi} = \mathbb{E}_{\mathbf{x}^{(i)} \sim p_{data}} [\phi(\mathbf{x}^{(i)})] \quad (225)$$

$$\text{and } \mathbf{x}_t^* = \arg \max_{\mathbf{x}} \mathbf{w}_t^T \phi(\mathbf{x}) \quad (226)$$

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

$$\mathbf{z}(\mathbf{x}, \mathbf{w}) = \arg \max_{\mathbf{z}'} \mathbf{w}^T \phi(\mathbf{x}, \mathbf{z}') \quad (227)$$

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

$$\bar{\phi} = \mathbb{E}_{\mathbf{x}^{(i)} \sim p_{data}} [\phi(\mathbf{x}^{(i)}, \mathbf{z}(\mathbf{x}^{(i)}, \mathbf{w}))] \quad (228)$$

$$\mathbf{x}_t^*, \mathbf{z}_t^* = \arg \max_{\mathbf{x}, \mathbf{z}} \mathbf{w}_t^T \phi(\mathbf{x}, \mathbf{z}) \quad (229)$$

**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

$$\mathbf{w}_{t+1} = \mathbf{w}_t + \frac{\eta}{|\mathcal{D}_t|} \sum_{(\mathbf{x}^{(i)}, \mathbf{y}^{(i)}) \in \mathcal{D}_t} (\phi(\mathbf{x}^{(i)}, \mathbf{y}^{(i)}, \mathbf{z}) - \phi(\mathbf{x}^{(i)}, \mathbf{y}^*, \mathbf{z}^*)) \quad (230)$$

# STRUCTURED PREDICTION

## CONTENTS

3.1	Joint Event Extraction via RNNs . . . . .	50
3.2	Event Extraction via Bidi-LSTM Tensor NNs . . . . .	52
3.3	Reasoning with Neural Tensor Networks . . . . .	54
3.4	Language to Logical Form with Neural Attention . . . . .	55
3.5	Seq2SQL: Generating Structured Queries from NL using RL . . . . .	57
3.6	SLING: A Framework for Frame Semantic Parsing . . . . .	60
3.7	DeepWalk: Online Learning of Social Representations . . . . .	62
3.8	Review of Relational Machine Learning for Knowledge Graphs . . . . .	64
3.9	Fast Top-K Search in Knowledge Graphs . . . . .	67
3.10	Dynamic Recurrent Acyclic Graphical Neural Networks (DRAGNN) . . . . .	69
3.10.1	More Detail: Arc-Standard Transition System . . . . .	72
3.11	Joint Extraction of Events and Entities within a Document Context . . . . .	73
3.12	Globally Normalized Transition-Based Neural Networks . . . . .	76
3.13	Structured Attention Networks . . . . .	79
3.14	Relation Extraction: A Survey . . . . .	81
3.15	Neural Relation Extraction with Selective Attention over Instances . . . . .	84

## 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  
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

$$x_i := [\text{GloVe}(w_i); \text{Embed}(\text{EntityType}(w_i)); \text{DepVec}(w_i)] \quad (231)$$

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:

$$R_i^{trig} := [h_i; L_i^{trg}; G_{i-1}^{trg}] \quad (232)$$

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:

$$P_{i;t}^{trg} := F_t^{trg}(R_i^{trg}) \quad (233)$$

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:
```

---

$$R_{ij}^{arg} := [h_i; h_{ij}; L_{ij}^{arg}; B_{ij}; G_{i-1}^{arg}[j]; G_{i-1}^{arg/trg}[j]] \quad (234)$$

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:

$$c_l(w_i) = \overrightarrow{LSTM}(c_l(w_{i-1}), e(w_{i-1})) \quad (235)$$

$$c_r(w_i) = \overleftarrow{LSTM}(c_r(w_{i+1}), e(w_{i+1})) \quad (236)$$

$$(237)$$

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  $x_i$ :

$$x_i := c_l(w_i) \oplus e(w_i) \oplus c_r(w_i) \oplus pi \oplus pe \oplus e(a) \oplus e(t) \quad (238)$$

**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>20</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 238. We split our sentence by trigger and candidate argument, then (confusingly) redefine our notation as

$$y_{1j}^{(1)} \leftarrow \begin{bmatrix} y_{1j}^{(1)} & y_{2j}^{(1)} \end{bmatrix} \quad (239)$$

$$y_{2j}^{(1)} \leftarrow \begin{bmatrix} y_{3j}^{(1)} & \dots & y_{7j}^{(1)} \end{bmatrix} \quad (240)$$

$$y_{3j}^{(1)} \leftarrow \begin{bmatrix} y_{8j}^{(1)} & \dots & y_{nj}^{(1)} \end{bmatrix} \quad (241)$$

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

$$y_{i,j}^{(2)} := \max(y_{i,j}^{(1)}) \quad 1 \leq i \leq 3, 1 \leq j \leq m \quad (242)$$

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,

$$O = W_2 y^{(2)} + b_2 \quad (243)$$

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>20</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:

$$g(e_1, R, e_2) = u_R^T \tanh \left( e_1^T W_R^{[1:k]} e_2 + V_R \begin{bmatrix} e_1 \\ e_2 \end{bmatrix} + b_R \right) \quad (244)$$

$$W_R^{[1:k]} \in \mathbb{R}^{d \times d \times k}$$

$$V_R \in \mathbb{R}^{k \times 2d}$$

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:

$$J(\Omega) = \sum_{i=1}^N \sum_{c=1}^C \max \left( 0, 1 - g(T^{(i)}) + g(T_c^{(i)}) \right) + \lambda \|\Omega\|_2^2 \quad (245)$$

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.

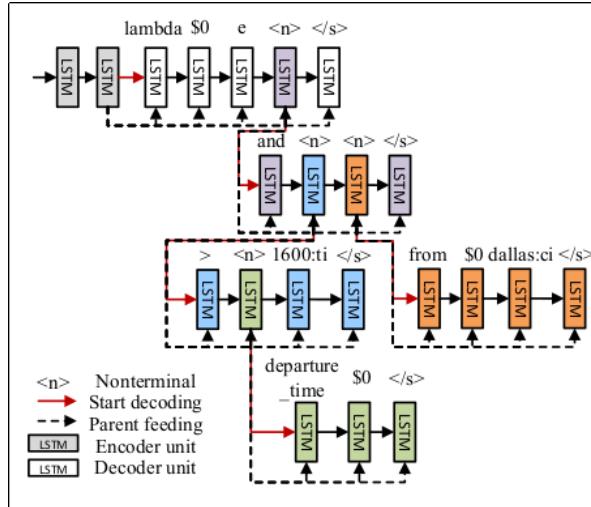
## Language to Logical Form with Neural Attention

Table of Contents Local

Written by Brandon McKinzie

Dong and Lapata, “Language to Logical Form with Neural Attention,” (2016)

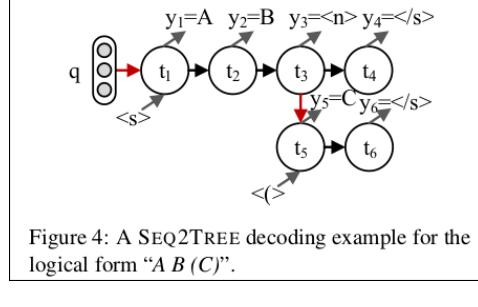
**Sequence-to-Tree Model.** Variant of Seq2Seq that is more faithful to the compositional nature of meaning representations. Its schematic is shown below. The authors define a “nonterminal”  $< n >$  token which indicates [the root of] a subtree.



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:

$$p(a | q) = p(y_1 y_2 y_3 y_4 | q) p(y_5 y_6 | y_{\leq 3}, q) \quad (246)$$

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.

$$x := [<\text{col}>; x_1^c; x_2^c; \dots; x_N^c; <\text{sql}>; x^s; <\text{question}>; x^q] \quad (247)$$

$$x_j^c \in \mathbb{R}^{T_j}$$

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).

$$[\text{ptr net}] y_s = \arg \max_t (\alpha_{s,t}^{ptr}) \quad \text{where} \quad \alpha_{s,t}^{ptr} = w^{ptr} \cdot \tanh(U^{ptr} g_s + V^{ptr} h_t) \quad (248)$$

**Seq2SQL.**

1. **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:

$$\alpha_t^{inp} = w^{inp} \cdot h_t^{enc} \quad (249)$$

$$\beta^{inp} = \text{softmax}(\alpha^{inp}) \quad (250)$$

$$\kappa^{agg} = \sum_t \beta_t^{inp} h_t^{enc} \quad (251)$$

$$\alpha^{agg} = W^{agg} \tanh(V^{agg} \kappa^{agg} + b^{agg}) + c^{agg} \quad (252)$$

$$\beta^{agg} = \text{softmax}(\alpha^{agg}) \quad (253)$$

$$W^{agg} \in \mathbb{R}^{C \times T}$$

$$\beta^{agg} \in \mathbb{R}^C$$

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>21</sup> as its representation;  $T_j$  denotes the number of words in the  $j$ th column.

$$e_j^c = h_{j,T_j}^c \quad \text{where} \quad h_{j,t}^c = \text{LSTM}\left(\text{emb}(x_{j,t}^c), h_{j,t-1}^c\right) \quad (254)$$

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:

$$\alpha_j^{sel} = W^{sel} \tanh\left(V^{sel} \kappa^{sel} + V^c e_j^c\right) \quad (255)$$

$$\beta^{sel} = \text{softmax}\left(\alpha^{sel}\right) \quad (256)$$

Similar to the aggregation, we train the SELECT network using cross entropy loss  $L^{sel}$ .

3. **WHERE Clause Pointer Decoder.** Recall from equation 248 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>22</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>21</sup>Yes, we encode each column with an LSTM separately.

<sup>22</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.

$$R(q(y), q_g) = \begin{cases} -2 & \text{if } q(y) \text{ is not a valid SQL query} \\ -1 & \text{if } q(y) \text{ is a valid SQL query and executes to an incorrect result} \\ +1 & \text{if } q(y) \text{ is a valid SQL query and executes to the correct result} \end{cases} \quad (257)$$

$$L^{whe} = -\mathbb{E}_y [R(q(y), q_g)] \quad (258)$$

$$\nabla L_\Theta^{whe} = -\nabla_\Theta (\mathbb{E}_{y \sim p_y} [R(q(y), q_g)]) \quad (259)$$

$$= -\mathbb{E}_{y \sim p_y} \left[ R(q(y), q_g) \nabla_\Theta \sum_t \log p_y(y_t; \Theta) \right] \quad (260)$$

$$\approx -R(q(y), q_g) \nabla_\Theta \sum_t \log p_y(y_t; \Theta) \quad (261)$$

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>23</sup>.*

---

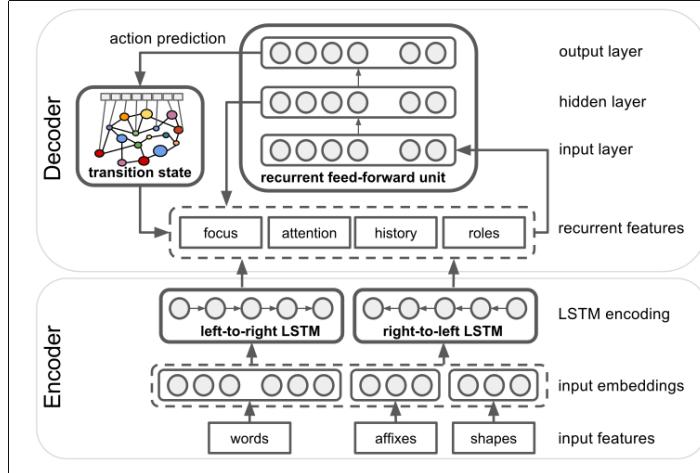
<sup>23</sup>Original: Seq2SQL takes as input a question and the columns of a table. It generates the corresponding SQL query, which, during training, is executed against a database. The result of the execution is utilized as the reward to train the reinforcement learning algorithm.

## 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.**

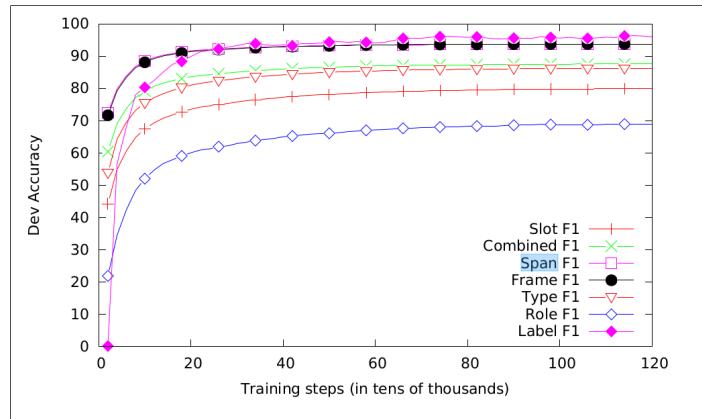
- **Inputs.** [words; affixes; shapes]
- **Encoder.**
  1. Embed.
  2. Bidirectional LSTM.
- **Inputs to TBRU.**
  - BLSTM [forward and backward] hidden state for the current token in the parser state.
  - **Focus.** Hidden layer activations corresponding to the transition steps that evoked/brought into **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>24</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>24</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>25</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>25</sup>This is important to keep in mind, since more than one sequence of actions can result in a given predicted frame graph.

## 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)

$$\Pr [v_i | \Phi(v_1), \Phi(v_2), \dots, \Phi(v_{i-1})] \tag{262}$$

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.

$$\min_{\Phi} -\log \Pr [v_{i-w}, \dots, v_{i-1}, v_{i+1}, \dots, v_{i+w} | \Phi(v_i)] \tag{263}$$

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  $\Phi \in \mathbb{R}^{|V| \times d}$

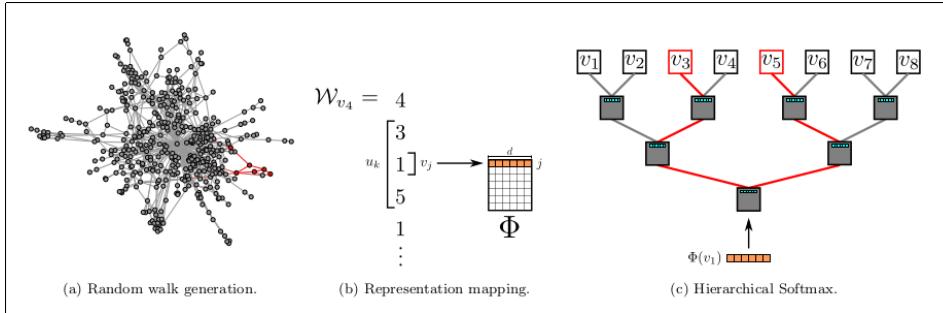
```

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

$$\Pr [\mathbf{Y} | \mathcal{D}, \Theta] = \prod_{i=1}^{N_e} \prod_{j=1}^{N_e} \prod_{k=1}^{N_r} \text{Ber}(y_{ijk} | \sigma(f(x_{ijk}; \Theta))) \quad (264)$$

where  $\text{Ber}$  is the Bernoulli distribution<sup>26</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>27</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

$$f_{ijk}^{\text{RESCAL}} := \mathbf{e}_i^T \mathbf{W}_k \mathbf{e}_j = \sum_{a=1}^{H_e} \sum_{b=1}^{H_e} w_{abk} e_{ia} e_{jb} \quad (266)$$

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

$$\mathbf{F}_k = \mathbf{E} \mathbf{W}_k \mathbf{E}^T \quad (267)$$

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

$$f_{ijk}^{\text{RESCAL}} := \mathbf{w}_k^T \phi_{i,j}^{\text{RESCAL}} \quad (268)$$

$$\phi_{i,j}^{\text{RESCAL}} := \mathbf{e}_j \otimes \mathbf{e}_i \quad (269)$$

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>26</sup>Notation used:

$$\text{Ber}(y | p) = \begin{cases} p & \text{if } y = 1 \\ 1 - p & \text{if } y = 0 \end{cases} \quad (265)$$

<sup>27</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:

$$f_{ijk}^{E-MLP} := \mathbf{w}_k^T \mathbf{g}(\mathbf{h}_{ijk}^a) \quad (270)$$

$$\mathbf{h}_{ijk}^a := \mathbf{A}_k^T \phi_{ij}^{E-MLP} \quad (271)$$

$$\phi_{ij}^{E-MLP} := [\mathbf{e}_i; \mathbf{e}_j] \quad (272)$$

**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  $e_i \xrightarrow{r_1} e_2 \xrightarrow{r_2} e_3 \dots \xrightarrow{r_L} e_j$ , 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:

$$\phi_{ijk}^{PRA} = [\Pr[\pi] : \pi \in \Pi_L(i, j, k)] \quad (273)$$

$$f_{ijk}^{PRA} := \mathbf{w}_k^T \phi_{ijk}^{PRA} \quad (274)$$

**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:

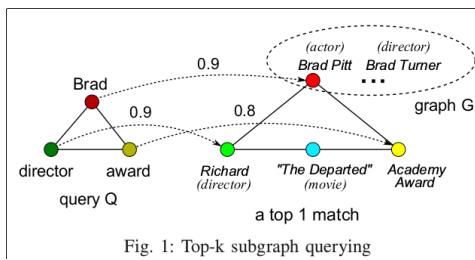


Fig. 1: Top-k subgraph querying

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

$$F(\phi(Q)) = \sum_{v \in V_Q} F_V(v, \phi(v)) + \sum_{e \in E_Q} F_E(e, \phi(e)) \quad (275)$$

$$F_V(v, \phi(v)) = \sum_i \alpha_i f_i(v, \phi(v)) \quad (276)$$

$$F_E(e, \phi(e)) = \sum_j \beta_j f_j(e, \phi(e)) \quad (277)$$

### 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>28</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>28</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>29</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> ) $([\sigma i j], B, A) \Rightarrow ([\sigma i], B, A \cup \{(i, l, j)\})$ (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:

$$d_i \leftarrow \arg \max_{d \in \mathcal{A}(s_i)} \mathbf{w}_d^T \mathbf{h}_i \quad (278)$$

$$s_{i+1} \leftarrow t(s_i, d_i) \quad (279)$$

**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>29</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:

$$L(\mathbf{x}, d_{N+1:N+n}^*; \theta) = \sum_i \log \Pr [d_{N+i}^* | d_{1:N}, d_{N+1:N+i-1}^*; \theta] \quad (280)$$

During training, the entire input sequence is unrolled and **backpropagation through structure** is used for gradient computation.

---

### 3.10.1 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.”

## 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

$$\Pr_{\boldsymbol{\theta}} [t_i, \mathbf{r}_i, \mathbf{a} | i, \mathcal{N}_i, x] \propto \quad (281)$$

$$\exp \left( \boldsymbol{\theta}_1^T \mathbf{f}_1(t_i) + \sum_{j \in \mathcal{N}_i} \left[ \boldsymbol{\theta}_2^T \mathbf{f}_2(r_{ij}) + \boldsymbol{\theta}_3^T f_3(t_i, r_{ij}) + \boldsymbol{\theta}_4^T \mathbf{f}_4(a_j) + \boldsymbol{\theta}_5^T f_5(r_{ij}, a_j) \right] \right) \quad (282)$$

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:

$$\boldsymbol{\theta}^* = \arg \max_{\boldsymbol{\theta}} \mathcal{L}(\boldsymbol{\theta}) - \lambda \|\boldsymbol{\theta}\|_2^2 \quad (283)$$

$$\mathcal{L}(\boldsymbol{\theta}) = \sum_i \log (\Pr_{\boldsymbol{\theta}} [t_i, \mathbf{r}_i, \mathbf{a} | i, \mathcal{N}_i, x]) \quad (284)$$

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

$$\Pr_\phi [t_i, t_{i'} | x, i, i'] \propto \exp (\phi^T g(t_i, t_{i'}, x, i, i')) \quad (285)$$

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:

$$\max_{\mathbf{t}, \mathbf{r}, \mathbf{a}} \sum_{i \in \mathcal{T}} E(t_i, \mathbf{r}_i, \mathbf{a}) + \sum_{i, i' \in \mathcal{T}} R(t_i, t_{i'}) + \sum_{j \in \mathcal{N}} D(a_j) \quad (286)$$

where

- The first term is the sum of confidence scores for individual event mentions from the within-event model.

$$E(t_i, \mathbf{r}_i, \mathbf{a}) = \log p_\theta(t_i) + \sum_{j \in \mathcal{N}_i} [\log p_\theta(r_{ij}) + \log p_\theta(r_{ij}, a_j)] \quad (287)$$

- The second term is the sum of confidence scores for event relations based on the pairwise event model.

$$R(t_i, t_{i'}) = \log p_\phi(t_i, t_{i'} | i, i', x) \quad (288)$$

- The third term is sum of confidence scores for entity mentions, where

$$D(a_j) = \log p_\psi(a_j | j, x) \quad (289)$$

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>30</sup>. To solve it efficiently, they find solutions for the relaxation of the problem using a **dual decomposition algorithm AD**<sup>3</sup>.

---

<sup>30</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:

$$\rho(s, d; \theta) = \phi(s; \theta^{(l)}) \cdot \theta^{(d)} \quad (290)$$

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  $s_j$ .

$$\Pr_L [d_{1:n}] = \prod_{j=1}^n \Pr [d_j | s_j; \theta] = \frac{\exp \left( \sum_{j=1}^n \rho(s_j, d_j; \theta) \right)}{\prod_{j=1}^n Z_L(s_j; \theta)} \quad (291)$$

$$Z_L(s_j; \theta) = \sum_{d' \in \mathcal{A}(s_j)} \exp (\rho(s_j, d'; \theta)) \quad (292)$$

Beam search can be used to attempt to find the action sequence with highest probability.

- **Global.** In contrast, a CRF defines:

$$\Pr_G [d_{1:n}] = \frac{\exp \left( \sum_{j=1}^n \rho(s_j, d_j; \theta) \right)}{Z_G(\theta)} \quad (293)$$

$$Z_G(\theta) = \sum_{d'_{1:n} \in \mathcal{D}_n} \exp \left( \sum_{j=1}^n \rho(s'_j, d'_j; \theta) \right) \quad (294)$$

where  $\mathcal{D}_n$  is the set of all valid sequences of decisions of length  $n$ . The inference problem is now to find

$$\arg \max_{d_{1:n} \in \mathcal{D}_n} \Pr_G [d_{1:n}] = \arg \max_{d_{1:n} \in \mathcal{D}_n} \sum_{j=1}^n \rho(s_j, d_j; \theta) \quad (295)$$

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.**

$$L_{local}(d_{1:n}^*; \theta) = -\ln \Pr_L [d_{1:n}^*; \theta] \quad (296)$$

$$= -\sum_{j=1}^n \rho(s_j^*, d_j^*; \theta) + \sum_{j=1}^n \ln Z_L(s_j^*; \theta) \quad (297)$$

- **Global.**

$$L_{global}(d_{1:n}^*; \theta) = -\ln \Pr_G [d_{1:n}^*; \theta] \quad (298)$$

$$= -\sum_{j=1}^n \rho(s_j^*, d_j^*; \theta) + \ln Z_G(\theta) \quad (299)$$

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>31</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:

$$L_{global-beam}(d_{1:j}^*, \theta) = -\sum_{t=1}^j \rho(d_{1:t-1}^*, d_t^*; \theta) - \ln \sum_{d'_{1:j} \in \mathcal{B}_j} \exp \left( \sum_{t=1}^j \rho(d'_{1:t-1}, d'_t; \theta) \right) \quad (300)$$

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>32</sup> using the local objective function, and then perform additional training steps using the global objective function..

---

<sup>31</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>33</sup>. Let  $\mathcal{P}_L$  denote the set of all possible distributions  $p_L(d_{1:n} \mid x_{1:n})$  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>34</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

$$p_G(d_{1:t-1}, d_t, x_{1:t}) = \ln p_L(d_t \mid d_{1:t-1}, x_{1:t}) \quad (301)$$

By definition, this means that  $p_G(d_{1:t} \mid x_{1:t}) = p_L(d_{1:t} \mid x_{1:t})$ .

**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:

$$\mathbf{x} = abc, \quad \mathbf{d} = ABC \quad (302)$$

$$\mathbf{x} = abe, \quad \mathbf{d} = ADE \quad (303)$$

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

$$\rho(d_{1:t-1}, d_t, x_{1:t}) = \alpha (\mathbb{1}_{(d_{t-1}, d_t) \in \mathcal{T}} + \mathbb{1}_{(x_t, d_t) \in \mathcal{E}}) \quad (304)$$

for all  $t$ . This results in the following distributions  $p_G$  and  $p_L$ , evaluating on input sequence of length 3

$$p_G(d_{1:3} \mid x_{1:3}) = \frac{\exp(\alpha \sum_{t=1}^3 (\mathbb{1}_{(d_{t-1}, d_t) \in \mathcal{T}} + \mathbb{1}_{(x_t, d_t) \in \mathcal{E}}))}{Z_G(x_{1:3})} \quad (305)$$

$$p_L(d_{1:3} \mid x_{1:3}) = p_L(d_1 \mid x_1)p_L(d_2 \mid d_1, x_{1:2})p_L(d_3 \mid d_{1:2}, x_{1:3}) \quad (306)$$

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  $\sum_{d_t} p_L(d_t \mid d_{1:t-1}, x_{1:t}) = 1$* . With this, the following comparison of  $p_G$  and  $p_L$  for large  $\alpha$  completes the proof of  $\mathcal{P}_G \not\subseteq \mathcal{P}_L$ :

$$\lim_{\alpha \rightarrow \infty} p_G(ABC \mid abc) = \lim_{\alpha \rightarrow \infty} p_G(ADE \mid abe) = 1 \quad (307)$$

$$p_L(ABC \mid abc) + p_L(ADE \mid abe) \leq 1 \quad (\forall \alpha) \quad (308)$$

$$\therefore \mathcal{P}_L \subsetneq \mathcal{P}_G.$$

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>33</sup>This is for conditional models only.

<sup>34</sup>Note that  $\subset$  and  $\subsetneq$  mean the same thing. Matter of notational preference/being explicit/etc.

## 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>35</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>36</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:

$$c = \sum_t^T p(z = t | x, q) \mathbf{x}_t \quad (309)$$

$$p(z = t | x, q) = \text{softmax}(\theta_t) \quad (310)$$

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:

$$c = \mathbb{E}_{z \sim p(z|x,q)} [f(x, z)] = \sum_C \mathbb{E}_{z \sim p(z_C|x,q)} [f_C(x, z_C)] \quad (311)$$

$$p(z | x, q) = \frac{1}{Z(x, q)} \prod_C \psi_C(z_C) \quad (312)$$

---

<sup>35</sup>Also called the “alignments”. It is the output of the softmax layer of attention scores in the majority of cases.

<sup>36</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>37</sup>. This yields the context vector,

$$c = \mathbb{E}_{z_1, \dots, z_T} [f(x, z)] = \sum_t^T p(z_t = 1 | x, q) \mathbf{x}_t \quad (313)$$

Although this looks similar to equation 309, we haven't yet revealed the functional form for  $p(z | x, q)$ . Two possible choices:

**Linear-Chain CRF:**  $p(z_1, \dots, z_T | x, q) = \frac{1}{Z(x, q)} \prod_t^T \psi_t(z_t, z_{t-1}) \quad (314)$

**Bernoulli:**  $p(z_1, \dots, z_T | x, q) = \prod_t^T p(z_t = 1 | x, q) = \prod_t^T \sigma(\psi_t(z_t)) \quad (315)$

The factor  $\psi$  for the CRF is NOT the same as the factor for the Bernoulli!

These show why equation 313 is fundamentally different than equation 309:

- 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>38</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>37</sup>Ok, so equivalently,  $z^T x$ , i.e. the indicator function can just be replaced by  $z_t$  here

<sup>38</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.

## 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>39</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	M1_leaders, M2_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>40</sup> (a relation instance is a triplet of the form  $(e_1, e_2)$ ), and SVM for the classification.

---

<sup>39</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>40</sup>Recall that kernel methods are for the general optimization problem

$$\min_w \sum_{i=1}^n \text{loss} \left( \sum_{j=1}^n \alpha_j \phi(x_j)^T \phi(x_i), y_i \right) + \lambda \sum_{j=1}^n \sum_{k=1}^n \alpha_j \alpha_k \phi(x_j)^T \phi(x_k) \quad (316)$$

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

$$u_i \in \mathbf{s}_i \quad \forall i \in ii \tag{317}$$

$$u_j \in \mathbf{t}_j \quad \forall j \in jj \tag{318}$$

$$\lambda^{l(ii)+l(jj)}, \text{ where } l(x) = \max(x) - \min(x) \text{ and } 0 < \lambda < 1. \text{ Sparser (more spaced out) subsequences get lower weight.}$$

$$\tag{319}$$

- 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

$$s(t_i, t_j) = \sum_{v_q \in \phi(t_i)} \sum_{v_r \in \phi(t_j)} \text{Compat}(v_q, v_r) \tag{320}$$

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  $T_i$ .

$$K(T_1, T_2) = \begin{cases} 0 & \text{if } m(r_1, r_2) = 0 \\ s(r_1, r_2) + K_c(r_1[\mathbf{c}], r_2[\mathbf{c}]) & \text{otherwise} \end{cases} \quad (321)$$

$$K_c(t_i[\mathbf{c}], t_j[\mathbf{c}]) = \sum_{\mathbf{a}, \mathbf{b}, l(\mathbf{a})=l(\mathbf{b})} \lambda^{d(\mathbf{a})+d(\mathbf{b})} K(t_i[\mathbf{a}], t_j[\mathbf{b}]) \quad (322)$$

$$a_1 \leq a_2 \leq \dots \leq a_n$$

$$d(\mathbf{a}) \triangleq a_n - a_1 + 1$$

$$0 < \lambda < 1$$

The interpretation is that, whenever a pair of matching nodes is found, *all* possible matching subsequences<sup>41</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>41</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  $\mathbf{w}_i \in \mathbb{R}^d$  ( $d = d_w + 2d_{pos}$ ),

$$\mathbf{w}_i := [\text{word2Vec}(w_i); \text{dist}(w_i, e_1); \text{dist}(w_i, e_2)] \quad (323)$$

where  $\text{dist}(a, b)$  returns the [embedded] relative distance (num tokens) between a and b in the given sentence (positive integer)<sup>42</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,

$$\mathbf{q}_i = \mathbf{w}_{i-\ell+1:i} \quad (1 \leq i \leq T + \ell - 1) \quad (324)$$

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

$$\mathbf{p}_k = [\mathbf{W}\mathbf{Q}^T]_k + \mathbf{b} \quad (325)$$

$$[\mathbf{x}]_k = [\max(\mathbf{p}_{k1}); \max(\mathbf{p}_{k2}); \max(\mathbf{p}_{k3})] \quad (326)$$

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>42</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

$$\mathbf{s} = \sum_i \alpha_i \mathbf{x}_i \quad (327)$$

$$\alpha_i = \text{softmax}(\mathbf{x}_i \mathbf{A} \mathbf{r}) \quad (328)$$

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

$$p(r | \mathcal{S}, (e_1, e_2)\theta) = \text{softmax}(\mathbf{M}\mathbf{s} + \mathbf{d})_r \quad (329)$$

where  $\mathcal{S}$  is the set of  $n$  sentences for the given entity pair  $(e_1, e_2)$ .

# LANGUAGE MODELING

## CONTENTS

4.1	Improving Language Understanding by Generative Pre-Training . . . . .	87
4.2	Exploring the Limits of Language Modeling . . . . .	88
4.3	Noise Contrastive Estimation . . . . .	90
4.3.1	Self-Normalized NCE . . . . .	92
4.4	Improving Neural Language Models with a Continuous Cache . . . . .	94
4.5	Context Dependent RNN Language Model . . . . .	95
4.6	Strategies for Training Large Vocabulary Neural Language Models . . . . .	96
4.7	Large Memory Layers with Product Keys . . . . .	98
4.8	XLNet . . . . .	100
4.9	Transformer-XL . . . . .	102
4.10	Efficient Softmax Approximation for GPUs . . . . .	103
4.11	Adaptive Input Representations for Neural Language Modeling . . . . .	104
4.12	Fine-Tuning Language Models from Human Preferences . . . . .	105
4.13	Cross-lingual Language Model Pretraining . . . . .	106
4.14	Learning and Evaluating General Linguistic Intelligence . . . . .	108
4.15	Scaling Laws for Neural Language Models . . . . .	111
4.16	Language Models are Few-Shot Learners . . . . .	113
4.17	Generating Long Sequences with Sparse Transformers . . . . .	118
4.18	The Curious Case of Neural Text Degeneration . . . . .	120

## 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:

$$L_1(\mathcal{U}) = \sum_i^n \log P(u_i | u_{i-k}, \dots, u_{i-1}; \Theta) \quad (330)$$

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:

$$P(y | x^1, \dots, x^m) = \text{softmax}(h_l^m W_y) \quad (331)$$

$$L_2(\mathcal{C}) = \sum_{(x,y)} \log P(y | x^1, \dots, x^m) \quad (332)$$

They also found that including a language modeling auxiliary objective helped learning,

$$L_3(\mathcal{C}) = L_2(\mathcal{C}) + \lambda L_1(\mathcal{C}) \quad (333)$$

...that’s it. Extremely simple, yet somehow effective.

## 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

$$p(Y=\text{true} | w) = \frac{p(w | \text{true})p(\text{true})}{p(w)} \quad (334)$$

$$= \frac{p_d(w)p(\text{true})}{p(w, \text{true}) + p(w, \text{false})} \quad (335)$$

$$= \frac{p_d(w)p(\text{true})}{p_d(w)p(\text{true}) + p_n(w)p(\text{false})} \quad (336)$$

$$= \frac{p_d(w)}{p_d(w) + kp_n(w)} \quad (337)$$

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,

$$p(Y=i | W) = \frac{p_d(w_i)}{p_n(w_i)} \frac{1}{\sum_{i'=1}^{k+1} p_d(w_{i'})/p_n(w_{i'})} \quad (338)$$

$$\propto_Y \frac{p_d(w_i)}{p_n(w_i)} \quad (339)$$

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

$$z_w = h^T CNN(chars_w) + h^T M corr_w \quad (340)$$

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.

## 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  $\alpha^{43}$ .

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  $J_T(\theta)^{44}$ :

$$J_T(\theta) = \frac{1}{2T} \sum_t \ln [h(\mathbf{x}_t)] + \ln [1 - h(\mathbf{y}_t)] \quad (341)$$

$$h(\mathbf{u}) = \sigma(G(\mathbf{u})) \quad (342)$$

$$G(\mathbf{u}) = \ln p_m(\mathbf{u}) - \ln p_n(\mathbf{u}) \quad (343)$$

$$\ln p_m(\cdot; \theta) := \ln \tilde{p}_m(\cdot; \alpha) + c \quad (344)$$

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>45</sup>.

<sup>43</sup>The implicit assumption here is that  $\exists \alpha^*$  such that  $p_d(\cdot) = p_m(\cdot; \alpha^*)$ .

<sup>44</sup>This all assumes of course that  $p_n$  is fully defined.

<sup>45</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 ( $p_n$ ).

$$P(C=1 \mid \mathbf{u} \in X \cup Y) = \frac{P(C=1)p(\mathbf{u} \mid C=1)}{p(\mathbf{u})} \quad (345)$$

$$= \frac{p_m(\mathbf{u})}{p_m(\mathbf{u}) + p_n(\mathbf{u})} \quad (346)$$

$$\equiv h(\mathbf{u}; \theta) \quad (347)$$

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

$$\ell(\theta) = \sum_t^{2T} [C_t \ln P(C_t=1 \mid \mathbf{u}_t; \theta) + (1 - C_t) \ln P(C_t=0 \mid \mathbf{u}_t; \theta)] \quad (348)$$

$$= \sum_t^{2T} [C_t \ln h(\mathbf{u}_t; \theta) + (1 - C_t) \ln [1 - h(\mathbf{u}_t; \theta)]] \quad (349)$$

$$= \sum_t^T [\ln h(\mathbf{x}_t; \theta) + \ln [1 - h(\mathbf{y}_t; \theta)]] \quad (350)$$

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

$$J(\theta) \triangleq \lim_{T \rightarrow \infty} J_T(\theta) = \frac{1}{2} \mathbb{E} [\ln h(\mathbf{x}; \theta) + \ln [1 - h(\mathbf{y}; \theta)]] \quad (351)$$

$$\tilde{J}(f) \triangleq \frac{1}{2} \mathbb{E} \left[ \ln \left[ \sigma \left( f(\mathbf{x}) - \ln p_n(\mathbf{x}) \right) \right] + \ln \left[ 1 - \sigma \left( f(\mathbf{y}) - \ln p_n(\mathbf{y}) \right) \right] \right] \quad (352)$$

### Theorem 1

$\tilde{J}$  attains exactly one maximum, located at  $f(\cdot) = \ln p_d(\cdot)$ , provided  $p_d(\cdot) \neq 0 \implies p_n(\cdot) \neq 0$ .

---

### 4.3.1 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.

$$\frac{\partial \log P(w)}{\partial \theta} = \frac{\partial}{\partial \theta} s(w, h) - \frac{\partial}{\partial \theta} \log \left[ \sum_{w'} e^{s(w', h)} \right] \quad (353)$$

$$= \frac{\partial}{\partial \theta} s(w, h) - \sum_{w'} P(w') \frac{\partial}{\partial \theta} s(w', h) \quad (354)$$

$$= \frac{\partial}{\partial \theta} s(w, h) - \mathbb{E}_{w \sim P_\theta^h} \left[ \frac{\partial}{\partial \theta} s(w, h) \right] \quad (355)$$

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:

$$\frac{\partial \log P(w)}{\partial \theta} = \frac{\partial}{\partial \theta} s(w, h) - \sum_{w'} P(w') \frac{\partial}{\partial \theta} s(w', h) \quad (356)$$

$$\approx \frac{\partial}{\partial \theta} s(w, h) - \frac{1}{V} \sum_{j=1}^k v(x_j) \frac{\partial}{\partial \theta} s(w', h) \quad (357)$$

$$\text{where } v(x) = \frac{e^{s_\theta(x, h)}}{Q^h(w=x)} \quad (358)$$

and we refer to  $v$  as the **importance weights**. In NLP, we typically set  $Q$  to the **Zipfian distribution**<sup>46</sup>

---

<sup>46</sup>TensorFlow seems to define this as

$$P_{Zipf}(w) = \frac{\log(w+2) - \log(w+1)}{\log(V+1)} \quad (359)$$

where  $V$  is the vocabulary size. I can’t seem to find this definition anywhere else though. A more common form seems to be

$$P(w) = \frac{\frac{1}{w}}{\sum_{w'} \frac{1}{w'^s}} \quad (360)$$

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,

$$P^h(D, w) = P(D=1)P_d^h(w) + P(D=0)P_n(w) \quad (361)$$

$$= \frac{1}{k+1}P_d^h(w) + \frac{k}{k+1}P_n(w) \quad (362)$$

Our goal is to learn the posterior distribution  $P^h(D=1 | w)$  (so we replace  $P_d$  with  $P_\theta$ ):

$$P^h(D=1 | w, \theta) = \frac{P_\theta^h(w)}{P_\theta^h(w) + kP_n(w)} \quad (363) \quad \text{NCE posterior}$$

In NCE, we re-parameterize  $P_\theta$  by treating  $-\log Z$  as a parameter itself,  $c^h$ <sup>47</sup>.

$$P_\theta^h(w) := P_{\theta^0}^h \exp(c^h) \quad (364)$$

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 363 using the unnormalized probabilities in place of  $P_\theta^h$ .* The NCE objective is to find<sup>48</sup> as follows, where I've shown each step of the derivation:

$$\theta^* = \arg \max_{\theta} J^h(\theta) \quad (365)$$

$$J^h(\theta) = \mathbb{E}_{(D, w) \sim P^h} [\log P(D | w, \theta)] \quad (366)$$

$$= \sum_{D=0}^1 \sum_w P^h(D, w) \log P^h(D | w, \theta) \quad (367)$$

$$= \sum_w P^h(0, w) \log P^h(0 | w) + \sum_w P^h(1, w) \log P^h(1 | w) \quad (368)$$

$$= \frac{1}{k+1} \sum_w [kP_n(w) \log P^h(0 | w) + P_d^h(w) \log P^h(1 | w)] \quad (369)$$

$$= \frac{1}{k+1} \left[ k\mathbb{E}_{P_n} [\log P^h(0 | w)] + \mathbb{E}_{P_d^h} [\log P^h(1 | w)] \right] \quad (370)$$

$$\propto k\mathbb{E}_{P_n} [\log P^h(0 | w)] + \mathbb{E}_{P_d^h} [\log P^h(1 | w)] \quad (371)$$

The gradient of the NCE objective is thus

$$\frac{\partial}{\partial \theta} J^h(\theta) = \sum_w \frac{kP_n(w)}{P_\theta^h(w) + kP_n(w)} (P_d^h(w) - P_\theta^h(w)) \frac{\partial}{\partial \theta} \log P_\theta^h(w) \quad (372)$$

**TODO:** incorporate more info from Chris Dyer's excellent notes.

---

<sup>47</sup>Reminder that the  $h$  is a reminder that  $Z$  is a function of the context  $h$ .

<sup>48</sup>I'll drop off  $\theta$  dependence wherever obvious for the sake of compactness.

## 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>49</sup> based on this representation.

$$p_{vocab}(w | \mathbf{x}_{\langle 1 \dots t \rangle}) \propto \exp(h_t^T o_w) \quad (373)$$

$$p_{cache}(w | \mathbf{h}_{\langle 1 \dots t \rangle}, \mathbf{x}_{\langle 1 \dots t \rangle}) \propto \sum_{i=1}^{t-1} \mathbb{1}_{x_{i+1}=w} \exp(\theta h_i^T h_i) \quad (374)$$

$$= \sum_{\substack{(x,h) \in \text{cache} \\ s.t. x=w}} \exp(\theta h_i^T h_i) \quad (375)$$

$$p(w | \mathbf{h}_{\langle 1 \dots t \rangle}, \mathbf{x}_{\langle 1 \dots t \rangle}) = (1 - \lambda)p_{vocab}(w | h_t) + \lambda p_{cache}(w | \mathbf{h}_{\langle 1 \dots t \rangle}, \mathbf{x}_{\langle 1 \dots t \rangle}) \quad (376)$$

where  $\theta$  is a scalar parameter that controls the flatness of the cache distribution.

---

<sup>49</sup>They say this, but everything else in the paper strongly suggests they mean the next gold-standard input instead.

## 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.

$$\mathbf{y}_t = \text{Softmax}(\mathbf{V}\mathbf{h}_t + \mathbf{G}\mathbf{f}_t) \quad (377)$$

$$\mathbf{h}_t = \sigma(\mathbf{U}\mathbf{x}_t + \mathbf{W}\mathbf{h}_{t-1} + \mathbf{F}\mathbf{f}_t) \quad (378)$$

**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>50</sup>:

$$N \sim \text{Poisson}(\xi) \quad (379) \quad \begin{matrix} \text{N: number of words} \\ \Theta_i \equiv p(\text{topic}[i]) \end{matrix}$$

$$\Theta \sim \text{Dir}(\alpha) \quad (380) \quad z_n: \text{topic of word } n$$

$$z_n \sim \text{Multinomial}(\Theta) \quad (381)$$

$$w_n \sim \text{Pr}[w_n | z_n, \beta] \quad (382)$$

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  $\beta$ .

$$\mathbf{f}_t = \frac{1}{Z} \prod_{i=0}^K \mathbf{t}_{x_{t-i}} \quad (383)$$

$$\mathbf{f}_t = \frac{1}{Z} \mathbf{f}_{t-1}^\gamma \mathbf{t}_{x_t}^{1-\gamma} \quad (384)$$

---

<sup>50</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:

$$h^{k+1} = W^{k+1}h^k + b^{k+1} \in \mathbb{R}^V \quad (385)$$

$$y = \frac{1}{Z} \exp \left\{ h^{k+1} \right\} \quad (386)$$

Using cross-entropy loss, the derivative of  $\log p(w_t=i)$  wrt the  $j$ th element of the logits is:

$$\frac{\partial \log y_i}{\partial h_j^{k+1}} = \frac{\partial}{\partial h_j^{k+1}} \left[ h_i^{k+1} - \log Z \right] \quad (387)$$

$$= \delta_{ij} - y_j \quad (388)$$

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>51</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>51</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

$$L_{NCE}(\{\mathbf{w}_1, \dots, \mathbf{w}_N\}) = \sum_{i=1}^N \sum_{t=1}^{len(\mathbf{w}^{(i)})} \left[ \log h(w_t^{(i)}) + \sum_{k=1}^K \log(1 - h(n_k^{(i)})) \right] \quad (389)$$

$$h(w_t) = \frac{P_{model}(w)}{P_{model}(w) + kP_{noise}(w)} \quad (390)$$

$$\approx \frac{\tilde{P}_{model}(w)}{\tilde{P}_{model}(w) + kP_{noise}(w)} \quad (391)$$

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.

## 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>52</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>53</sup>.

The *standard* way of doing key assignment/weighting is as follows:

$$[\text{KNN}] \quad \mathcal{I} \triangleq \text{TopK}(\mathbf{q}(\mathbf{x})^T \mathbf{k}_i) \quad 1 \leq i \leq \mathcal{K} \quad (392)$$

$$[\text{Normalize}] \quad \mathbf{w} = \text{Softmax}(\mathcal{I}) \quad (393)$$

$$[\text{Aggregate}] \quad \mathbf{m}(\mathbf{x}) = \sum_{i \in \mathcal{I}} w_i \mathbf{v}_i \quad (394)$$

where equation 392 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

$$\mathcal{K} \triangleq \{(\mathbf{c}, \mathbf{c}') \mid c \in \mathcal{C}, c' \in \mathcal{C}'\} \quad (395)$$

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>52</sup>They choose  $d_q = 512$  as the output dimensionality of their query network.

<sup>53</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

## 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

$$[\text{BERT}] \quad \max_{\theta} \log p_{\theta}(\bar{\mathbf{x}} | \hat{\mathbf{x}}) \approx \sum_{\bar{x} \in \bar{\mathbf{x}}} \log p_{\theta}(\bar{x} | \hat{\mathbf{x}}) \quad (396)$$

$$p(\bar{x} | \hat{\mathbf{x}}) = \text{Softmax} \left( H_{\theta}(\hat{\mathbf{x}})^T e(\bar{x}) \right) \quad (397)$$

**Objective & Architecture.** Their proposed **permutation language modeling objective** is:

$$\max_{\theta} \mathbb{E}_{\mathbf{z} \sim \mathcal{Z}_T} \left[ \sum_{t=1}^T \log p_{\theta}(x_{z_t} | \mathbf{x}_{\langle z_1 \dots z_{t-1} \rangle}) \right] \quad (398)$$

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**:

$$p_{\theta}(X_{z_t} = x | \mathbf{x}_{\langle z_1 \dots z_{t-1} \rangle}) = \text{Softmax} \left( g_{\theta} \left( \mathbf{x}_{\langle z_1 \dots z_{t-1} \rangle}, \mathbf{z}_t \right)^T e(x) \right) \quad (399)$$

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:

$$g_{z_t}^{(m)} \leftarrow \text{Attention}(Q=g_{z_t}^{(m-1)}, K=V=\mathbf{h}_{z_{<t}}^{(m-1)}) \quad (400)$$

$$h_{z_t}^{(m)} \leftarrow \text{Attention}(Q=h_{z_t}^{(m-1)}, K=V=\mathbf{h}_{z_{\leq t}}^{(m-1)}) \quad (401)$$

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:

$$h_{z_t}^{(m)} \leftarrow \text{Attention}\left(Q=h_{z_t}^{(m-1)}, K=V=\left[\tilde{\mathbf{h}}^{(m-1)}; \mathbf{h}_{z_{\leq t}}^{(m-1)}\right]\right) \quad (402)$$

## 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:

$$\mathbf{h}_{\tau+1}^n = \text{TransformerLayer}(\mathbf{q}_{\tau+1}^n, \mathbf{k}_{\tau+1}^n, \mathbf{v}_{\tau+1}^n) \quad (403)$$

$$= \text{TransformerLayer}(\mathbf{h}_{\tau+1}^{n-1} \mathbf{W}_q^T, \tilde{\mathbf{h}}_{\tau+1}^{n-1} \mathbf{W}_k^T, \tilde{\mathbf{h}}_{\tau+1}^{n-1} \mathbf{W}_v^T) \quad (404)$$

$$\tilde{\mathbf{h}}_{\tau+1}^{n-1} = [\text{SG}(\mathbf{h}_\tau^{n-1}); \mathbf{h}_{\tau+1}^{n-1}] \quad (405)$$

where the concat in 405 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.

$$A_{i,j}^{abs} = e_{x_i}^T W_q^T W_k e_{x_j} + e_{x_i}^T W_q^T W_k u_j + u_i^T W_q^T W_k e_{x_j} + u_i^T W_q^T W_k u_j \quad (406)$$

$$A_{i,j}^{res} = \underbrace{e_{x_i}^T W_q^T W_k, E e_{x_j}}_{\text{cont-based addr}} + \underbrace{e_{x_i}^T W_q^T W_k, R r_{i-j}}_{\text{cont-dep pos bias}} + \underbrace{\mathbf{u}^T W_k, E e_{x_j}}_{\text{global cont bias}} + \underbrace{\mathbf{v}^T W_k, R r_{i-j}}_{\text{global pos bias}} \quad (407)$$

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:

$$\Pr[w | h] = \begin{cases} P_{\mathcal{V}_h}(w | h) & \text{if } w \in \mathcal{V}_h \\ P_{\mathcal{V}_t}(w | h)P_{\mathcal{V}_h}(\text{tail} | h) & \text{otherwise} \end{cases} \quad (408)$$

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

$$\mathcal{V} = \mathcal{V}_1 \cup \mathcal{V}_2 \cup \dots \cup \mathcal{V}_n \quad (409)$$

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  $\mathcal{V}_1$ ).

## Fine-Tuning Language Models from Human Preferences

Table of Contents Local

Written by Brandon McKinzie

Ziegler et al., “Fine-Tuning Language Models from Human Preferences” *OpenAI*, (Sep 2019).

**Methods** (2). Setup: vocab  $\Sigma$ , language model  $\rho$ . Want to model probability of response sequence  $y$  given input sequence  $x$ . They initialize a **policy**  $\pi = \rho$ , then fine-tune  $\pi$  with RL. Instead of *defining* a reward function  $r : X \times Y \mapsto \mathbb{R}$ , the *learn* one using a dataset  $S$  of human annotations. Each element of  $S$  is a tuple  $(x, y_0, y_1, y_2, y_3, b)$  – each  $y_i$  is a multiple-choice option presented to the human labeler, and the human selects option  $0 \leq b < 3$ . They train the reward model  $r$  with loss

$$\text{loss}(r) = \mathbb{E}_{x, \{y_i\}, b \sim S} \left[ \log \frac{e^{r(x, y_b)}}{\sum_i e^{r(x, y_i)}} \right] \quad (410)$$

They initialize  $r$  as a random linear function of the final embedding output of  $\rho$ . They fine-tune  $\pi$  to optimize  $r$ , performing RL on the modified reward:

$$R(x, y) = r(x, y) - \beta \log \frac{\pi(y | x)}{\rho(y | x)} \quad (411)$$

how is this any different from just continuing training the LM on the human annotations?

## Cross-lingual Language Model Pretraining

Table of Contents Local

Written by Brandon McKinzie

G. Lample and A. Conneau, “Cross-lingual Language Model Pretraining” *FAIR*, (Jan 2019).

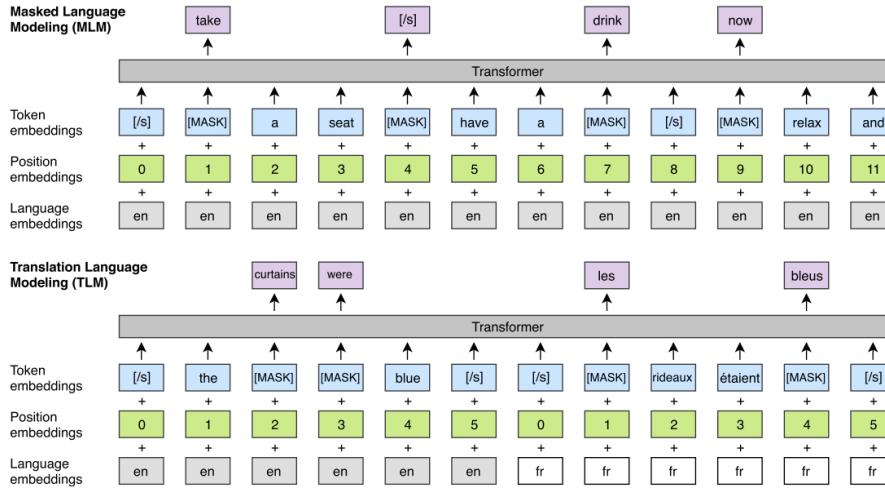
**Cross-lingual Language Models** (3). Use single shared vocabulary across all languages using BPE. Learn BPE splits by sampling sentences randomly from each corpus<sup>54</sup>  $C_i$  containing  $n_i$  sentences, according to a multinomial distribution with probabilities  $\{q_i\}_{i=1}^N$ :

$$q_i = \frac{p_i^\alpha}{\sum_{j=1}^N p_j^\alpha} \quad \text{with} \quad p_i = \frac{n_i}{\sum_{k=1}^N n_k} \quad (412) \quad \text{Authors set } \alpha = 0.7$$

*Sampling with this distribution increases the number of tokens associated to low-resource languages and alleviates the bias towards high-resource languages. In particular, this prevents words of low-resource languages from being split at the character level.*

They experiment with the following language modeling objectives:

- **CLM**: causal language modeling. Predict next word given context.
- **MLM**: masked language modeling. Predict masked-out tokens.
- **TLM**: translation language modeling. Extension of MLM when parallel data available. Concatenate parallel sentences and predict masked-out tokens.



<sup>54</sup>Each corpus is associated with a unique language. There are  $N$  languages (and thus corpora) total.

**Cross-lingual Language Model Pretraining** (4). After LM pretraining, add linear classifier on top of first-state final-layer of transformer and finetune all params on *English* NLI training data. They show how the model can still be used for non-English NLI queries at evaluation. This is what they are referring to by “zero-shot” classification – the fact that they can classify on non-English even though they only finetuned on English.

**Results & Analysis** (5). Simply training on multiple language can improve perplexity on low-resource languages. The table below shows how perplexity on Nepali improves by including various other languages during training.

Training languages	Nepali perplexity
Nepali	157.2
Nepali + English	140.1
Nepali + Hindi	115.6
Nepali + English + Hindi	<b>109.3</b>

Table 4: **Results on language modeling.** Nepali perplexity when using additional data from a similar language (Hindi) or a distant one (English).

## Learning and Evaluating General Linguistic Intelligence

Table of Contents Local

Written by Brandon McKinzie

Yogatama et al., “Learning and Evaluating General Linguistic Intelligence” *DeepMind*, (Jan 2019).

**Introduction** (1). Authors argue for an evaluation paradigm that rewards abilities to

1. Deal with the full complexity of natural language across a variety of tasks.
2. Effectively store/reuse representations, combinatorial modules, and previously acquired knowledge to avoid **catastrophic forgetting**.
3. Adapt to new linguistic tasks/environments with little experience.

**Tasks** (2). Below are the tasks & their respective datasets we’ll explore.

- **Reading Comprehension.** SQuAD 1.1, TriviaQA, QuAC.
- **Semantic Role Labeling.** QA-SRL Bank 2.0.
- **Relation Extraction.** ZeroShot.
- **Natural Language Inference.** MNLI, SNLI.

**Models** (3). They focus on two classes of models:

- Self-attention models. They use BERT<sub>BASE</sub>.
- Recurrent neural networks. Augment pretrained ELMo with a 300-dim BiLSTM<sup>55</sup>. Then use BiDAF to “to aggregate context (premise) and question (hypothesis) representations to predict an answer span (a label) in question answering (natural language inference).”

**Evaluating Transfer** (4). Say there exists some set of  $N$  examples  $\mathcal{A} = \{(x_i, y_i)\}_{i=1}^N$  that we want to learn rapidly. We’ll denote slices of this with the shorthand  $\mathcal{A}_i \triangleq \{(x_j, y_j)\}_{j=1}^i$ . Let  $\mathbf{W}$  denote the model parameters, and  $\hat{\mathbf{W}}_{\mathcal{A}_i}$  the model parameters trained on  $\mathcal{A}_i$ . Consider an online setup where the **codelength** is computed as:

$$\ell(\mathcal{A}) = \lg |\mathcal{Y}| - \sum_{i=2}^N \lg p(y_i | x_i; \hat{\mathbf{W}}_{\mathcal{A}_{i-1}}) \quad (413)$$

$|\mathcal{Y}|$  is the number of classes

---

<sup>55</sup>So they just stick another BLSTM at the outputs? There are already 2 BiLSTMs though... why do this?

My interpretation of the authors' [vague] interpretation of  $\ell$  is as follows: Alice wants to communicate all  $y_i$ 's in  $\mathcal{A}$  to Bob, who only has access to the  $x_i$ 's. They both agree on a model, initial seed, and learning algorithm. Alice first communicates  $y_1$  to Bob with a uniform code (I guess this would be a binary code of  $\lg |\mathcal{Y}|$  bits). They then train their model on  $\mathcal{A}_1 = (x_1, y_1)$ . Alice uses this model to communicate  $y_2$  (presumably a smaller code than the initial uniform code). I have no idea in practice how this pretend scenario would be implemented, but the gist is just:

*... a model that performs well with a limited number of training examples will be rewarded by having a shorter codelength.*

**Experiments** (5). Here I'll list the main takeaways:

- **Unsupervised Pretraining** (5.1)
  - ✗ Both BERT and ELMo required  $\sim 40k$  training examples to converge in finetuning on SQuAD1.1 and MNLI, even though they were already pretrained.
  - ✓ BERT was slightly faster than ELMo.
  - ✗ Training from scratch for 500k iterations yields much poor results than the pre-trained+finetuned models of  $\sim 40k$  iterations.
  - ✓ Bert has smaller codelengths than *ELMo* (learns from fewer examples better).
- **Beyond Unsupervised Pretraining** (5.2). They also include supervised pretraining into the mix and observe impacts on exact match, F1, and codelength.
  - ✓ Slightly improves EM and F1.
  - ✓ Significantly improves codelength. Kind of misleading/trivial though since the supervised pretraining trains the final layer prediction head that's also used for the final training/eval on SQuAD. Basically gives the models a head-start (nice pun).
- **Generalization** (5.3). Do these models generalize to other datasets from the same task?
  - ✗ High-performing SQuAD models do *not* perform well on other [QA] datasets w/o being provided w/training examples from these datasets.
- **Curriculum and Catastrophic Forgetting** (5.4). Take the best unsupervised+SQuAD-trained models from 5.1 and observe how performance on SQuAD changes as they subsequently finetune on MNLI or TriviaQA.
  - ✗ None of the models that obtain good results on the new task are able to maintain performance on SQuAD.
  - ✓ Training on all tasks simultaneously<sup>56</sup> can result in a model that performs well on all of them (especially BERT).

---

<sup>56</sup>By randomly sampling examples from each task uniformly throughout training.

## **Discussion (6).**

- Need better *transfer and continual learning methods*, such as EWC and progress-and-compress.
- Underexplored promising directions:
  - *memory module* that rapidly adapts to domain shifts and generalizes across tasks.
  - *meta-learning* for NLP that doesn't require knowledge of all tasks we want to evaluate on beforehand.
- *Training curriculum* can have a considerable effect on model performance.

## Scaling Laws for Neural Language Models

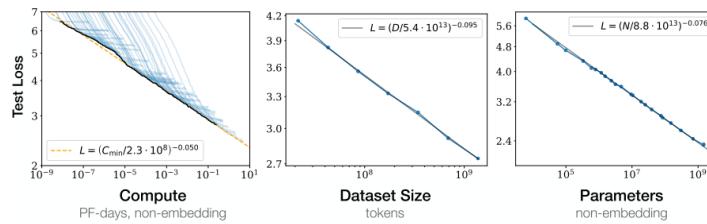
Table of Contents Local

Written by Brandon McKinzie

Kaplan et al., “Scaling Laws for Neural Language Models” *OpenAI*, (Jan 2020).

**Introduction.** Key findings for transformer language models:

- **Performance depends strongly on scale, weakly on model shape.** Performance depends strongly on  $(N, DC)$ :
  - $N$ : number of model parameters.
  - $D$ : size of the dataset.
  - $C$ : amount of compute used for training.
- **Smooth power laws.** Test loss decreases reliably as a power law wrt the 3 aforementioned scale factors.



Notice how the plot above contains predictive fits that we can use for estimating performance gain from increasing each.

- **Universality of overfitting.**
- **Universality of training.** By extrapolating the early part of a training curve, we can roughly predict the loss that would be achieved if we trained for much longer.
- **Transfer improves with test performance.** Transfer to a different distribution incurs a constant penalty but otherwise improves roughly in line with performance on the training set.
- **Sample efficiency.** Larger models are better.
- **Convergence is inefficient.** Just train a fat model and stop way before convergence.
- **Optimal batch size** can be computed via gradient noise scale.

The main scaling law equations are as follows.

$$L(N) = \left( \frac{N_c}{N} \right)^{\alpha_N} \quad \alpha_N \sim 0.076, \quad N_c \sim 8.8e13 \text{ non-embed params} \quad (414)$$

$$L(D) = \left( \frac{D_c}{D} \right)^{\alpha_D} \quad \alpha_D \sim 0.095, \quad D_c \sim 5.4e13 \text{ tokens} \quad (415)$$

$$L(C_{min}) = \left( \frac{C_c^{min}}{C_{min}} \right)^{\alpha_C^{min}} \quad \alpha_C^{min} \sim 0.050, \quad C_c^{min} \sim 3.1e8 \text{ PF-days} \quad (416)$$

$$B_{crit}(L) = \frac{B_*}{L^{\frac{1}{\alpha_B}}} \quad B_* \sim 2e8 \text{ tokens}, \quad \alpha_B \sim 0.21 \quad (417)$$

$$L(N, D) = \left[ \left( \frac{N_c}{N} \right)^{\frac{\alpha_N}{\alpha_D}} + \frac{D_c}{D} \right]^{\alpha_D} \quad (418)$$

Other noteworthy quotes I found throughout the paper:

- We found that results at convergence were largely independent of learning rate schedule. (2.2)
- This suggests that the embedding matrix can be made smaller without impacting performance, as has been seen in recent work. (3.2)
- LSTMs perform as well as Transformers for tokens appearing early in the context, but cannot match the Transformer performance for later tokens. (3.2.1)
- We also observe no dependence [on out-of-domain generalization] on model depth (3.2.2)
- Datasets with about 20 million tokens may represent a different regime for language modeling – overfitting happens very early. (4.2)
- Lower bound on the step that early stopping should occur (5.3)

$$S_{stop}(N, D) \geq \frac{S_c}{[L(N, D) - L(N, \infty)]^{\frac{1}{\alpha_S}}} \quad (419)$$

## Language Models are Few-Shot Learners

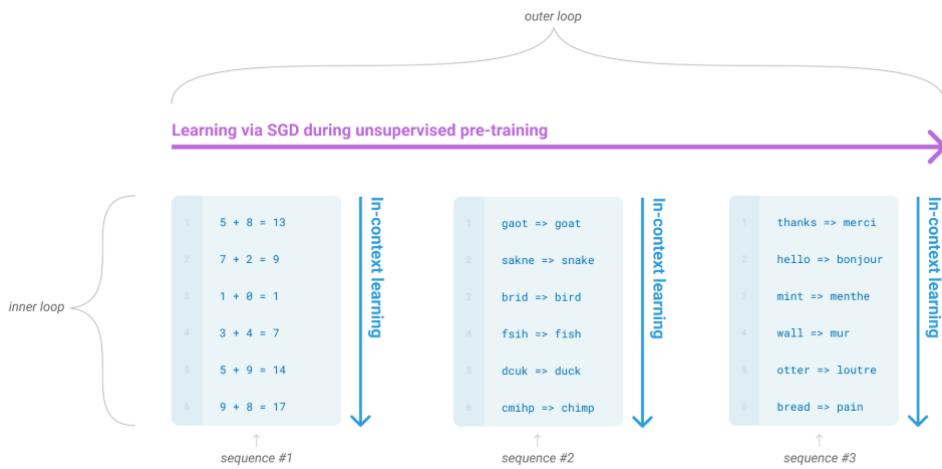
Table of Contents Local

Written by Brandon McKinzie

Brown et al., “Language Models are Few-Shot Learners” *OpenAI*, (May 2020).

**Introduction** (1). Major limitation is needing task-specific data and finetuning. Especially since large pre-trained models can exploit spurious correlations during finetuning. Humans typically just need a *directive* in natural language<sup>57</sup> to learn a new task. Meta-learning, or **in-context learning** for the LM case, may address some of these issues:

*...the model is conditioned on a natural language instruction and/or a few demonstrations of the task and is then expected to complete further instances of the task simply by predicting what comes next.*



Authors train **GPT-3**, a 175 billion parameter LM, and measure:

- **Few-shot learning:** in-context learning. Allow as many demonstrations as will fit into model’s context window<sup>58</sup>.
- **One-shot learning:** one demonstration.
- **Zero-shot learning:** no demonstrations – only a natural-language instruction.

Authors find that larger models perform significantly better, although still struggle on some tasks like ANLI (natural language inference), RACE/QuAC(reading comprehension). Finally, authors include a systematic study of data contamination.

<sup>57</sup>For example, “Tell me if this sentence is positive or negative.”

<sup>58</sup>This is typically about 10 to 100 demonstrations.

**Thoughts:** I wonder whether in-context learning is an inherent property of the model, or whether it's crucial that the training data itself has examples of this type of structure (multiple demonstrations in sequence).

**Approach (2).** Below is the set of models they trained. They are the same architecture as GPT-2, “with the exception that we use alternating dense and locally banded sparse attention patterns in the layers of the transformer, similar to the Sparse Transformer.”

Model Name	$n_{\text{params}}$	$n_{\text{layers}}$	$d_{\text{model}}$	$n_{\text{heads}}$	$d_{\text{head}}$	Batch Size	Learning Rate
GPT-3 Small	125M	12	768	12	64	0.5M	$6.0 \times 10^{-4}$
GPT-3 Medium	350M	24	1024	16	64	0.5M	$3.0 \times 10^{-4}$
GPT-3 Large	760M	24	1536	16	96	0.5M	$2.5 \times 10^{-4}$
GPT-3 XL	1.3B	24	2048	24	128	1M	$2.0 \times 10^{-4}$
GPT-3 2.7B	2.7B	32	2560	32	80	1M	$1.6 \times 10^{-4}$
GPT-3 6.7B	6.7B	32	4096	32	128	2M	$1.2 \times 10^{-4}$
GPT-3 13B	13.0B	40	5140	40	128	2M	$1.0 \times 10^{-4}$
GPT-3 175B or “GPT-3”	175.0B	96	12288	96	128	3.2M	$0.6 \times 10^{-4}$

$d_{ff} = 4 \cdot d_{model}$

$n_{ctx} = 2048$

**Table 2.1:** Sizes, architectures, and learning hyper-parameters (batch size in tokens and learning rate) of the models which we trained. All models were trained for a total of 300 billion tokens.

I found it surprising how the learning rate decreases as model size (and batch size) increases, which is typically the opposite.

**Training Dataset (2.2).** Based on Common Crawl, with some processing:

1. Filtered based on similarity to a range of high-quality reference corpora.
2. Fuzzy deduplication at the document level<sup>59</sup>.
3. Added known high-quality reference corpora to the training mix (WebText, Books1, Books2, English-language Wikipedia).

During training, they sample from the different datasets *proportional to their quality*, as show in the table below.

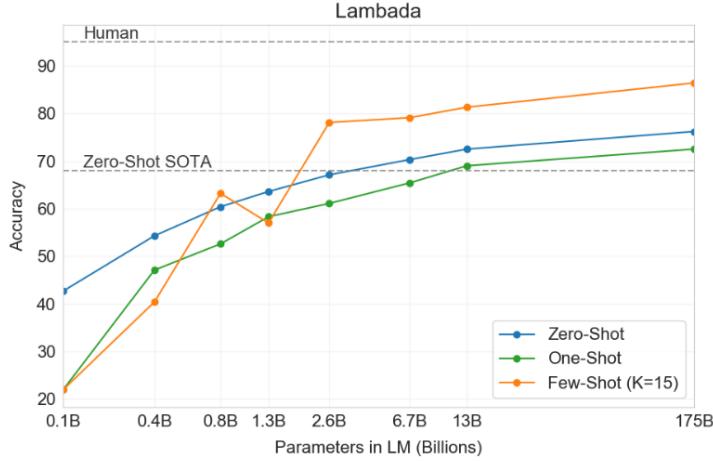
Dataset	Quantity (tokens)	Weight in training mix	Epochs elapsed when training for 300B tokens
Common Crawl (filtered)	410 billion	60%	0.44
WebText2	19 billion	22%	2.9
Books1	12 billion	8%	1.9
Books2	55 billion	8%	0.43
Wikipedia	3 billion	3%	3.4

**Table 2.2: Datasets used to train GPT-3.** “Weight in training mix” refers to the fraction of examples during training that are drawn from a given dataset, which we intentionally do not make proportional to the size of the dataset. As a result, when we train for 300 billion tokens, some datasets are seen up to 3.4 times during training while other datasets are seen less than once.

**Evaluation (2.4).** For few-shot learning, they evaluate each example in the evaluation set by randomly drawing K examples from that task’s *training set* as conditioning, delimited by 1 or 2 newlines depending on the task.

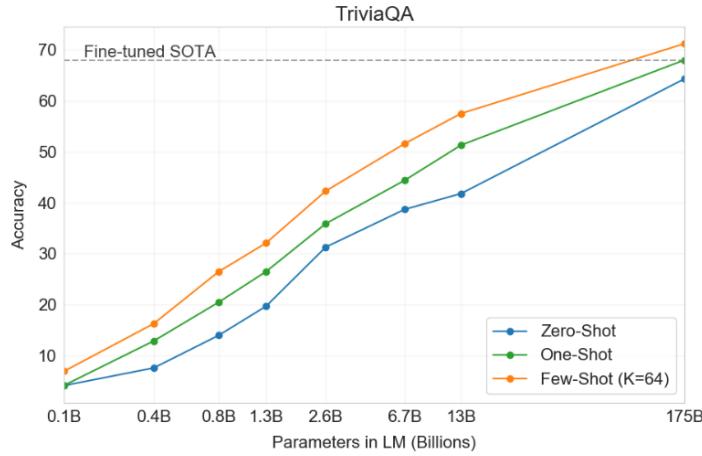
<sup>59</sup>May want to try this technique on some of my LM datasets.

**Language Modeling, Cloze, and Completion Tasks** (3.1). New SOTA perplexity on PTB of 20.5, beating previous SOTA of 35.8. For LAMBADA, model is asked to predict last word of sentences which require reading a paragraph of context.



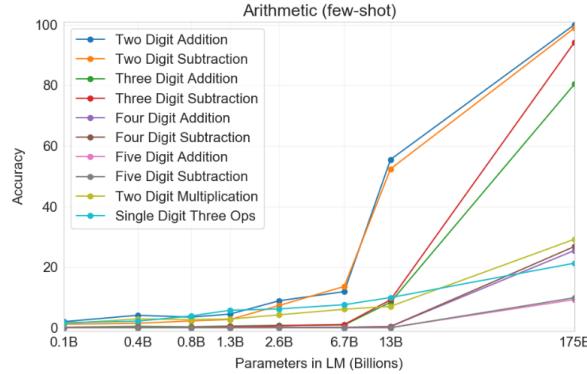
Notice how one-shot performs worse than zero-shot. This is because the demonstration for fill-in-the-blank is a bit confusing and it takes a few demonstrations for the model to see the pattern. GPT3 did *not* obtain SOTA on HellaSwag and StoryCloze 2016 (picking best ending to story/instruction).

**Closed Book Question Answering.** “Closed-book” QA means the model has to answer a question without conditioning on any auxiliary information. The results below show how increasing model size allows the model to effectively absorb more information required to answer questions.



**Translation** (3.3). The authors made no effort to select for/against foreign languages in the training data. Still, it ended up being 93% English and 7% other languages. Results indicate a tendency for translation *into* English to be stronger than translation *from* English.

**Synthetic and Qualitative** (3.9). Below we see that a sharp increase in performance on arithmetic occurs around the 7B parameter point.



## Limitations (5).

*GPT-3 samples still sometimes repeat themselves semantically at the document level, start to lose coherence over sufficiently long passages, contradict themselves, and occasionally contain non-sequitur sentences or paragraphs.*

- Difficulty with “common sense physics.”
- In-context learning does little better than chance when evaluated one-shot or even few-shot on some “comparison” tasks, such as determining if two words are used the same way in a sentence (WIC), or if one sentence implies another (ANLI).
- Bidirectionality may improve downstream results.
- Limits of the pretraining objective itself.
- No grounding in other modalities like video/real-world interactions.
- Poor sample-efficiency compared to humans.
- Expensive/inconvenient to perform inference on.

*... useful language systems (for example virtual assistants) might be better thought of as taking goal-directed actions rather than just making predictions.*

**Further Reading.** Below are a list of references from the paper that I'd like to do a deeper dive on (and why).

- Introduction
  - Larger models  $\neq$  better OOD generalization: [HLW+20].
  - Benchmark performance may exaggerate actual performance on underlying task [GSL+18, NK19].
  - In-context learning [RWC+19].
  - Log loss, which correlates well with many downstream tasks, follows a smooth trend of improvement with scale [KMH+20].
- Approach
  - Sparse transformer [CGRS19].
  - Larger models can typically use a larger batch size, but require a smaller learning rate. [KMH+20, MKAT18]
  - Beam search for free-form completion [RSR+19].
- Results
  - LMs can answer questions about broad factual knowledge without conditioning on auxiliary information (like IR system outputs) [RRS20].
- Limitations
  - Improved fine-tuning performance of BERT vs GPT [RSR+19].
  - Learning the objective function from humans [ZSW+19a].

## Generating Long Sequences with Sparse Transformers

Table of Contents   Local

*Written by Brandon McKinzie*

R. Child, S. Gray, A. Radford, and I. Sutskever, “Generating Long Sequences with Sparse Transformers” *OpenAI*, (Apr 2019).

**Self-Attention Review.** Recall the defining equations for multihead self-attention in transformers, where  $Q=K=V$  are in  $\mathbb{R}^{T \times d_{model}}$ :

$$\text{MultiHead}(Q, K, V) = \text{Concat}(\text{head}_1, \dots, \text{head}_h)W^O \quad (420)$$

$$\text{where } \text{head}_i = \text{Attention}(QW_i^Q, KW_i^K, VW_i^V) \quad (421)$$

$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V \quad (422)$$

In practice, this is implemented in a heavily vectorized way, with a single 1x1 convolution on the input, instead of  $H$  fully connected layers on each of the three  $Q, K, V$ . Roughly speaking, the actual implementation does the following sequence of operations:

1. Using a single weight matrix (I'll omit mention of biases)  $\mathbf{W} \in \mathbb{R}^{M \times 3M}$  (where  $M \equiv d_{model}$ ), we project the inputs  $\mathbf{x}$  from  $\mathbb{R}^{T \times M}$  to  $\mathbb{R}^{T \times 3M}$ . Note that the result is NOT  $Q, K$ , and  $V$ . Conceptually,  $Q, K$ , and  $V$  are each just identical copies of  $x$ , so you could say we “skipped” a step. Rather, we have the *projected versions* of  $Q, K$ , and  $V$ , which contains the slices (heads) we'll pass to the attention functions (next step).
2. Split the last dimension of the *projected versions* of  $Q, K$ , and  $V$  (which is  $M$ ) into  $[H, d_{qkv}]$ . In practice, the practitioner defines the number of heads  $H$  they want, and  $d_{qkv} := M/H$ .
3. Feed each of the respective  $H$  slices into their own Attention function. We can see now that, conceptually, this *is* the same as if we'd taken the original [not projected]  $Q, K$ , and  $V$  and projected them down to  $d_{qkv}$  for each of the  $H$  heads and fed them through, but that's clearly less efficient than the single matrix multiply version seen in implementations.

**Interpreting Visualizations.** The attention maps shown in visualizations are showing the softmax outputs for specific attention heads. These are matrices  $A^{(h)}$  with shape  $(T_q, T_{kv})$ , and normalized as  $1 = \sum_{i=1}^{T_{kv}} A_{q,i}^{(h)}$ .

**Factorized Self-Attention** (4.1). Let  $S = \{S_1, \dots, S_T\}$  denote a **connectivity map** of  $T$  output vectors, where  $S_i$  is itself a set of *input* indices to which the  $i$ th *output* vector attends. For example, traditional decoder LMs define  $S_i = \{j : j \leq i\}$ . Instead of using a single connectivity map for all  $H$  attention heads, **factorized self-attention** has  $P$  heads with different connectivity maps, each a subset of the traditional one:

$$A_i^{(p)} \subset \{j : j \leq i\} \quad \text{with} \quad |A_i^{(p)}| \propto \sqrt[3]{T} \quad (423)$$

and lets  $S_i = A_i^{(p)}$  (**so  $S_i$  is now a set of sets??**)

**Two-Dimensional Factorized Attention** (4.3). Here,  $P=2$ . Let  $\ell$  denote the **stride**, where typically  $\ell \approx \sqrt{T}$ .

$$A_i^{(1)} = \{t, t+1, \dots, i\} \quad \text{for} \quad t = \max(0, i - \ell) \quad (424)$$

$$A_i^{(2)} = \{j : (i - j) \bmod \ell = 0\} \quad (425)$$

**Sparse Transformer** (5). **TODO**

## The Curious Case of Neural Text Degeneration

Table of Contents   Local

*Written by Brandon McKinzie*

Holtzman et al., “The Curious Case of Neural Text Degeneration” *Allen Institute of Artificial Intelligence*, (Feb 2020).

**Nucleus Sampling** (3.1). Given a distribution  $P(x \mid x_{1:i-1})$ , define its top- $p$  vocabulary  $V^{(p)} \subset V$  as the smallest set such that

$$\sum_{x \in V^{(p)}} P(x \mid x_{1:i-1}) \geq p \quad (426)$$

Nucleus sampling renormalizes over this subset and samples from the resulting distribution.

*For high values of  $p$ , this is a small subset of vocabulary that takes up vast majority of the probability mass – the **nucleus**.*

### Further Reading.

- *Finding the optimum argmax sequence from recurrent neural language models or Transformers is not tractable.* Yining Chen, Sorcha Gilroy, Andreas Maletti, Jonathan May, and Kevin Knight. Recurrent neural networks as weighted language recognizers.
- *Maximization-based decoding does not lead to high quality text.* (1) Angela Fan, Mike Lewis, and Yann Dauphin. Hierarchical neural story generation. Ari Holtzman, Jan Buys, Maxwell Forbes, Antoine Bosselut, (2) David Golub, and Yejin Choi. Learning to write with cooperative discriminators.

# EMBEDDINGS AND TRANSFER LEARNING

## CONTENTS

5.1	Simple Baseline for Sentence Embeddings . . . . .	122
5.2	Deep Sentence Embedding Using LSTMs . . . . .	124
5.3	Supervised Universal Sentence Representations (InferSent) . . . . .	127
5.4	Dist. Rep. of Sentences from Unlabeled Data (FastSent) . . . . .	128
5.5	Attention Is All You Need . . . . .	130
5.6	Poincaré Embeddings for Learning Hierarchical Representations . . . . .	134
5.7	Enriching Word Vectors with Subword Information (FastText) . . . . .	136
5.8	Deep Contextualized Word Representations . . . . .	138
5.9	BERT . . . . .	140
5.10	Wasserstein is all you need . . . . .	142
5.11	On the Dimensionality of Word Embedding . . . . .	144
5.12	Universal Language Model Fine-Tuning for Text Classification . . . . .	145
5.13	Hopfield Networks is All You Need . . . . .	147
5.13.1	Review: Hopfield Networks . . . . .	150

## 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:

$$\frac{1}{N} \sum_i^N \frac{a}{a + p(\mathbf{w}_i)} \mathbf{w}_i \quad (427)$$

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:

$$\Pr [w_t | c_t] \propto e^{c_t \cdot v_w} \quad (428)$$

- **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>60</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 428.**
  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 (?)):

$$\Pr [w | c_s] = \alpha p(w) + (1 - \alpha) \frac{e^{\tilde{c}_s \cdot v_w}}{Z_{\tilde{c}_s}} \quad (429)$$

$$\tilde{c}_s = \beta c_0 + (1 - \beta) c_s \quad (430)$$

with  $c_0 \perp c_s$  and  $Z_{\tilde{c}_s}$  is a normalization constant, taken over all  $w \in V$ .

---

<sup>60</sup>Review of MAP:

$$\theta_{MAP} = \arg \max_{\theta} \sum_i \log (p_X(x | \theta)p(\theta))$$

## 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>61</sup>).

**Introduction.** Sentence embeddings are learned using a loss function defined on *sentence pairs*. For example, the well-known Paragraph Vector<sup>62</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>63</sup>

---

<sup>61</sup>Note that this similar to topic extraction.

<sup>62</sup>Q. V. Le and T. Mikolov, “Distributed representations of sentences and documents.”

<sup>63</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>64</sup>  $\mathbf{H}$ :

$$\mathbf{l}(t) = \mathbf{H}\mathbf{x}(t) \quad (431)$$

3. Compute the hidden state  $\mathbf{h}(t)$ , which is the sentence embedding for  $t = T$ , the length of the sentence.

$$\mathbf{h}(t) = \tanh(\mathbf{U}\mathbf{l}(t) + \mathbf{W}\mathbf{h}(t-1) + \mathbf{b}) \quad (432)$$

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:

$$f_i^{(t)} = \sigma \left( b_i^f + \sum_j U_{i,j}^f x_j^{(t)} + \sum_j W_{i,j}^f h_j^{(t-1)} + \sum_j P_{i,j}^f s_j^{(t-1)} \right) \quad (433)$$

$$s_i^{(t)} = f_i^{(t)} \odot s_i^{(t-1)} + g_i^{(t)} \odot \sigma \left( b_i + \sum_j U_{i,j} x_j^{(t)} + \sum_j W_{i,j} h_j^{(t-1)} \right) \quad (434)$$

$$g_i^{(t)} = \sigma \left( b_i^g + \sum_j U_{i,j}^g x_j^{(t)} + \sum_j W_{i,j}^g h_j^{(t-1)} + \sum_j P_{i,j}^g s_j^{(t-1)} \right) \quad (435)$$

$$q_i^{(t)} = \sigma \left( b_i^o + \sum_j U_{i,j}^o x_j^{(t)} + \sum_j W_{i,j}^o h_j^{(t-1)} + \sum_j P_{i,j}^o s_j^{(t)} \right) \quad (436)$$

The final hidden state can then be computed via

$$h_i^{(t)} = \tanh(s_i^{(t)}) \odot q_i^{(t)} \quad (437)$$

---

<sup>64</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:

$$L(\boldsymbol{\Lambda}) = \min_{\boldsymbol{\Lambda}} \left\{ -\log \prod_{r=1}^N \Pr \left[ D_r^+ \mid Q_r \right] \right\} = \min_{\boldsymbol{\Lambda}} \sum_{r=1}^N l_r(\boldsymbol{\Lambda}) \quad (438)$$

$$l_r(\boldsymbol{\Lambda}) = \log \left( 1 + \sum_{j=1}^n e^{-\gamma \cdot \Delta_{r,j}} \right) \quad (439)$$

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>65</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>65</sup>Note that  $\Delta_{r,j} \in [-2, 2]$ . We use  $\gamma$  as a scaling factor so as to expand this range.

## Supervised Universal Sentence Representations (InferSent)

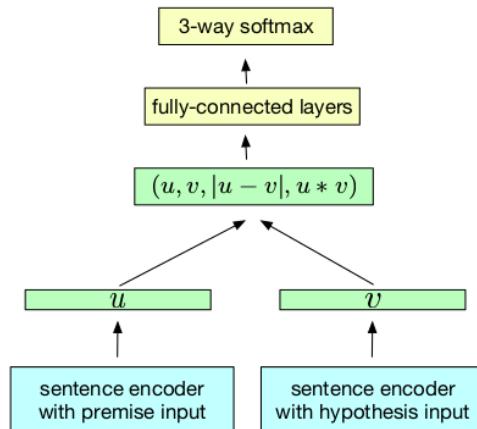
Table of Contents Local

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>66</sup>.

<sup>66</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.

## Attention Is All You Need

Table of Contents Local

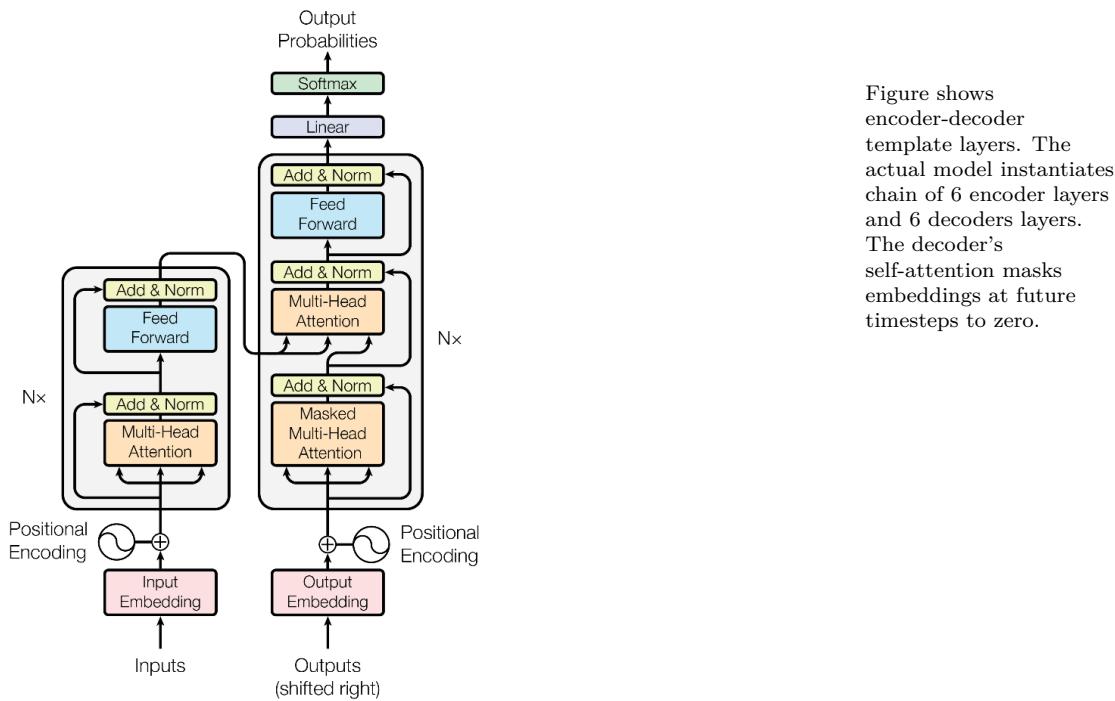
Written by Brandon McKinzie

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.



**Attention.** An attention function can be described as a mapping:

$$\text{Attn}(\text{query}, \{(k_1, v_1), \dots\}) \Rightarrow \sum_i f_n(\text{query}, k_i) v_i \quad (440)$$

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:

$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V \quad (441)$$

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>67</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)

$$\text{Attention}(Q, K, V)_t = \left[ \text{softmax}\left(QK^T\right)V \right]_t \quad (445)$$

$$= \sum_{t'}^T \frac{e^{Q_t \cdot K_{t'}}}{\sum_{t''}^T e^{Q_t \cdot K_{t''}}} V_{t'} \quad (446)$$

---

<sup>67</sup>Assume that  $q$  and  $k$  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

$$\mathbb{E}[q \cdot k] = \mathbb{E}\left[\sum_i^d q_i k_i\right] = \sum_i^d \mathbb{E}[q_i k_i] = \sum_i^d \mathbb{E}[q_i] \mathbb{E}[k_i] = 0 \quad (442)$$

$$\text{Var}[q \cdot k] = \text{Var}\left[\sum_i^d q_i k_i\right] = \sum_i^d \text{Var}[q_i k_i] = \sum_i^d \mathbb{E}[q_i^2 k_i^2] - \mathbb{E}[q_i k_i]^2 \quad (443)$$

$$= \sum_i^d \mathbb{E}[q_i^2] \mathbb{E}[k_i^2] = \sum_i^d \text{Var}[q_i] \text{Var}[k_i] = d \quad (444)$$

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

$$\frac{\partial \text{Softmax}_d(\mathbf{x})}{\partial x_j} = \text{Softmax}_d(\mathbf{x}) (\delta_{dj} - \text{Softmax}_d(\mathbf{x})) \quad (447)$$

- **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.

$$\text{MultiHead}(Q, K, V) = \text{Concat}(\text{head}_1, \dots, \text{head}_h)W^O \quad (448)$$

$$\text{where } \text{head}_i = \text{Attention}(QW_i^Q, KW_i^K, VW_i^V) \quad (449)$$

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:

$$\text{FFN}(x) = \max(0, xW_1 + b_1)W_2 + b_2 \quad (450)$$

The FFN is linear  $\rightarrow$  ReLU  $\rightarrow$  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>68</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>69</sup>:

$$\text{PE}(pos, 2i) = \sin\left(\frac{pos}{10000^{2i/d_{model}}}\right) \quad (451)$$

$$\text{PE}(pos, 2i + 1) = \cos\left(\frac{pos}{10000^{2i/d_{model}}}\right) \quad (452)$$

For inputs to encoder/decoder, the embedding weights are multiplied by  $\sqrt{d_{model}}$

The authors justify this choice:

<sup>68</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>69</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  $\epsilon_{ls} = 0.1$ .

## 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>70</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:

$$H = \{P \mid ||PF_2| - |PF_1|| = 2a\} \quad (453)$$

$$\frac{x^2}{a^2} - \frac{y^2}{b^2} = 1 \quad (454)$$

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

$$H^n = \{x \in \mathbb{R}^{n+1} : x * x = -1\} \quad (455)$$

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>70</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>71</sup>  $(\mathcal{B}^d, g_{\mathbf{x}})$ , where

$$g_{\mathbf{x}} = \left( \frac{2}{1 - \|\mathbf{x}\|^2} \right)^2 g^E \quad (456)$$

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

$$d(\mathbf{u}, \mathbf{v}) = \operatorname{arccosh} \left( 1 + 2 \frac{\|\mathbf{u} - \mathbf{v}\|^2}{(1 - \|\mathbf{u}\|^2)(1 - \|\mathbf{v}\|^2)} \right) \quad (457)$$

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

$$\Theta' \leftarrow \arg \min_{\Theta} \mathcal{L}(\Theta) \quad s.t. \quad \forall \boldsymbol{\theta}_i \in \Theta : \|\boldsymbol{\theta}_i\| < 1 \quad (458)$$

**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

$$\boldsymbol{\theta}_{t+1} = \mathfrak{R}_{\boldsymbol{\theta}_t}(-\eta_t \nabla_R \mathcal{L}(\boldsymbol{\theta}_t)) \quad (459)$$

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>71</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>72</sup> of these representations.

**Skipgram with Negative Sampling.** Since this is based on skipgram, recall the objective of skipgram which is to maximize:

$$\sum_{t=1}^T \sum_{c \in \mathcal{C}_t} \log \Pr [w_c | w_t] \quad (460)$$

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  $s : (w_t, w_c) \mapsto \mathbb{R}$ ,

$$\Pr [w_c | w_t] = \frac{e^{s(w_t, w_c)}}{\sum_{j=1}^W e^{s(w_t, j)}} \quad (461)$$

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:

$$\sum_{t=1}^T \left[ \sum_{c \in \mathcal{C}_t} \ell(s(w_t, w_c)) + \sum_{n \in \mathcal{N}_{t,c}} \ell(-s(w_t, n)) \right] \quad (462)$$

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>72</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  $n = 3$ :

$$\text{where} \longrightarrow \{<wh, whe, her, ere, re>, <where>\} \quad (463)$$

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

$$s(w, c) = \sum_{g \in \mathcal{G}_w} \mathbf{z}_g^T \mathbf{v}_c \quad (464)$$

## 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

$$p(t_1, \dots, t_N) = \prod_{k=1}^N p(t_k | t_1, \dots, t_{k-1}) \quad (465)$$

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

$$\begin{aligned} \sum_{k=1}^N & \left[ \log p(t_k | t_1, \dots, t_{k-1}; \Theta_x, \vec{\Theta}_{LSTM}, \Theta_s) \right. \\ & \left. + \log p(t_k | t_{k+1}, \dots, t_N; \Theta_x, \vec{\Theta}_{LSTM}, \Theta_s) \right] \end{aligned} \quad (466)$$

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.

$$\mathbf{ELMo}_k^{task} = \gamma^{task} \sum_{j=0}^L s_j^{task} \mathbf{h}_{k,j}^{LM} \quad (467)$$

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>73</sup>.
- Halved all embedding and hidden dimensions from the CNN-BIG-LSTM model in *Exploring the Limits of Language Modeling*.
- The  $\mathbf{x}_k$  token representations use 2048 character n-gram convolutional filters followed by two **highway layers**.

---

<sup>73</sup>So the output of some layer, instead of being  $\text{LSTM}(\mathbf{x})$ , becomes  $(\mathbf{x} + \text{LSTM}(\mathbf{x}))$

## BERT

Table of Contents Local

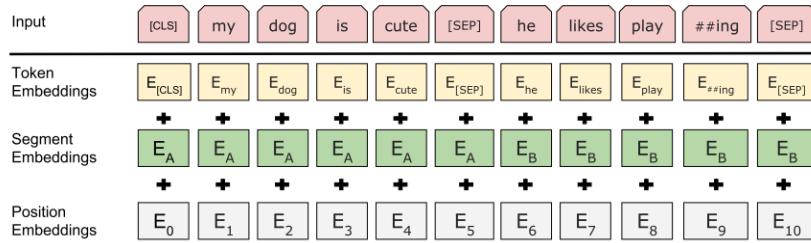
Written by Brandon McKinzie

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>74</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>74</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>75</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>75</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**:

$$\text{OT}(\mu, \nu; \mathbf{M}) = \min_{\mathbf{T} \in \mathbb{R}_+^{n \times m}} \sum_{ij} T_{ij} M_{ij} \quad \text{s.t.} \quad (\forall i) \sum_j T_{ij} = a_i, \quad (\forall j) \sum_i T_{ij} = b_j \quad (468)$$

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

$$W_p(\mu, \nu) \triangleq \text{OT}(\mu, \nu; D_\Omega^p)^{1/p} \quad (469)$$

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

$$P_{\mathbf{V}}^w \triangleq \sum_{c \in \mathcal{C}} (H^w)_c \delta(\mathbf{v}_c) \quad (470)$$

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>76</sup>:

$$\text{OT}(P_{\mathbf{V}}^{w_i}, P_{\mathbf{V}}^{w_j}; D_\Omega^p) := W_p^\lambda(P_{\mathbf{V}}^{w_i}, P_{\mathbf{V}}^{w_j})^p \quad (471)$$

**Concrete Framework** (5). The authors make use of the **shifted positive pointwise mutual information** (SPPMI),  $\mathbf{S}_s^\alpha$ , for computing the word histograms:

$$(H^w)_c := \frac{\mathbf{S}_s^\alpha(w, c)}{\sum_{c' \in \mathcal{C}} \mathbf{S}_s^\alpha(w, c')} \quad (472)$$

$$\mathbf{S}_s^\alpha(w, c) \triangleq \max \left[ \log \left( \frac{\text{Count}(w, c) \sum_{c'} \text{Count}(c')^\alpha}{\text{Count}(w) \text{Count}(c)^\alpha} \right) - \log(s), 0 \right] \quad (473)$$

---

<sup>76</sup>I'm not sure whether the rightmost exponent of  $p$  is a typo in the paper, but that is how it is written.

## 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

$$\mathbf{E} = \mathbf{U}_{1:k} \mathbf{D}_{1:k, 1:k}^\alpha \quad (474)$$

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

$$\text{PIP}(\mathbf{E}) \triangleq \mathbf{E}\mathbf{E}^T \quad (475)$$

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:

$$\|\text{PIP}(\mathbf{E}_1) - \text{PIP}(\mathbf{E}_2)\|_F = \sqrt{\sum_{i,j} \left( \langle \mathbf{w}_i^{(1)}, \mathbf{w}_j^{(1)} \rangle - \langle \mathbf{w}_i^{(2)}, \mathbf{w}_j^{(2)} \rangle \right)^2} \quad (476)$$

## 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.

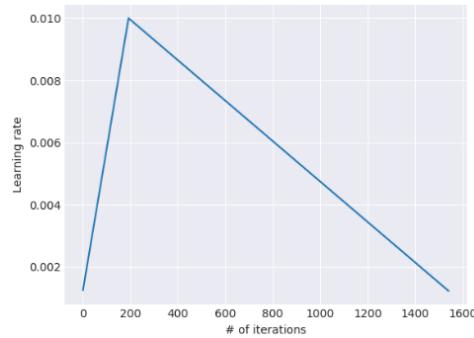
**Target Task LM Fine-tuning** (3.2). For step 2, the authors propose what they call **discriminative fine-tuning** and **slanted triangular learning rates**.

- **Discriminative fine-tuning.** Tune each layer with different learning rates:

$$\theta_t^\ell = \theta_{t-1}^\ell - \eta^\ell \cdot \nabla_{\theta^\ell} J(\theta) \quad (477)$$

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>77</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  $t$ :

$$\eta_t = \eta_{max} \cdot \frac{1 + p \cdot (ratio - 1)}{ratio} \quad (478)$$

$$p = \begin{cases} \frac{t}{cut} & \text{if } t < cut \\ 1 - \frac{t-cut}{cut \cdot (1/cfrac-1)} & \text{otherwise} \end{cases} \quad (479)$$

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>78</sup>:

$$\mathbf{h}_c = [\mathbf{h}_t, \text{maxpool}(\mathbf{H}), \text{meanpool}(\mathbf{H})] \quad (480)$$

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>79</sup> of size  $b$ . They initialize the  $i$ th section with the final state of the model run on section  $i - 1$ .

---

<sup>77</sup>Steps-per-epoch times number of epochs.

<sup>78</sup>Or just as much as we can fit into GPU memory.

<sup>79</sup>Not a fan of how they overloaded this term here.

## Hopfield Networks is All You Need

Table of Contents Local

Written by Brandon McKinzie

Ramsauer et al., “Hopfield Networks is All You Need” *ELLIS Unit Linz and LIT AI Lab*, (Jul 2020).

[See the review sub-section at this end of this note for a review on Hopfield networks]

**From Binary Hopfield Networks to the Transformer.** First we define what the authors mean by “binary modern Hopfield networks.” Let  $\mathbf{X} = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}$  ( $\mathbf{X} \in \{0, 1\}^{d \times N}$ ) be a set of desired binary vectors with  $x_{ij} \in \{-1, +1\} \forall i, j$ . Let  $\boldsymbol{\xi}$  denote the actual state of the units in the Hopfield model. Define the **energy function**  $E$  with **interaction function**  $F$ :

| Name                      | Energy Function                                        | Interaction Function | Storage Capacity |
|---------------------------|--------------------------------------------------------|----------------------|------------------|
| Modern Hopfield           | $-\sum_{i=1}^N F(\boldsymbol{\xi} \cdot \mathbf{x}_i)$ | $x^n$                | $d^{n-1}$        |
| Demircigil Generalization | $-\exp(\text{lse}(1, \mathbf{X}^T \boldsymbol{\xi}))$  | $\exp(x)$            | $2^{d/2}$        |

**Table 1**

where

$$\text{lse}(\beta, \mathbf{x}) \triangleq \beta^{-1} \log \left( \sum_{i=1}^N \exp(\beta x_i) \right) \quad (481)$$

Let  $M = \max_i \|\mathbf{x}_i\|$ . The authors generalize this energy function to continuous-valued patterns:

$$E = -\text{lse}(\beta, \mathbf{X}^T \boldsymbol{\xi}) + \frac{1}{2} \boldsymbol{\xi}^T \boldsymbol{\xi} + \beta^{-1} \log N + \frac{1}{2} M^2 \quad (482)$$

$$\boldsymbol{\xi}^{new} = f(\boldsymbol{\xi}) = \mathbf{X} \mathbf{p} = \mathbf{X} \text{softmax}(\beta \mathbf{X}^T \boldsymbol{\xi}) \quad (483)$$

Novel energy function  
and update rule

with interpretation that  $\mathbf{p} \in \mathbb{R}^N$  contains the softmax probabilities over each of the  $N$  desired memories  $\mathbf{x}_i$ . Notice that the “retrieved” memory  $\boldsymbol{\xi}^{new}$  (after one update) has elements  $\xi_d^{new} = \sum_n p_n(\mathbf{x}_n)_d$ , i.e. a convex combination over the desired memories.

Below I provide a summary of the theorems presented (and proved) by the authors:

### Theorems & Definitions

- **Convergence of E.** Let  $\xi^*$  be a fixed point of the energy function. As  $t \rightarrow \infty$ , the energy converges  $E(\xi^t) \rightarrow E(\xi^*)$ , where  $\xi^{t+1} \leftarrow f(\xi^t)$ .
- **Convergence of  $\xi$ .** As  $t \rightarrow \infty$ ,  $\|\xi^{t+1} - \xi^t\| \rightarrow 0$ .
- **Storage & Retrieval of Patterns.**
  - Assume:  $(\forall \mathbf{x}_i) \exists S_i$ , where  $S_i$  is a sphere around  $\mathbf{x}_i$ .
  - Storage:  $\mathbf{x}_i$  is stored if  $\exists \mathbf{x}_i^* \in S_i$  to which all  $\xi \in S_i$  converge.
  - Retrieval: Given some  $\xi$ , we say  $\mathbf{x}_i$  is retrieved if  $f(\xi)$  converges to  $\mathbf{x}_i^*$ . The retrieval error is  $\|\mathbf{x}_i - \mathbf{x}_i^*\|$ .
- **Separation of Patterns.** Define the separation  $\Delta_i$  of a pattern  $\mathbf{x}_i$  to other patterns by:

$$\Delta_i \triangleq \min_{j,j \neq i} (\mathbf{x}_i^T \mathbf{x}_i - \mathbf{x}_i^T \mathbf{x}_j) = \mathbf{x}_i^T \mathbf{x}_i - \max_{j,j \neq i} \mathbf{x}_i^T \mathbf{x}_j \quad (484)$$

- **Convergence for Well-Separated Patterns.** For some query  $\xi$ , after one update the distance of the new point  $f(\xi)$  to the fixed point  $\mathbf{x}_i^*$  is exponentially small in the separation  $\Delta_i$ . The retrieval error also decreases exponentially with the separation  $\Delta_i$ .

**Hopfield Update is Transformer Attention.** Now consider the case of transformers. Let some input  $\mathbf{y}_i \in \mathbb{R}^{d_{model}}$  be mapped to the Hopfield space of dimension  $d_k$  to represent the query/keys for the transformer attention mechanism:

$$\mathbf{x}_i := \mathbf{W}_K^T \mathbf{y}_i \quad (485)$$

$$\xi_i := \mathbf{W}_Q^T \mathbf{y}_i \quad (486)$$

Denote the full sequence of  $N$  inputs as the matrix  $\mathbf{Y} = [\mathbf{y}_1; \dots; \mathbf{y}_N]^T$ , so  $\mathbf{Y} \in \mathbb{R}^{N \times d_{model}}$ . We can now define the query/key/value matrices:

$$\mathbf{Q} = \mathbf{Y} \mathbf{W}_Q \quad (487)$$

$$\mathbf{K} = \mathbf{X}^T = \mathbf{Y} \mathbf{W}_K \quad (488)$$

$$\mathbf{V} = \mathbf{Y} \mathbf{W}_K \mathbf{W}_V = \mathbf{X}^T \mathbf{W}_V \quad (489)$$

Finally, letting  $\beta := \frac{1}{\sqrt{d_k}}$ , we arrive at the formula for the transformer attention mechanism as our update rule.

**Analysis of Transformer and BERT Models.** The attention mechanism has three fixed point dynamics:

- Global fixed point:** if all  $\mathbf{x}_i$  are not well separated (i.e. if  $p_i \approx \frac{1}{N}$  for all  $i$ ), then  $\xi$  approaches a vector close to the arithmetic mean of  $\{\mathbf{x}_i\}$ .
- Fixed point close to a single pattern:** if all  $\mathbf{x}_i$  are well-separated,  $\xi$  will converge to the most similar pattern  $\mathbf{x}_i$  and  $\mathbf{p}$  will approach  $\mathbf{e}_i$ .
- Metastable state:** when some vectors  $\mathbf{x}_i$  aren't well-separated from the rest, iterates near their associated “metastable state” will converge to it.

All of these can be thought of as special cases of (c), so the authors only consider metastable states, further categorizing into 4 classes:

- (I) Averaging over a *very* large number of patterns (global fixed point).
- (II) Averaging over a large number of patterns (large metastable state).
- (II) Averaging over a medium number of patterns (medium metastable state).
- (IV) Averaging over a small number of patterns (fixed point close to a single pattern).

The authors fed a sequence of  $N$  tokens through BERT and plotted the distribution over  $k$  – the number of patterns required to sum up the softmax values to 0.9 – for each of the  $N$  tokens in each attention head, and classified the heads into one of the four categories above.



---

### 5.13.1 REVIEW: HOPFIELD NETWORKS

---

(Mix of notes from Ch. 2 of Hertz et al. and Ch. 42 of Mackay)

First, recall the general problem/task of **associative memory**: store a set of  $p$  binary patterns  $\{\mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \dots, \mathbf{x}^{(p)}\}$ , where each  $\mathbf{x}^{(i)} \in \{-1, +1\}^N$  is a binary vector<sup>80</sup>, such that when presented with some new pattern  $\hat{\mathbf{x}}$ , the network responds by producing whichever stored pattern most closely resembles  $\hat{\mathbf{x}}$ <sup>81</sup>. Below is the definition of the Hopfield model for accomplishing this behavior.

#### Hopfield Network

Let  $w_{ij}$  denote connection from neuron  $j$  to neuron  $i$ , and let  $x_i$  denote the activity of neuron  $i$ .

- **Architecture.** Consists of  $I$  fully-connected neurons with bidirectional  $((\forall w_{ij}) \exists w_{ji})$  symmetric ( $w_{ij} = w_{ji}$ ) connections<sup>82</sup>. Biases  $w_{i0}$  may be included.
- **Activity rule.** Each neuron updates its state via a threshold:

$$x(a) = \Theta(a) \equiv \begin{cases} 1 & a \geq 0 \\ -1 & a < 0 \end{cases} \quad (490)$$

Due to feedback (fully-connected with symmetric weights), we have to specify an order for computing these updates.

- **Synchronous.** All neurons compute their activations,  $a_i = \mathbf{w}_i \cdot \mathbf{x}$ , then update their states simultaneously,  $\mathbf{x} := \Theta(\mathbf{a})$ .
- **Asynchronous.** One neuron at a time computes its activations and updates its state. The order may be fixed or random.
- **Learning rule.** Our goal is to make a set of desired “memories”  $\{\mathbf{x}^{(n)}\}$  be **stable states** (configurations) of [the activity rule of] our network.

$$[\text{Hebb rule}] \quad w_{ij} = \eta \sum_n x_i^{(n)} x_j^{(n)} \quad (491)$$

---

<sup>80</sup>with  $+1/-1$  instead of  $1/0$  for convenience as we'll see later.

<sup>81</sup>This notion of setting a network's activations in order for it to then re-adjust them to a different desired set of activations is referred to as a **content-addressable** memory.

<sup>82</sup>There are no self-connections, i.e.  $w_{ii} = 0 \forall i$ .

**One Pattern.** Consider the base case where we only want to learn *one pattern*  $\mathbf{x}^{(1)}$ . In other words, if we set the network activations  $\mathbf{x} := \mathbf{x}^{(1)}$ , and subsequently let the network apply its activity rule one or more times (analogous to a forward pass), we want it to still remain in the configuration equal to  $\mathbf{x}^{(1)}$ . This can easily be accomplished if we define the connections to satisfy

$$w_{ij} \propto x_i^{(1)} x_j^{(1)} \quad (492)$$

since this would result in

$$a_i := \sum_j w_{ij} x_j^{(1)} = \sum_j x_i^{(1)} x_j^{(1)} x_j^{(1)} = N x_i^{(1)} \propto x_i^{(1)} \quad (493)$$

$$x_i := \Theta(a_i) = \Theta(x_i^{(1)}) = x_i^{(1)} \quad (494)$$

Furthermore, if fewer than half of the bits of the starting pattern are wrong<sup>83</sup>, it will still result in the correct response. For example, say  $N_c$  bits are correct, and  $N_w = N - N_c$  bits are wrong in the starting configuration ( $N_c > N_w$ ).

$$a_i = N_c x_i^{(1)} + \sum_{x_j \in \text{wrong neurons}} (x_i^{(1)} x_j^{(1)}) x_j \quad (495)$$

$$= N_c x_i^{(1)} - N_w x_i^{(1)} \quad (496)$$

$$= (N_c - N_w) x_i^{(1)} \quad (497)$$

and thus it still results in the correct sign! In such a scenario, we say that  $\mathbf{x}$  is an **attractor**. Also notice that, if  $N_w > N_c$  instead, we'd end up in the **reversed state**  $-\mathbf{x}^{(1)}$ .

**Multiple Patterns.** We can naturally generalize our observations regarding the situation when fewer than half of the bits are wrong to the scenario where we want to store a set of  $p > 1$  patterns  $\{\mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \dots, \mathbf{x}^{(p)}\}$ . Consider what happens if we define our weight update rule as

$$w_{ij} = \frac{1}{N} \sum_{\mu=1}^p x_i^{(\mu)} x_j^{(\mu)} \quad (498)$$

Notice how, if we set the starting configuration of the network to be one of the desired patterns, then this expands out to the same situation as having fewer than half of the bits begin wrong in the previous subsection (think about it).

---

<sup>83</sup>i.e. if the activation of neuron  $i$  is not equal to  $x_i^{(1)}$

# VISION AND LANGUAGE

## CONTENTS

|        |                                                                                                                    |     |
|--------|--------------------------------------------------------------------------------------------------------------------|-----|
| 6.1    | Show, Ask, Attend, and Answer . . . . .                                                                            | 153 |
| 6.2    | Neural Module Networks . . . . .                                                                                   | 154 |
| 6.3    | Learning to Compose Neural Networks for QA . . . . .                                                               | 156 |
| 6.4    | End-to-End Module Networks for VQA . . . . .                                                                       | 158 |
| 6.5    | Fast Multi-language LSTM-based Online Handwriting Recognition . . . . .                                            | 160 |
| 6.6    | Multi-Language Online Handwriting Recognition . . . . .                                                            | 161 |
| 6.7    | Learning Two-Branch Neural Networks for Image-Text Matching Tasks . . . . .                                        | 163 |
| 6.8    | Learning by Abstraction: The Neural State Machine . . . . .                                                        | 165 |
| 6.9    | Tags and Attributes based Visual Commonsense Reasoning Baselines . . . . .                                         | 168 |
| 6.10   | From Recognition to Cognition: Visual Commonsense Reasoning . . . . .                                              | 170 |
| 6.11   | GQA . . . . .                                                                                                      | 172 |
| 6.12   | A negative case analysis of visual grounding methods for VQA . . . . .                                             | 175 |
| 6.13   | Bottom-Up and Top-Down Attention for Image Captioning and Visual Question Answering . . . . .                      | 177 |
| 6.14   | ViLBERT . . . . .                                                                                                  | 179 |
| 6.14.1 | 12-in-1: Multi-Task Vision and Language Representation Learning . . . . .                                          | 180 |
| 6.15   | Language-Conditioned Graph Networks for Relational Reasoning . . . . .                                             | 181 |
| 6.16   | Towards VQA Models That Can Read . . . . .                                                                         | 183 |
| 6.17   | Accuracy vs. Complexity: A Trade-off in VQA Models . . . . .                                                       | 184 |
| 6.18   | Iterative Answer Prediction with Multimodal Transformers for TextVQA . . . . .                                     | 186 |
| 6.19   | CVPR 2020 . . . . .                                                                                                | 189 |
| 6.19.1 | Unified Visual-Semantic Embeddings: Bridging Vision and Language with Structured Meaning Representations . . . . . | 189 |
| 6.19.2 | Counterfactual Samples Synthesizing for Robust Visual Question Answering . . . . .                                 | 190 |
| 6.19.3 | What Makes Training Multi-modal Classification Networks Hard? . . . . .                                            | 191 |

## 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

**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>84</sup> in the training data:

$$\hat{a} = \arg \max_a \Pr [a | I, q] \quad (499)$$

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>85</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 in parallel, like multi-head attention). This represents the scores for the attention function. The attention output, as usual, is computed as

$$\mathbf{x}_c = \sum_{\ell} \alpha_{c,\ell} \phi_{\ell} \quad (500)$$

$$\alpha_{c,\ell} \propto \exp F_c(s, \phi_{\ell}) \quad (501)$$

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  $a_i$ ,  $i \in [1..M]$ .

**Dataset.** Although 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.

---

<sup>84</sup>Same approach as how we define vocabulary in language modeling tasks.

<sup>85</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.

## Neural Module Networks

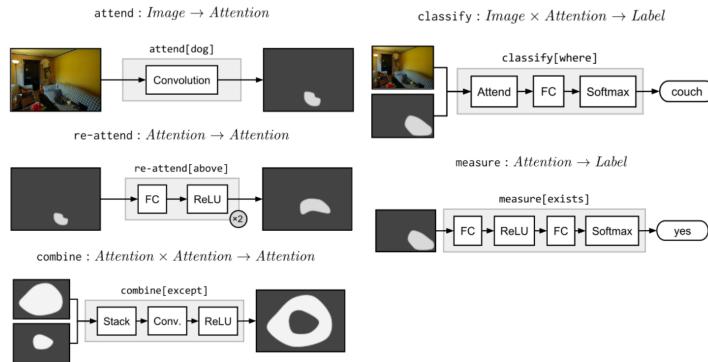
Table of Contents Local

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>86</sup>. The final prediction is a geometric average of the LSTM output probabilities and the root module output probabilities.

---

<sup>86</sup>This is the same distribution that the root module is trying to predict

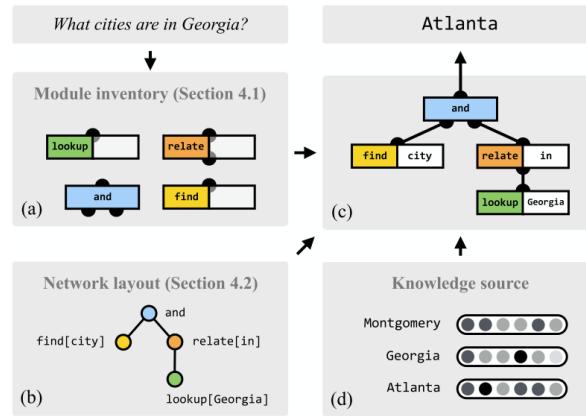
## Learning to Compose Neural Networks for QA

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

$$p_z(y | w) = (\llbracket z \rrbracket_w)_y \quad (502)$$

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 ( $\sigma \equiv \text{ReLU}$ ,  $sm \equiv \text{softmax}$ ):

$$\llbracket \text{lookup}[i] \rrbracket = e_{f(i)} \quad (503)$$

$$\llbracket \text{find}[i] \rrbracket = sm(a \odot \sigma(Bv^i \oplus CW \oplus d)) \quad (504)$$

$$\llbracket \text{relate}[i](h) \rrbracket = sm(a \odot \sigma(Bv^i \oplus CW \oplus D\bar{w}(h) \oplus e)) \quad (505)$$

$$\llbracket \text{and}(h^1, h^2, \dots) \rrbracket = h^1 \odot h^2 \odot \dots \quad (506)$$

$$\llbracket \text{describe}[i](h) \rrbracket = sm(A\sigma(B\bar{w}(h)) + v^i) \quad (507)$$

$$\llbracket \text{exists}(h) \rrbracket = sm\left(\left(\max_k h_k\right)a + b\right) \quad (508)$$

$$\bar{w}(h) \triangleq \sum_k h_k w^{(k)}$$

To train, maximize

$$\sum_{(w,y,z)} \log p_z(y | w; \theta_e) \quad (509)$$

**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  $m$ :

$$x_{txt}^{(m)} = \sum_{i=1}^T \alpha_i^{(m)} w_i \quad (510)$$

**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>87</sup>.

---

<sup>87</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>88</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>89</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>88</sup> **Segmentation/cut point**: a point at which another character may start. **Segment**: the (partial) strokes between 2 consecutive segmentation points.

<sup>89</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.

## Learning Two-Branch Neural Networks for Image-Text Matching Tasks

Table of Contents   Local

*Written by Brandon McKinzie*

Wang et al., “Learning Two-Branch Neural Networks for Image-Text Matching Tasks” *UoI Urbana-Champaign, Georgia Tech*, (May 2018).

**Introduction.** How can we measure semantic similarity between visual data and text data?  
Authors propose two approaches:

- **Embedding Network:** map image/text to joint latent space of common dimensionality.  
Train with bidirectional ranking loss augmented with neighborhood information.
- **Similarity Network:** binary classifier outputting +1 for positive pairs and -1 for negative pairs.

### Overview of Image-Text Tasks (3.1).

- **Phrase Localization:** (image, entity mention)  $\mapsto$  bounding box.
- **Bidirectional Image-Sentence Retrieval:** either (image, sentence database)  $\mapsto$  best matching sentence OR (sentence, image database)  $\mapsto$  best matching image.

**Embedding Network** (3.2). Let  $V$  and  $T$  denote the collections of training images/regions and sentences/phrases. For a given visual input<sup>90</sup>  $v$ , let  $T^+$  and  $T^-$  denote matching/non-matching text samples<sup>91</sup>. Let  $N(v)$  denote the **neighborhood** of image/region  $v$ : the set of images/regions described by the same text as  $v$  (and defined similarly for the text case). They define their margin-based loss function with neighborhood-preserving constraints as:

They set  $m = 0.05$

---

<sup>90</sup>Visual inputs can be either full images or individual regions. For phrase localization, it's a region.

<sup>91</sup>I'll denote individual samples from these sets with lowercase, like  $t^+$  for a positive example, and with the subscript  $i$  implied/obvious given the context.

$$L_{marg}(a, b^+, b^-) \triangleq \left[ m + d(a, b^+) - d(a, b^-) \right]_+ \quad (511)$$

$$L_1(v) = \sum_{t^+, t^-} L_{marg}(v, t^+, t^-) \quad (512)$$

$$L_2(t) = \sum_{v^+, v^-} L_{marg}(t, v^+, v^-) \quad (513)$$

$$L_3(v) = \sum_{v_j \in N(v), v_k \notin N(v)} L_{marg}(v, v_j, v_k) \quad (514)$$

$$L_4(t) = \sum_{t_j \in N(t), t_k \notin N(t)} L_{marg}(t, t_j, t_k) \quad (515)$$

$$L(V, T) = \sum_{v \in V} [\lambda_1 L_1(v) + \lambda_3 L_3(v)] + \sum_{t \in T} [\lambda_2 L_2(t) + \lambda_4 L_4(t)] \quad (516)$$

**TODO:** write details of triplet sampling algorithm(s).

## Learning by Abstraction: The Neural State Machine

Table of Contents   Local

*Written by Brandon McKinzie*

D. Hudson and C. Manning, “Learning by Abstraction: The Neural State Machine” *Stanford*, (Nov 2019).

### Neural State Machine

Defined as a tuple  $(C, S, E, \{r_i\}_{i=0}^N, p_0, \delta)$

- $C$ : model’s **alphabet**, consisting of a set of concepts, embedded as learned vectors.
- $S$ : states.
- $E$ : directed edges that specify valid transitions between states.
- $r_i$ : sequence of **instructions**, each of dimension  $d$ , subsequently passed as input to  $\delta$ .
- $p_0: S \mapsto [0, 1]$ . Initial-state probability distribution.
- $\delta_{S,E}: p_i \times r_i \mapsto p_{i+1}$ . **State transition function**: “a neural module that at each step  $i$  considers the distribution  $p_i$  over the states as well as an input instruction  $r_i$ , and uses it to redistribute the probability along the edges, yielding an updated state distribution  $p_{i+1}$ .”

**Concept Vocabulary** (3.1). The authors define  $L$  types of **attributes** and also include identity and relations (the “+2”) to define their **concept vocabulary**  $C$ :

- $C_0$ : object identity (e.g. *cat, shirt*)
- $\cup_{i=1}^L C_i$ : attributes (e.g. *colors, materials*)
- $C_{L+1}$ : relations (e.g. *holding, behind*)

for a total of  $L + 2$  property types. They store an associated embedding for each property type and denote this set of embeddings as  $D$ .

**States and Edge Transitions** (3.2). Construct a **probabilistic scene graph**:

- Nodes (states) represent objects in an image. Each node in  $S$  is accompanied by:
  - A bounding box.
  - **Dense visual features**.
  - Collection of discrete probability distributions  $\{P_i\}_{i=0}^L$  for each of the object’s  $L + 1$  **semantic properties**<sup>92</sup>, defined over the concept vocabulary  $\{C_i\}_{i=0}^L$ .
- Edges (valid transitions) represent relations between objects. Each edge has an associated probability distribution  $P_{L+1}$  of its **semantic type**<sup>93</sup> among the concepts in  $C_{L+1}$ .

---

<sup>92</sup>Such as color, material, shape, etc.”

<sup>93</sup>Such as “on top of”, “eating”

Then, compute structured embedded representations for each state and transition edge. For each state  $s \in S$ , define a set of  $L + 1$  **property variables**  $\{s^j\}_{j=0}^L$ :

$$s^j = \sum_{c_k \in C_j} P_j(k) c_k \quad (517)$$

where  $c_k \in C_j$  denotes each embedded concept of the  $j$ th property type, and  $P_j$  refers to the corresponding distribution over these concepts [for state  $s$ ].

Similarly, for each edge  $e \in E$ , compute relation embedding:

$$e' = \sum_{c_k \in C_{L+1}} P_{L+1}(k) c_k \quad (518)$$

**Reasoning Instructions** (3.3). Translate the input questions into a sequence of reasoning instructions:

1. Embed each question word with GloVe.
2. Compute a similarity-based distribution

$$P_i = \text{softmax} \left( w_i^T \mathbf{W} C \right) \quad (519)$$

where  $\mathbf{W}$  is initialized to the identity matrix and  $C$  denotes the “matrix of all embedded concepts along with an additional learned default embedding  $c'$  to account for structural or other non-content words”.

3. Translate each word into a concept-based representation<sup>94</sup>:

$$v_i = P_i(c') w_i + \sum_{c \in C \setminus \{c'\}} P_i(c) c \quad (520)$$

4. Feed through attention encoder-decoder to obtain instruction sequence.
  - (a) Feed sequence of  $P$  normalized words  $V^{P \times d} = \{v_i\}_{i=1}^P$  through LSTM encoder to obtain final state  $q$  (representing the question).
  - (b) Run decoder for fixed number of steps  $N+1$  (?), yielding  $N+1$  hidden states  $\{h_i\}_{i=0}^N$ .
  - (c) Transform each decoder hidden state  $h_i$  into a corresponding **reasoning instruction**:

$$r_i = \text{softmax}(h_i V^T) V \quad (521)$$

---

<sup>94</sup>They also subsequently refer to these as the “normalized” words.

### Model Simulation (3.4).

1. Determine each *instruction's type* by computing distribution over the  $L + 2$  embedded property types in  $D$ :

$$R_i = \text{softmax}(r_i^T \circ D) \quad (522)$$

and denote  $r_{i'} \triangleq R_i(L + 1)$  that corresponds to the relation property.

2. Compare each instruction  $r_i$  for all  $s \in S$  and  $e \in E$  via a **relevance score**:

$$\gamma_i(s) = \sigma \left( \sum_{j=0}^L R_i(j) (r_i \circ \mathbf{W}_j s^j) \right) \quad (523)$$

$$\gamma_i(e) = \sigma (r_i \circ \mathbf{W}_{L+1} e') \quad (524)$$

3. Shift attention  $p_i$  from current nodes  $s \in S$  to their most relevant neighbors:

$$p_{i+1}^s = \text{softmax}_{s \in S} (\mathbf{W}_s \cdot \gamma_i(s)) \quad (525)$$

$$p_{i+1}^r = \text{softmax}_{s \in S} \left( \mathbf{W}_r \cdot \sum_{(s',s) \in E} p_i(s') \cdot \gamma_i((s',s)) \right) \quad (526)$$

$$p_{i+1} = r_i' \cdot p_{i+1}^r + (1 - r_i') \cdot p_{i+1}^s \quad (527)$$

4. Repeat over  $N$  steps.

Finally to predict the answer, feed a 2-layer FC-softmax classifier the concatenated vector  $[q; m]$ , where

$$m = \sum_{s \in S} p_N(s) \left( \sum_{j=0}^L R_N(j) \cdot s^j \right) \quad (528)$$

“reflects the information extracted from the final states as guided by the final reasoning instruction  $r_N$ .”

## Tags and Attributes based Visual Commonsense Reasoning Baselines

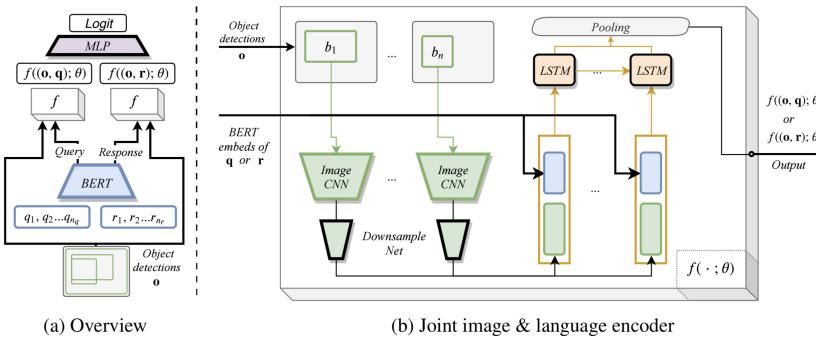
Table of Contents Local

Written by Brandon McKinzie

Lin et al., “TAB-VCR: Tags and Attributes based Visual Commonsense Reasoning Baselines” *UoI Urbana-Champaign*, (Jan 2020).

**TL;DR:** Smaller/simpler model can get SOTA on VCR. Architecture summary:

1. Feed collection of object detections through **Faster RCNN**.
2. Feed query and responses [separately] through **BERT**.
3. Concat word embeddings with their associated image embeddings.
4. Pass the concatenated image/word embeddings through an LSTM<sup>95</sup>.
5. Pool outputs of LSTM to get a joint image+sentence embedding.



**Attribute-Based VCR** (3). Each data point is denoted by a triplet  $(\mathbf{o}, \mathbf{q}, \{\mathbf{r}_1, \dots, \mathbf{r}_4\})$ :

- Image data denoted by the set  $\mathbf{o} = \{o_i\}_{i=1}^{n_o}$  where each  $o_i$  consists of a bounding box  $b_i$  and class label  $l_i \in \mathcal{L}$ .
- Query sequence  $\mathbf{q} = \{q_i\}_{i=1}^{n_q}$  where each  $q_i$  is either a word in  $\mathcal{V}$  or a **tag** referring to a bounding box in  $\mathbf{o}$ .
- Four responses, each denoted as  $\mathbf{r} = \{r_i\}_{i=1}^{n_r}$ , where each  $r_i$  is either a word in  $\mathcal{V}$  or a **tag** referring to a bounding box in  $\mathbf{o}$ .

<sup>95</sup>This decision makes little sense to me, considering there is no meaningful temporal dependence of the data we are sequentially feeding to the LSTM.

**Improving Visual Representation & Image-Text Grounding** (3.3). Authors use a pretrained attribute classifier which augments every detected bounding box  $b_i$  (provided by dataset) with a set of attributes (color, texture, size, emotions). Authors also develop the algorithm below for finding new text-image groundings (*tags*) missing from the original VCR dataset.

---

**Algorithm 1** Finding new tags

---

```

Forward pass through image CNN to obtain object detections  $\hat{o}$ 
 $\hat{L} \leftarrow \text{set}(\text{all class labels in } \hat{o})$ 
for  $w \in \mathbf{w}$  where  $w \in \{\mathbf{q}, \mathbf{r}\}$  do
    if  $w$  is tag then  $w \leftarrow \text{remap}(w)$ 
     $\text{new\_tags} \leftarrow \{\}$ 
    for  $w \in \mathbf{w}$  where  $w \in \{\mathbf{q}, \mathbf{r}\}$  do
        if  $(\text{pos\_tag}(w|w) \in \{\text{NN}, \text{NNS}\})$  and  $(\text{wsd\_synset}(w, w) \text{ has a noun})$  then
            if  $w \in \hat{L}$  then ▷ Direct match between word and detections
                 $\text{new\_detections} \leftarrow \text{detections in } \hat{o} \text{ corresponding to } w$ 
                add  $(w, \text{new\_detections})$  to  $\text{new\_tags}$ 
            else ▷ Use word sense to match word and detections
                 $\text{max\_wup} \leftarrow 0$ 
                 $\text{word\_lemma} \leftarrow \text{lemma}(w)$ 
                 $\text{word\_sense} \leftarrow \text{first\_synset}(\text{word\_lemma})$ 
                for  $\hat{l} \in \hat{L}$  do
                    if  $\text{wup\_similarity}(\text{first\_synset}(\hat{l}), \text{word\_sense}) > \text{max\_wup}$  then
                         $\text{max\_wup} \leftarrow \text{wup\_similarity}(\text{first\_synset}(\hat{l}), \text{word\_sense})$ 
                         $\text{best\_label} \leftarrow \hat{l}$ 
                if  $\text{max\_wup} > k$  then
                     $\text{new\_detections} \leftarrow \text{detections in } \hat{o} \text{ corresponding to best\_label}$ 
                    add  $(w, \text{new\_detections})$  to  $\text{new\_tags}$ 

```

---

## From Recognition to Cognition: Visual Commonsense Reasoning

Table of Contents Local

Written by Brandon McKinzie

Zellers et al., “From Recognition to Cognition: Visual Commonsense Reasoning” *UW and Allen Institute for Artificial Intelligence*, (Mar 2019).

**Task Overview** (2). The primary task is **staged answering and justification** in a multiple-choice setting ( $Q \rightarrow AR$ ):

*Given a question along with four answer choices, a model must first select the right answer. If its answer was correct, then it is provided four rationale choices (that could purportedly justify its correct answer), and it must select the correct rationale.*

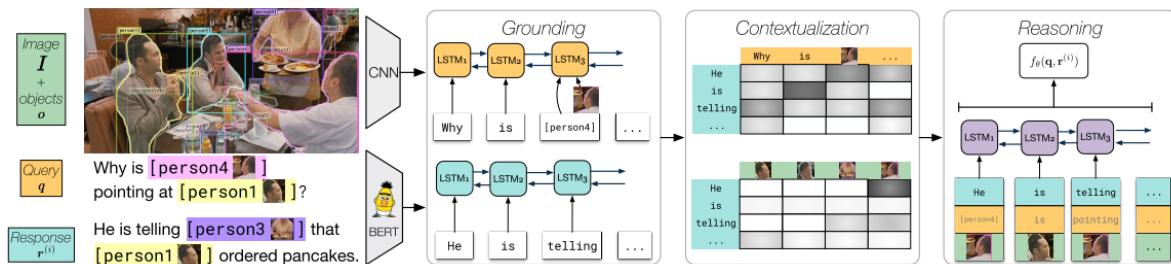
Authors present the following multiple-choice *sub-tasks*:

- **Answering**  $Q \rightarrow A$ : Given query, select the correct answer.
- **Justification**  $QA \rightarrow R$ : Given concatenated [query, correct answer], select the correct rationale.

**VCR subtask**

A single example consists of an image  $I$  and:

- $\mathbf{o}$ : sequence of object detections. Each  $o_i$  consists of:
  - bounding box  $b$
  - segmentation mask  $m$
  - class label  $\ell_i \in \mathcal{L}$
- $\mathbf{q}$ : query sequence, where each  $q_i$  is either a word in  $\mathcal{V}$  or a tag referring to an object in  $\mathbf{o}$ .
- Set of  $N = 4$  responses  $\mathbf{r}^{(i)}$ .



**Recognition to Cognition Networks** (5). Authors propose their three-stage R2C model for training on VCR.

1. **Grounding:**

- (a) CNN learns object-level features (embeds each of the objects  $o_i$  in  $\mathbf{o}$ ). BERT is used for getting contextualized word embeddings.
- (b) Bidirectional LSTM is feed the question and current response separately (but shared weights). It seems like the CNN outputs are only used when feeding in a question/response token referring/pointing to an object in  $\mathbf{o}$ ? Very unclear/poorly written.

2. **Contextualization:** Given LSTM output sequences for response and query be  $\mathbf{r}$  and  $\mathbf{q}$ , respectively. They compute an attended query representation  $\hat{\mathbf{q}}_i$  and attended object representation  $\hat{\mathbf{o}}_i$ :

$$\hat{\mathbf{q}}_i = \sum_j \alpha_{ij} \mathbf{q}_j \quad \text{with} \quad \alpha_{ij} = \text{softmax}_j(\mathbf{r}_i \mathbf{W} \mathbf{q}_j) \quad (529)$$

$$\hat{\mathbf{o}}_i = \sum_j \beta_{ij} \mathbf{o}_j \quad \text{with} \quad \beta_{ij} = \text{softmax}_j(\mathbf{r}_i \mathbf{W} \mathbf{o}_j) \quad (530)$$

with interpretation that  $\hat{\mathbf{q}}_i$  is the contextualized version of response token  $\mathbf{r}_i$  attended over all *query* tokens. Similarly,  $\hat{\mathbf{o}}_i$  gives an alternative representation of response token  $\mathbf{r}_i$  attended over all *object* tokens.

3. **Reasoning:**

- Feed the element-wise concatenated sequence of  $[\mathbf{r}_i, \hat{\mathbf{q}}_i, \hat{\mathbf{o}}_i]$  to another BiLSTM.
- Concatenate the sequence output from this BiLSTM with the original question/answer representations  $\mathbf{q}$  and  $\mathbf{r}$ .
- Max-pool this fat concatenated sequence (**along time dim presumably?**).
- Feed to MLP which predicts a logit for query-response compatibility.

**Results** (6). This paper's results are particularly interesting because they highlight both (1) how well humans perform on this dataset, (2) how poorly VQA models perform, and (3) how text-only models like BERT can substantially outperform VQA models in some scenarios!

| Model                | $Q \rightarrow A$ |             | $QA \rightarrow R$ |             | $Q \rightarrow AR$ |             |
|----------------------|-------------------|-------------|--------------------|-------------|--------------------|-------------|
|                      | Val               | Test        | Val                | Test        | Val                | Test        |
| Chance               | 25.0              | 25.0        | 25.0               | 25.0        | 6.2                | 6.2         |
| BERT                 | 53.8              | 53.9        | 64.1               | 64.5        | 34.8               | 35.0        |
| BERT (response only) | 27.6              | 27.7        | 26.3               | 26.2        | 7.6                | 7.3         |
| ESIM+ELMo            | 45.8              | 45.9        | 55.0               | 55.1        | 25.3               | 25.6        |
| LSTM+ELMo            | 28.1              | 28.3        | 28.7               | 28.5        | 8.3                | 8.4         |
| RevisitedVQA [38]    | 39.4              | 40.5        | 34.0               | 33.7        | 13.5               | 13.8        |
| BottomUpTopDown[4]   | 42.8              | 44.1        | 25.1               | 25.1        | 10.7               | 11.0        |
| MLB [42]             | 45.5              | 46.2        | 36.1               | 36.8        | 17.0               | 17.2        |
| MUTAN [6]            | 44.4              | 45.5        | 32.0               | 32.2        | 14.6               | 14.6        |
| <b>R2C</b>           | <b>63.8</b>       | <b>65.1</b> | <b>67.2</b>        | <b>67.3</b> | <b>43.1</b>        | <b>44.0</b> |
| Human                | 91.0              |             | 93.0               |             | 85.0               |             |

However, reading the details in this section suggests (to me) that the comparisons here aren't that fair. Take with a grain of salt...

## GQA

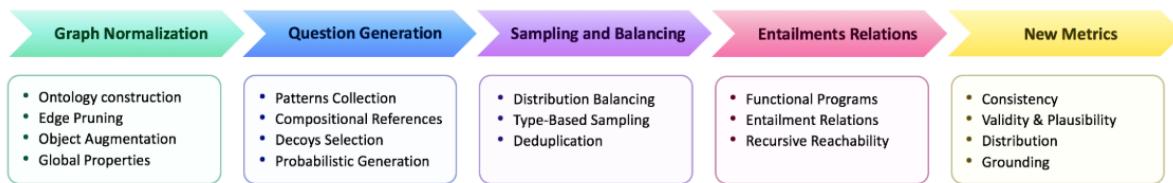
Table of Contents Local

Written by Brandon McKinzie

D. Hudson and C. Manning, “GQA: A New Dataset for Real-World Visual Reasoning and Compositional Question Answering” *Stanford University*, (May 2019).

**TL;DR:** New dataset for visual reasoning.

- 113K images and 22M questions. Leverages **Visual Genome** scene graph structures.
- New metrics: consistency, grounding, plausibility.
- Method for generating semantically varied questions.



**The GQA Dataset** (3). Images, questions, and answers are accompanied by matching semantic representations:

- **Images** are annotated with a dense **Scene Graph**.
- **Questions** are associated with a **functional program** which lists the series of reasoning steps needed to be performed to arrive at the answer.
- **Answers** are augmented with textual and visual justifications.

The four-step dataset construction pipeline is defined as follows:

1. **Scene Graph Normalization** (3.1). Visual Genome Scene Graph annotations<sup>96</sup>. The apply a significant amount of cleaning/normalization/rules to these.
2. **The Question Engine** (3.2). Generation harnesses (1) the normalized scene graphs and (2) structural patterns. Looks like sophisticated template/slot filling.
3. **Functional Representation and Entailment** (3.3). Questions can be represented by programs composed of atomic operations such as object selection, traversal along a relation edge, or an attribute verification. These are chained together to create challenging reasoning questions. A functional program can also **entail** another program<sup>97</sup>.

*We define direct entailment relations between the various functional programs and use these to recursively compute all the questions that can be entailed from a given source.*

<sup>96</sup>Covers 113K images from **COCO** and **Flickr**.

<sup>97</sup>For example, [knowing the answer to] *What color is the apple?* can determine entailment for *Is the apple red?*

4. **Sampling and Balancing** (3.4). Use the functional programs attached to each question to smooth out the answer distribution.

- Derive global and local labels for each question (using its functional program).
- Using the labels, partition the questions into *groups*.
- For each group, compute its **answer distribution**  $P$ .
- Use rejection-sampling to fit a smoother distribution  $Q$ .

*We iterate over the answers of the group in decreasing frequency order, and reweight  $P$ 's head up to the current iteration to make it more comparable to the tail size. While repeating this operation as we go through the answers, iteratively “moving” probability from the head into the tail, we also maintain minimum and maximum ratios between each pair of subsequent answers (sorted by frequency).*

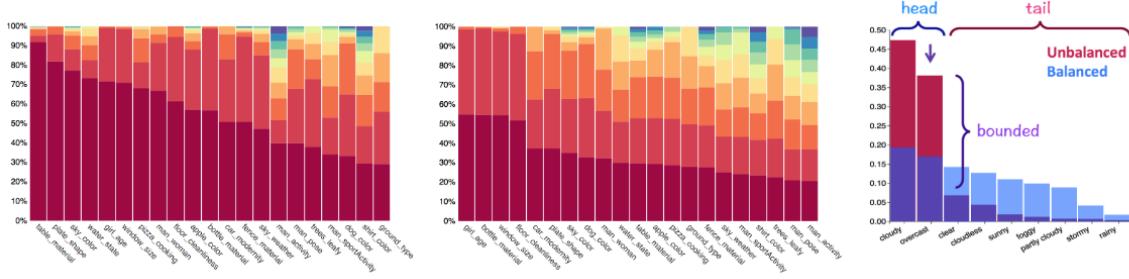


Figure 5: Visualization of the balancing process. The conditional answer distribution before (**left**) and after (**middle**) the balancing for a selection of question groups. We show the top 10 answers, where the column height corresponds to the relative frequency of each answer. We can see that on the left the distributions are heavily biased, while on the middle it is more uniform and with heavier tails, while intentionally retaining the original real-world tendencies up to a tunable degree. **Right:** An illustration of the balancing process.

**Analysis and Baseline Experiments** (4). Each question has two types:

- **Structural:** derived from final operation in question's functional program.
  - **Verify:** yes/no questions.
  - **Query:** open questions.
  - **Choose:** questions presenting two alternative choices.
  - **Logical:** involve logical inference.
  - **Compare:** comparisons between two or more objects.
- **Semantic:** refers to main subject of question.
  - **Object:** existence questions.
  - **Attribute:** consider properties/position of object.
  - **Category:** object identification within some class.
  - **Relation:** questions asking about subject/object of a described relation.
  - **Global:** overall properties of the scene (e.g. weather, place).

The authors also introduce the following metrics (4.4):

- **Consistency:** measures how consistent a model’s answers are wrt each other. For each question-answer pair (from dataset), define  $E_q \triangleq \{q_1, \dots, q_n\}$  as the set of *entailed questions*<sup>98</sup>. Then consistency can be computed as:

$$\frac{1}{|Q|} \sum_{q \in Q} \text{Acc}(E_q) \quad (531)$$

where  $Q$  is the set of questions for which the model’s answer  $\hat{a}$  was the correct answer  $a$ .

- **Validity:** fraction of answers that were at least within the question’s *scope*. For example, responding some color to a color question.
- **Plausibility:** measures whether the answer occurs at least once in relation with the question’s subject, across the whole dataset. For example, this would make “purple” and implausible (but valid) response for an apple’s color.
- **Distribution:** **TODO**

---

<sup>98</sup>Entailed questions are those for which the answer can be unambiguously inferred given  $(q, a)$ .

## A negative case analysis of visual grounding methods for VQA

Table of Contents Local

Written by Brandon McKinzie

R. Shrestha, K. Kafle, and C. Kanan, “A negative case analysis of visual grounding methods for VQA” *Rochester Institute of Technology*, (April 2020).

**Existing VQA Methods (3).**

$$[\text{task}] \quad P(a | \mathcal{Q}, \mathcal{I}) = f_{VQA}(v, \mathcal{Q}) \quad (532)$$

$$[\text{sensitivity}] \quad S(a, v_i) \triangleq (\nabla_{v_i} P(a | \mathcal{I}, \mathcal{Q}))^T \mathbf{1} \quad (533)$$

Existing approaches develop training objectives using  $S$  in order to improve visual grounding<sup>99</sup>. Here, the authors investigate: *is [the improved performance] actually due to better visual grounding?*

**Why Did the Performance Improve? (4).** Main experiments/findings:

- Training on *irrelevant* visual regions yields similar results as relevant regions<sup>100</sup>
- Training *random* visual cues also performs similarly, even if the random cues are re-randomly sampled each epoch!
- Use Welch’s t-test to confirm that differences in performance of each method is statistically insignificant.
- Roughly 90% overlap of predictions made on test set, confirming that all approaches learn similar things.
- HINT and SCR help forget linguistic priors, which is beneficial for VQA-CPv2 but not for VQAv2.
- Although HINT/SCR improve visual grounding, only  $\approx 25\%$  of their correct answers were most sensitive to a region in the top-3 most relevant ground truth regions.

---

<sup>99</sup>The two works being investigated are HINT and SCR.

<sup>100</sup>Relevant regions are roughly determined as those having high human-based importance scores  $S_h$ . Here, the authors define irrelevant as essentially the opposite:  $1 - S_h$ .

**Embarrassingly Simple Regularizer** (5). They show that simply degrading training performance by introducing a loss that has only zero-vector labels achieves near SOTA on test sets<sup>101</sup>.

$$L := \text{BCE}(P(\mathcal{A}), \mathcal{A}_{gt}) + \lambda \text{BCE}(P(\mathcal{A}), \mathbf{0}) \quad (534) \quad \text{Authors set } \lambda=1$$

---

<sup>101</sup>Ok this is actually hilarious.

## Bottom-Up and Top-Down Attention for Image Captioning and Visual Question Answering

Table of Contents Local

Written by Brandon McKinzie

Anderson et al., “Bottom-Up and Top-Down Attention for Image Captioning and Visual Question Answering”  
*Australian National University, JD AI, Microsoft*, (March 2018).

**Bottom-Up Attention Model** (3.1). Use Faster R-CNN to obtain bounding boxes with object labels. The  $k$  input image features,  $V = \{\mathbf{v}_1, \dots, \mathbf{v}_k\}$ ,  $\mathbf{v}_i \in \mathbb{R}^D$ , are obtained for each of the  $k$  regions output by R-CNN, then mean-pooling over the convolutional feature for the region.

**Captioning Model** (3.2). Consists of two LSTM layers:

1. Top-down visual attention model.

- (a) Compute mean-pooled image feature  $\bar{\mathbf{v}} = \frac{1}{k} \sum_i \mathbf{v}_i$ .
- (b) Lookup word embedding for input word<sup>102</sup> at timestep t and denote as  $\tilde{\mathbf{x}}_t^1$ .
- (c) Feed to LSTM:

$$\mathbf{x}_t^1 = [\mathbf{h}_{t-1}^2, \bar{\mathbf{v}}, \tilde{\mathbf{x}}_t^1] \quad (535)$$

where  $\mathbf{h}_{t-1}^2$  is the output of the second LSTM (below) at previous timestep.

- (d) Compute normalized attention weights  $\alpha_{i,t}$ :

$$a_{i,t} = \mathbf{w}_a^T \tanh(\mathbf{W}_{va} \mathbf{v}_i + \mathbf{W}_{ha} \mathbf{h}_t^1) \quad (536) \quad \mathbf{W}_{va} \in \mathbb{R}^{H \times V}$$

$$\boldsymbol{\alpha}_t = \text{softmax}(\mathbf{a}_t) \quad (537)$$

$$\hat{\mathbf{v}}_t = \sum_{i=1}^K \alpha_{i,t} \mathbf{v}_i \quad (538)$$

where  $\mathbf{w}_a$  is learned<sup>103</sup>.

2. Language model.

- (a) Input is the attention vector and output from LSTM 1:  $\mathbf{x}_t^2 = [\hat{\mathbf{v}}_t, \mathbf{h}_t^1]$
- (b) As usual, predictive next-token distribution is defined as  $p(y_t | y_{1:t-1}) = \text{softmax}(\mathbf{W}_p \mathbf{h}_t^2 + \mathbf{b}_p)$

<sup>102</sup>The most recently generated word from the 2nd LSTM (LM).

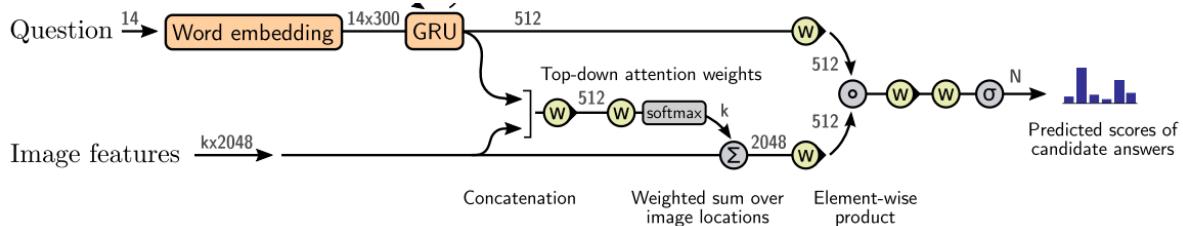
<sup>103</sup>Basically like an importance vector.

### VQA Model (3.3).

1. Encode question as hidden state  $\mathbf{q}$  of a GRU.
2. Generate unnormalized attention weights over  $k$  image features  $\mathbf{v}_i$ :

$$a_i = \mathbf{w}_a^T f_a ([\mathbf{v}_i, \mathbf{q}]) \quad (539)$$

and normalize to obtain  $\hat{\mathbf{v}}$  as defined earlier.



## ViLBERT

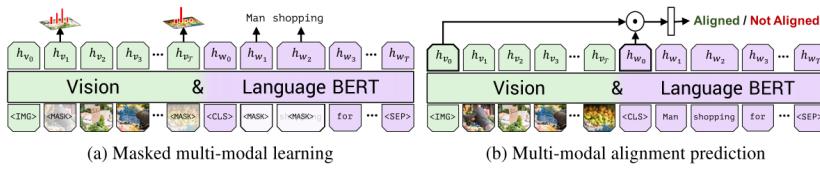
Table of Contents Local

Written by Brandon McKinzie

J. Lu, D. Batra, D. Parikh, and S. Lee, “ViLBERT: Pretraining Task-Agnostic Visiolinguistic Representations for Vision-and-Language Tasks” *GA Tech, FAIR, Oregon State*, (August 2019).

[Online Demo]

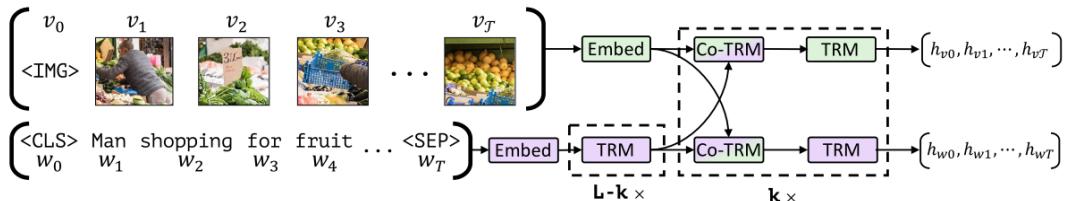
**TL;DR:** A model for learning task-agnostic joint representations of image context and natural language. Extends BERT with multi-modal two-stream attention. Pretrained on Conceptual Captions. Evaluated on VQA, VCR, referring expressions, caption-based image retrieval.



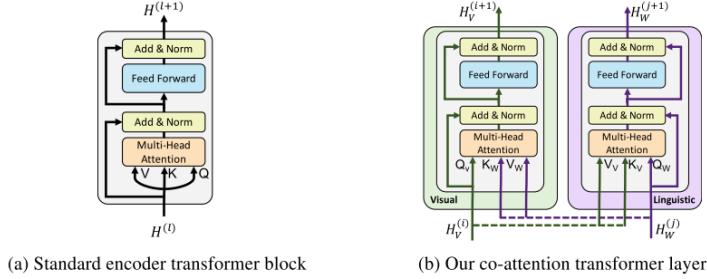
**ViLBERT (2.2).** The ViLBERT architecture consists of two parallel BERT-style models.

**TODO: need to understand in exact detail things like how the image regions are “embedded”, how the two-stream attn is impl, ,etc**

- **Inputs:** a set of region features  $\{v_1, \dots, v_T\}$  from image  $I$  and text input  $\{w_1, \dots, w_T\}$ .
- **Outputs:** representations  $\{h_{v_1}, \dots, h_{v_T}\}$  and  $\{h_{w_1}, \dots, h_{w_T}\}$ .



where Co-TRM is a modified attention block to obtain multi-modal attention:



**Vision-and-Language Transfer Tasks** (3.2). How to adapt the base ViLBERT architecture for various downstream tasks.

- **VQA.** Learn 2-layer MLP on top of the element-wise product of image and text representations  $h_{IMG}$  and  $h_{CLS}$ , mapping to 3,129 possible answers.
- **VCR.** Concatenate question and each possible response to form four different text inputs and pass each through ViLBERT along with the image. We learn a linear layer on top of the post-elementwise product representation to predict a score for each pair. The final prediction is a softmax over these four scores.
- **RefCOCO+.** Pass the final representation  $h_{v_i}$  for each image region  $i$  into a learned linear layer to predict a matching score. Label each proposal box by computing the IoU with the ground truth box and thresholding at 0.5. Train with binary cross-entropy.
- **Caption-Based Image Retrieval.** Train in a 4-way multiple-choice setting by randomly sampling three distractors for each image-caption pair – substituting a random caption, a random image, or a hard negative from among the 100 nearest neighbors of the target image. Compute the alignment score (as in alignment prediction pretraining) for each and apply a softmax. Train with cross-entropy loss to select the true image-caption pair.

---

#### 6.14.1 12-IN-1: MULTI-TASK VISION AND LANGUAGE REPRESENTATION LEARNING

Lu et al., “12-in-1: Multi-Task Vision and Language Representation Learning” *GA Tech, FAIR, Oregon State*, (April 2020).

**TL;DR:** Follow-up paper by authors of ViLBERT. Whereas before they seemingly finetune all parameters of the model for each task separately, they now present a multi-task training setup that allows a *single* ViLBERT model to achieve SOTA on 12 tasks.

## Language-Conditioned Graph Networks for Relational Reasoning

Table of Contents Local

Written by Brandon McKinzie

Hu et al., “Language-Conditioned Graph Networks for Relational Reasoning” *UC Berkeley, Boston University*, (August 2019).

**Context-Aware Scene Representation** (3.1). Let  $N$  denote number of entities in image  $I$ , and let  $x_i^{loc}$  ( $1 \leq i \leq N$ ) denote the feature representation of entity  $i$ .

- **Goal:** output a context-aware representation  $x_i^{out}$  for each entity  $i$  conditioned on the textual input  $Q$ .
- **Approach:** iterative message passing over  $T$  iterations on fully-connected graph of entity nodes. In addition to [the fixed]  $x_i^{loc}$ , each node has a *context feature*  $x_{i,t}^{ctx}$  that is updated during each iteration  $t$ .

**Textual command extraction.** Let  $Q$  denote question containing  $S$  words. At each iteration  $t$ , the **textual command**  $\mathbf{c}_t$  is obtained as follows:

$$\{\mathbf{h}_1, \dots, \mathbf{h}_S\} = \text{BiLSTM}(Q) \quad (540)$$

$$\mathbf{q} = [\mathbf{h}_1; \mathbf{h}_S] \quad (541)$$

$$\mathbf{c}_t = \sum_{s=1}^S \alpha_{t,s} \mathbf{h}_s \quad \text{where} \quad (542)$$

$$\alpha_{t,s} = \text{Softmax}_s \left( \mathbf{W}_1 \left( \mathbf{h}_s \odot \left( \mathbf{W}_2^{(t)} \text{ReLU}(\mathbf{W}_3 \mathbf{q}) \right) \right) \right) \quad (543)$$

**Language-Conditioned Message Passing.**

1. Compute “joint representation” for each node:

$$\tilde{x}_{i,t} = \left[ x_i^{loc}; x_{i,t-1}^{ctx}; \mathbf{W}_4 x_i^{loc} \odot \mathbf{W}_5 x_{i,t-1}^{ctx} \right] \quad (544)$$

2. Compute directed connection weights from node  $j$  (sender) to node  $i$  (receiver).

$$w_{j \rightarrow i}^{(t)} = \text{Softmax}_j \left( (\mathbf{W}_6 \tilde{x}_{i,t})^T (\mathbf{W}_7 \tilde{x}_{j,t} \odot \mathbf{W}_8 \mathbf{c}_t) \right) \quad (545)$$

3. Each node sends message to each other node.

$$m_{j \rightarrow i}^{(t)} = w_{j \rightarrow i}^{(t)} \cdot (\mathbf{W}_9 \tilde{x}_{j,t} \odot \mathbf{W}_{10} \mathbf{c}_t) \quad (546)$$

$$x_{i,t}^{ctx} = \mathbf{W}_{11} \left[ x_{i,t-1}^{ctx}; \sum_{j=1}^N m_{j \rightarrow i}^{(t)} \right] \quad (547)$$

4. Final representation after  $T$  iterations:

$$x_i^{out} = \mathbf{W}_{12} [x_i^{loc}; x_{i,T}^{ctx}] \quad (548)$$

**Application to VQA and REF** (3.2). They define a **single-hop answer classifier** for VQA as follows:

$$\beta_i = \text{Softmax}_i \left( \mathbf{W}_{13} (x_i^{out} \odot \mathbf{W}_{14} \mathbf{q}) \right) \quad (549)$$

$$y = \mathbf{W}_{15} \text{ReLU} \left( \mathbf{W}_{16} \left[ \sum_{i=1}^N \beta_i x_i^{out}; \mathbf{q} \right] \right) \quad (550)$$

where  $y$  is the scores (logits) over all possible answers.

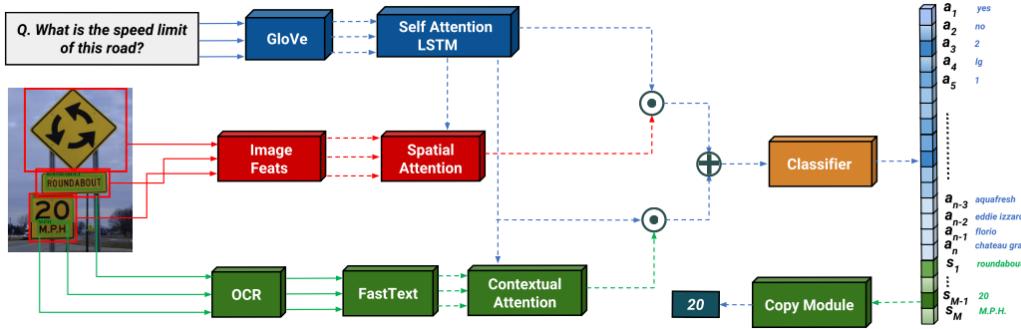
## Towards VQA Models That Can Read

Table of Contents Local

Written by Brandon McKinzie

Singh et al., "Towards VQA Models That Can Read" *FAIR, GA Tech*, (May 2019).

**TL;DR:** New **TextVQA** dataset that requires reasoning about images containing text. Introduce **LoRRA** model that can better handle OCR in the setting of VQA.



Denote the  $M$  words extracted from the OCR model over the image as  $s = \{s_1, \dots, s_M\}$ .

$$f_{VQA}(v, q) = f_{comb}^{(1)}(f_A^{(1)}(f_I(v), f_Q(q)), f_Q(q)) \quad (551)$$

$$f_{OCR}(s, q) = f_{comb}^{(2)}(f_A^{(2)}(f_O(s), f_Q(q)), f_Q(q)) \quad (552)$$

$$f_{LoRRA}(v, s, q) = f_{MLP}([f_{VQA}(v, q); f_{OCR}(s, q)]) \quad (553)$$

## Accuracy vs. Complexity: A Trade-off in VQA Models

Table of Contents Local

Written by Brandon McKinzie

Farazi et al., “Accuracy vs. Complexity: A Trade-off in Visual Question Answering Models” *ANU, CSIRO, IAAI*, (Jan 2020).

**Feature Extraction Meta-Architecture** (2.1). Consists of both a *visual* and *language* feature extraction.

- **Visual Feature.**

$$\mathbf{v} = \text{CNN}(\mathbf{I}) = \{\mathbf{v}_i \mid i \in [1, G]\} \quad (554) \quad \mathbf{v}_i \in \mathbb{R}^{d_v}$$

where  $G$  denote the total number of image locations in a grid. The extracted features can be categorized into three main types:

1. **Image Level (IL):** represent features of the whole image ( $G=1$ ).
  2. **Spatial Grid (SG):** spatial grid where each of the  $G$  feature maps corresponds to a uniform grid location on the input image.
  3. **Bottom-Up (BU):** each feature map corresponds to an object proposal region ( $G=N$ ).
- **Language Feature:** question is tokenized into words and fed through and LSTM or transformer to obtain contextualized representations of each word. These are aggregated in some way to obtain question vector  $\mathbf{q}$ .

**Fusion Model Meta-Architecture** (2.2). Jointly embed into a common space, learning multimodel embedding function  $\Phi$ :

$$\mathbf{z} = \Phi(\mathbf{q}, \mathbf{v}) \quad (555)$$

In the most general case, we can compute  $\mathbf{z} = \mathcal{W}[\mathbf{v} \otimes \mathbf{q}]$ , but this is parameter overkill. The authors outline a variety of instances:

- **Linear:**  $\mathbf{z} = \text{FC}(\text{FC}(\mathbf{v}) + \text{FC}(\mathbf{q}))$
- **C-MLP:**  $\mathbf{z} = \text{MLP}([\mathbf{v}; \mathbf{q}])$

**Attention-Based Meta-Architecture** (2.3). Question-specific attention over image regions (if SG) or object proposals (if BU).

**Experiments and Results** (5). Main insights:

- Set visual feature dimension closer to language feature embedding [by modifying the language dim size].
- PolyNet with smaller grid sizes performs surprisingly well.
- Resnet152 has largest performance boost when going from IL to SG.
- SeNet154 perfroms better on datasets with less language bias (e.g. VQA-CP instead of VQAv2).
- Using Bottom-Up features provides a consistent accuracy gain.
- VQA models are less sensitive to change in batch size.

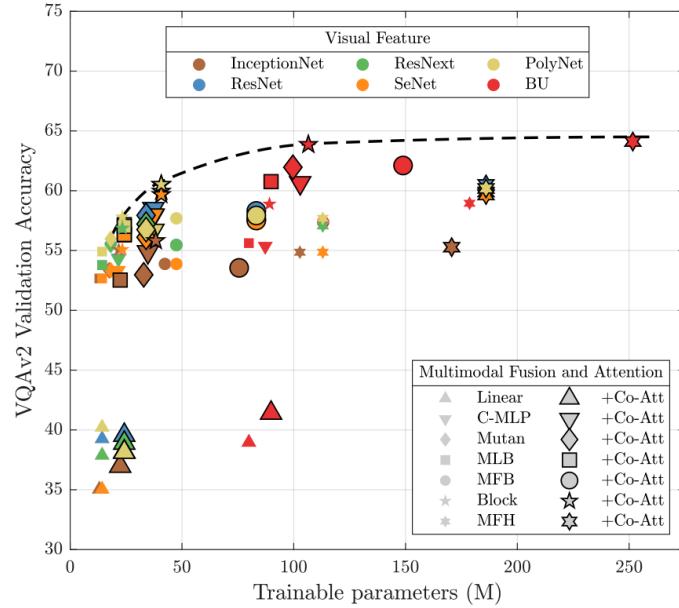


Figure 6: The trade-off between VQA v2 validation accuracy vs. the number of trainable parameters.

## Iterative Answer Prediction with Multimodal Transformers for TextVQA

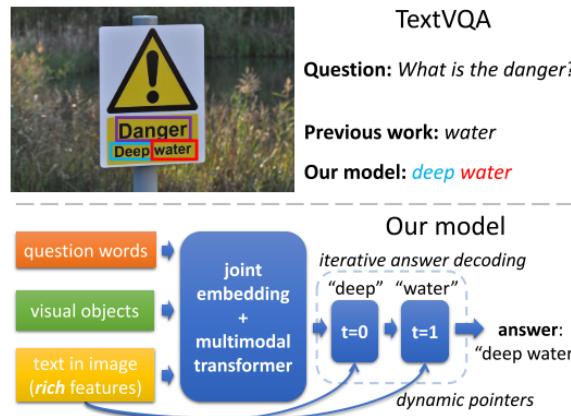
Table of Contents Local

Written by Brandon McKinzie

R. Hu, A. Singh, T. Darrell, and M. Rohrbach, “Iterative Answer Prediction with Pointer-Augmented Multi-modal Transformers for TextVQA” *FAIR, UC Berkeley*, (Mar 2020).

**Introduction** (1). Current approaches on the TextVQA task...

- “rely on custom pairwise multimodal fusion mechanisms between two modalities, which limit[s] the types of possible interactions between modalities” (**how is this limiting?**)
- Treat VQA as multi-label classification, which obviously limits the output space of the model.
- Make poor use of OCR by just concatenating detected token embeddings together (ignoring spacial location, color, etc.).



**Common Embedding Space** (3.1).

- **Question words.** Sequence of  $K$  question words are fed through pretrained BERT.<sup>104</sup>
- **Detected objects.** Obtain set of  $M$  visual objects via **Faster R-CNN**<sup>105</sup>.

$$x_m^{obj} = LN \left( \mathbf{W}_1 x_m^{fr} \right) + LN \left( \mathbf{W}_2 x_m^b \right) \quad (556)$$

- **Embedding of OCR tokens.** Obtain a set of  $N$  OCR tokens (words) through external OCR system. For each token  $n$ :
  1. FastText vector  $x_n^{ft}$
  2. Appearance feature  $x_n^{fr}$  from same Faster R-CNN detector.

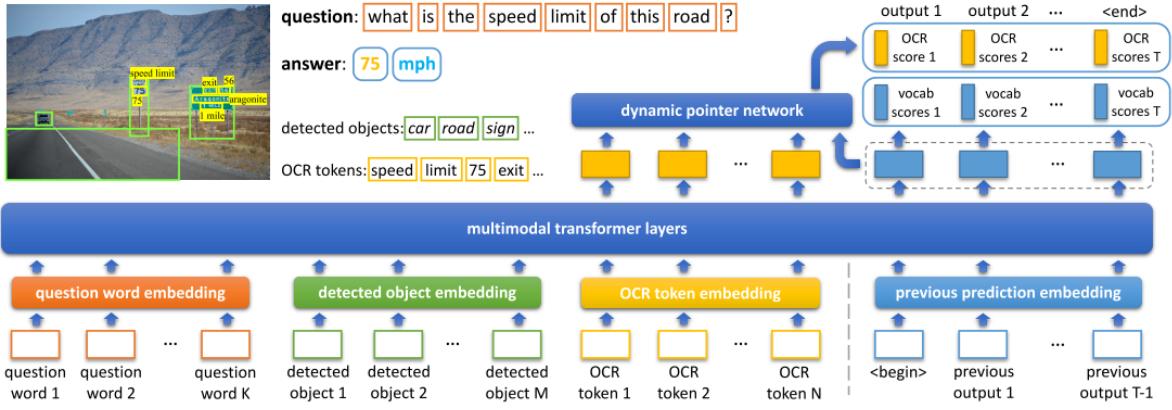
<sup>104</sup>During training, the BERT parameters are fine-tuned using the question answering loss.

<sup>105</sup>They fine-tune the last layer of Faster-RCNN during training.

3. PHOC vector  $x_n^p$ .
4. 4-dimensional location feature  $x_n^b$  based on OCR token's relative bounding box coordinates.

$$x_n^{ocr} = LN \left( \mathbf{W}_3 x_n^{ft} + \mathbf{W}_4 x_n^{fr} + \mathbf{W}_5 x_n^p \right) + LN \left( \mathbf{W}_6 x_n^b \right) \quad (557)$$

**Multimodel fusion and iterative answer prediction with pointer-augmented transformers** (3.2). Apply  $L$  transformer layers with hidden dim  $d$  over the list of all  $K + M + N$  entities from  $\{x_k^{ques}\}$ ,  $\{x_m^{obj}\}$ , and  $\{x_n^{ocr}\}$  (prev section).



Let  $\{z_1^{ocr}, \dots, z_N^{ocr}\}$  be the  $d$ -dimensional transformer outputs of the  $N$  OCR tokens in the image. As shown in the illustration above, they initialize decoding by feeding in special <begin> as  $x_1^{dec}$ , with corresponding output  $z_1^{dec}$ , then decode until  $z_t^{dec}$  corresponds to <end>.

$$y_{t,i}^{voc} = (w_i^{voc})^T z_t^{dec} + b_i^{voc} \quad (558)$$

$$y_{t,n}^{ocr} = (\mathbf{W}^{ocr} z_n^{ocr} + b^{ocr})^T (\mathbf{W}^{dec} z_t^{dec} + b^{dec}) \quad (559)$$

They predict by taking an argmax of  $y_t^{all} = [y_t^{voc}; y_t^{ocr}]$ . Train via ‘‘multi-label sigmoid loss’’ (instead of softmax loss) over scores  $y_t^{all}$ .

## Q & A

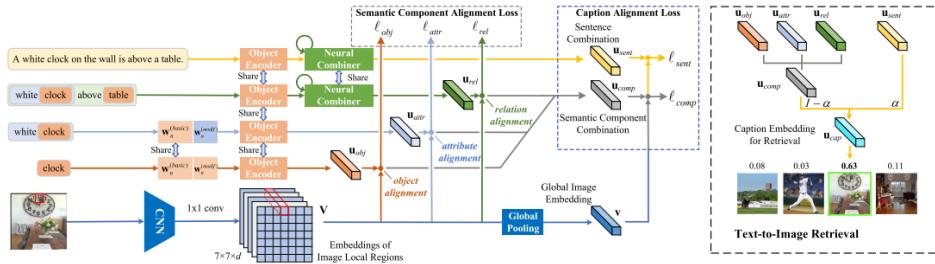
- **Q:** What happens if an image without scene text is fed as input?
  - **A:** From the code at least, it would just become a vector of zeros (padding).
- **Q:** What is PHOC?
  - **A:** Pyramidal Histogram of Characters. It's from the paper "Word Spotting and Recognition with Embedded Attributes." From the paper: *The proposed approach embeds text strings into a d-dimensional binary space. In a nutshell, this embedding – which we dubbed pyramidal histogram of characters or PHOC – encodes if a particular character appears in a particular spatial region of the string*
- **Q:** What is the backbone image model being used?
  - **A:** Faster R-CNN pretrained on Visual Genome. Based on the MMF code, it is **detectron**.
- **Q:** What is the detectron model exactly?
  - **A:**
- **Q:** What is the training dataset?
  - **A:** Best guess is textVQA because their example snippet here uses it.
- **Q:** The phrase "single-hop" attention seems to just literally mean "attention" – what would be an example of "multi-hop" attention?
  - **A:**

Quick summaries of various papers presented at CVPR 2020.

### 6.19.1 UNIFIED VISUAL-SEMANTIC EMBEDDINGS: BRIDGING VISION AND LANGUAGE WITH STRUCTURED MEANING REPRESENTATIONS

Wu et al., “Unified Visual-Semantic Embeddings: Bridging Vision and Language with Structured Meaning Representations” *Fudan University*, (2020).

*Given an image-caption pair, we first parse the caption into a structured meaning representation, composed by a set of semantic components: object nouns, prenominal modifiers, and relational dependencies. We encode different types of semantic components with type-specific encoders. A caption encoder combines the embedding of the semantic components into a **caption semantic embedding**. Jointly, we encode images with a convolutional neural network (CNN) into the same, unified VSE space. The distance between the image embedding and the sentential embedding measures the semantic similarity between the image and the caption. We employ a multi-task learning approach for the joint learning of embeddings for semantic components (as the “basis” of the VSE space) as well as the caption encoder (as the combiner of semantic components).*



They align modalities via a bidirectional margin-based ranking loss. They employ a semantic parser on the image captions.

---

### 6.19.2 COUNTERFACTUAL SAMPLES SYNTHESIZING FOR ROBUST VISUAL QUESTION ANSWERING

---

Chen et al., “Counterfactual Samples Synthesizing for Robust Visual Question Answering” *Zhejiang University*, (Mar 2020).

*We propose a model-agnostic Counterfactual Samples Synthesizing (CSS) training scheme. The CSS generates numerous counterfactual training samples by masking critical objects in images or words in questions, and assigning different ground-truth answers. After training with the complementary samples (i.e., the original and generated samples), the VQA models are forced to focus on all critical objects and words, which significantly improves both visual-explainable and question-sensitive abilities.*

**V-CSS** (4.2.1). Algorithm below consists of four main steps:

1. Initial objects selection (IO\_SEL).
2. Object local contributions calculation.
3. Critical objects selection (CO\_SEL).
4. Dynamic answer assigning (DS\_ASS).

---

#### Algorithm 2 Counterfactual Samples Synthesizing

---

```

1: function  $\mathcal{CSS}(I, Q, a)$ 
2:    $\mathbf{V}, \mathbf{Q}, P_{vqa}(\mathbf{a}) \leftarrow \mathcal{VQA}(I, Q, a, \text{True})$ 
3:    $cond \sim U[0, 1]$ 
4:   if  $cond \geq \delta$  then            $\triangleright$  execute V-CSS
5:      $\mathcal{I} \leftarrow \text{IO\_SEL}(I, Q)$ 
6:      $s(a, \mathbf{v}_i) \leftarrow \mathcal{S}(P_{vqa}(a), \mathbf{v}_i)$ 
7:      $I^+, I^- \leftarrow \text{CO\_SEL}(\mathcal{I}, \{s(a, \mathbf{v}_i)\})$ 
8:      $a^- \leftarrow \text{DA\_ASS}(I^+, Q, \mathcal{VQA}, a)$ 
9:      $\mathcal{VQA}(I^-, Q, a^-, \text{False})$ 
10:    else                       $\triangleright$  execute Q-CSS
11:       $s(a, \mathbf{w}_i) \leftarrow \mathcal{S}(P_{vqa}(a), \mathbf{w}_i)$ 
12:       $Q^+, Q^- \leftarrow \text{CW\_SEL}(\{s(a, \mathbf{w}_i)\})$ 
13:       $a^- \leftarrow \text{DA\_ASS}(I, Q^+, \mathcal{VQA}, a)$ 
14:       $\mathcal{VQA}(I, Q^-, a^-, \text{False})$ 
15:    end if
16: end function

```

---

---

### 6.19.3 WHAT MAKES TRAINING MULTI-MODAL CLASSIFICATION NETWORKS HARD?

---

W. Wang, D. Tran, and M. Feiszli, “What Makes Training Multi-modal Classification Networks Hard?” *FAIR*, (June 2020).

# GENERATIVE MODELS

## CONTENTS

|     |                                                              |     |
|-----|--------------------------------------------------------------|-----|
| 7.1 | WaveNet . . . . .                                            | 193 |
| 7.2 | Latent Dirichlet Allocation . . . . .                        | 197 |
| 7.3 | Generative Adversarial Nets . . . . .                        | 200 |
| 7.4 | Modular Generative Adversarial Networks . . . . .            | 203 |
| 7.5 | Transfer Learning from Speaker Verification to TTS . . . . . | 205 |
| 7.6 | Tacotron 2 . . . . .                                         | 206 |
| 7.7 | Glow . . . . .                                               | 208 |
| 7.8 | WaveGlow . . . . .                                           | 209 |

## 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,

$$p(x) = \prod_{t=1}^T p(x_t | x_1, \dots, x_{t-1}) \quad (560)$$

which are modeled by a stack of convolutional layers (no pooling). 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>106</sup>.

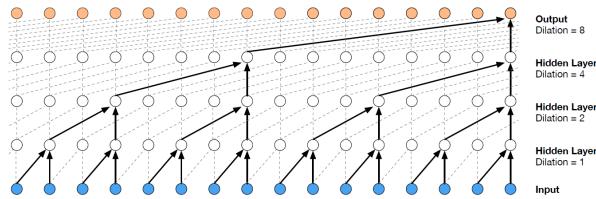


Figure 3: Visualization of a stack of *dilated* causal convolutional layers.

**Dilated Convolution**

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.

<sup>106</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”

**Softmax distributions.** To deal with the fact that there are  $2^{16}$  possible values, first apply a  $\mu$ -law companding transformation<sup>107</sup> to data, and then quantize it to 256 possible values:

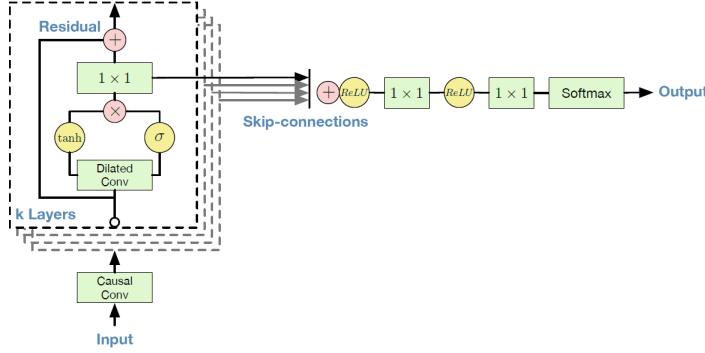
$$f(x_t) = \text{sign}(x_t) \frac{\ln(1 + \mu|x_t|)}{\ln(1 + \mu)} \quad -1 < x_t < 1 \text{ and } \mu = 255 \quad (561)$$

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

**Gated Activation and Residual/Skip Connections.** Use the same gated activation unit as PixelCNN:

$$z = \tanh(W_{f,k} * x) \odot \sigma(W_{g,k} * x) \quad (562)$$

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.



The 1x1 blocks are 1x1 convolutions (i.e. position-wise dense layers).

---

<sup>107</sup>In telecommunication and signal processing **companding** (occasionally called compansion) is a method of mitigating the detrimental effects of a channel with limited dynamic range.

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

$$p(\mathbf{x} \mid \mathbf{h}) = \prod_{t=1}^T p(x_t \mid x_1, \dots, x_{t-1}, h) \quad (563)$$

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

$$z = \tanh(W_{f,k} * \mathbf{x} + V_{f,k}^T \mathbf{h}) \odot \sigma(W_{g,k} * \mathbf{x} + V_{g,k}^T \mathbf{h}) \quad (564)$$

→ **Local conditioning.** Have a second time-series  $h_t$ . They first transform this  $h_t$  using a **transposed CNN (learned upsampling)** that maps it to a new time-series  $\mathbf{y} = f(\mathbf{h})$  w/same resolution as  $\mathbf{x}$ .

$$z = \tanh(W_{f,k} * \mathbf{x} + V_{f,k} * \mathbf{y}) \odot \sigma(W_{g,k} * \mathbf{x} + V_{g,k}^T \mathbf{y}) \quad (565)$$

## Experiments.

- **Multi-Speaker Speech Generation.** Dataset: multi-speaker corpus of 44 hours of data from 109 different speakers<sup>108</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.”

---

<sup>108</sup>Speakers encoded as ID in form of a one-hot vector

**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.”

## 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>109</sup><sup>110</sup>:

$$\text{Dir}(\boldsymbol{x}; \boldsymbol{\alpha}) = \frac{1}{B(\boldsymbol{\alpha})} \prod_{i=1}^K \boldsymbol{x}_i^{\boldsymbol{\alpha}_i - 1} \quad \text{where} \quad B(\boldsymbol{\alpha}) = \frac{\prod_{i=1}^K \Gamma(\boldsymbol{\alpha}_i)}{\Gamma(\sum_{i=1}^K \boldsymbol{\alpha}_i)} \quad (566) \quad \sum_{i=1}^K \boldsymbol{x}_i = 1 \\ (\forall i \in [1, K]) : \boldsymbol{x}_i \geq 0$$

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>111</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>109</sup>Recall that for positive integers  $n$ ,  $\Gamma(n) = (n - 1)!$ .

<sup>110</sup>The Dirichlet distribution is **conjugate** to the multinomial distribution. TODO: Review how to interpret this.

<sup>111</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. 566),  $\text{Dir}(\alpha)$ .  $k$  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>112</sup> The parameter  $\beta$  gives the distribution of words given a topic:

$$\beta_{ij} = \Pr [w_j | z_i] \quad (567)$$

In other words, we really sample  $w_n \sim \beta_{i,:}$

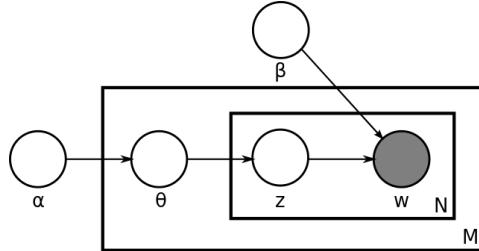
The defining equations for LDA are thus:

$$\Pr [\theta, \mathbf{z}, \mathbf{w} | \alpha, \beta] = \Pr [\theta | \alpha] \prod_{n=1}^N \Pr [z_n | \theta] \Pr [w_n | z_n, \beta] \quad (568)$$

$$\Pr [\mathbf{w} | \alpha, \beta] = \int \Pr [\theta' | \alpha] \left( \prod_{n=1}^N \sum_{z'_n} \Pr [z'_n | \theta'] \Pr [w_n | z'_n, \beta] \right) d\theta' \quad (569)$$

$$\Pr [\mathcal{D} = \{\mathbf{w}^{(1)}, \dots, \mathbf{w}^{(M)}\} | \alpha, \beta] = \prod_{d=1}^M \Pr [\mathbf{w}^{(d)} | \alpha, \beta] \quad (570)$$

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>113</sup>.

---

<sup>112</sup>**TODO:** interpret meaning of the multinomial distributions here. Seems a bit different than standard interp...

<sup>113</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  $w$ :

$$\Pr [\theta, z | w, \alpha, \beta] = \frac{\Pr [\theta, z, w | \alpha, \beta]}{\Pr [w | \alpha, \beta]} \quad (571)$$

Computing the denominator exactly is intractable. Common approximate inference algorithms for LDA include Laplace approximation, variational approximation, and Markov Chain Monte Carlo.

## 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>114</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:

$$\min_G \max_D V(\mathbf{D}, \mathbf{G}) = \mathbb{E}_{\mathbf{x} \sim p_{data}} [\log \mathbf{D}(\mathbf{x})] + \mathbb{E}_{\mathbf{z} \sim p_n} [\log (1 - \mathbf{D}(\mathbf{G}(\mathbf{z})))] \quad (572)$$

**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>114</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

$$D_G^*(\mathbf{x}) = \frac{p_{data}(\mathbf{x})}{p_{data}(\mathbf{x}) + p_g(\mathbf{x})} \quad (573)$$

**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

$$\mathbb{E}[G(\mathbf{z})] = \mathbb{E}_{\mathbf{z} \sim p_n}[G(\mathbf{z})] \quad (574)$$

$$= \int_{\mathbf{z}} p_n(\mathbf{z}) G(\mathbf{z}) d\mathbf{z} \quad (575)$$

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  $p_g(\mathbf{x})$ .

$$\mathbb{E}[G] = \mathbb{E}_{\mathbf{x} \sim p_g}[\mathbf{x}] \quad (576)$$

$$= \int_{\mathbf{x}} p_g(\mathbf{x}) \mathbf{x} d\mathbf{x} \quad (577)$$

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  $V(G, D)$ :

$$V(G, D) = \mathbb{E}_{\mathbf{x} \sim p_{data}} [\log D(\mathbf{x})] + \mathbb{E}_{\mathbf{z} \sim p_n} [\log (1 - D(G(\mathbf{z})))] \quad (578)$$

$$= \mathbb{E}_{\mathbf{x} \sim p_{data}} [\log D(\mathbf{x})] + \mathbb{E}_{\mathbf{x} \sim p_g} [\log (1 - D(\mathbf{x}))] \quad (579)$$

$$= \int_{\mathbf{x}} [p_{data}(\mathbf{x}) \log(D(\mathbf{x})) + p_g(\mathbf{x}) \log(1 - D(\mathbf{x}))] d\mathbf{x} \quad (580)$$

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  $C(G) \triangleq V(G, D_G^*)$ :

$$C(G) = \mathbb{E}_{\mathbf{x} \sim p_{data}} \left[ \log \frac{p_{data}(\mathbf{x})}{p_{data}(\mathbf{x}) + p_g(\mathbf{x})} \right] + \mathbb{E}_{\mathbf{x} \sim p_g} \left[ \log \frac{p_g(\mathbf{x})}{p_{data}(\mathbf{x}) + p_g(\mathbf{x})} \right] \quad (581)$$

### 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 581, do some substitutions, and find that

$$C(G) = 2 \cdot JSD(p_{data} || p_g) - \log 4 \quad (582)$$

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

## Modular Generative Adversarial Networks

Table of Contents Local

Written by Brandon McKinzie

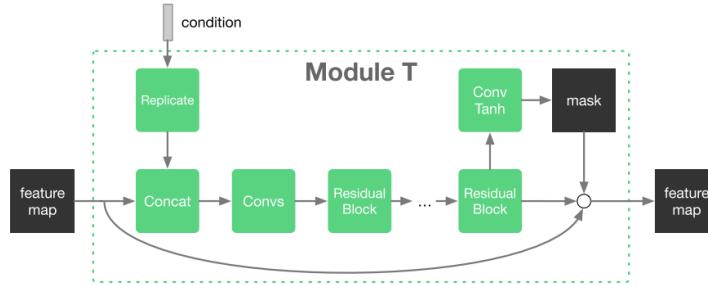
Zhao et al., “Modular Generative Adversarial Networks” *UBC, Tencent AI Lab* (April 2018).

**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  $E$** :  $x \mapsto E(x)$
  - [generation] **generator  $G$** :  $(z, a_0) \mapsto G(z, a_0)$  where  $z$  is random noise and  $a_0$  is a condition vector representing auxiliary information.
2. **transformer(s)  $T_i$**  :  $E(x) \mapsto T_i(E(x), a_i)$ . Modifies repr of attrib  $a_i$  in the FM.
3. **reconstructor  $R$**  :  $(T_i, T_j, \dots) \mapsto y$ . Reconstructs image from an intermediate FM.
4. **discriminator  $D_i$** :  $R \mapsto \{0, 1\} \times \text{Val}(a_i)$ . 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).**

$$\mathcal{L}_D(\mathbf{D}) = -\sum_{i=1}^n \mathcal{L}_{adv_i} + \lambda_{cls} \sum_{i=1}^n \mathcal{L}_{cls_i}^r \quad (583)$$

$$\mathcal{L}_G(\mathbf{E}, \mathbf{T}, \mathbf{R}) = \sum_{i=1}^n \mathcal{L}_{adv_i} + \lambda_{cls} \sum_{i=1}^n \mathcal{L}_{cls_i}^f + \lambda_{cyc} \left( \mathcal{L}_{cyc}^{\mathbf{ER}} + \sum_{i=1}^n \mathcal{L}_{cyc}^{\mathbf{T}_i} \right) \quad (584)$$

$$\begin{aligned} \mathcal{L}_{adv_i}(\mathbf{E}, \mathbf{T}_i, \mathbf{R}, \mathbf{D}_i) &= \mathbb{E}_{y \sim p_{data}(y)} [\log \mathbf{D}_i(y)] + \\ &\quad \mathbb{E}_{x \sim p_{data}(x)} [\log (1 - \mathbf{D}_i(\mathbf{R}(\mathbf{T}_i(\mathbf{E}(x)))))] \end{aligned} \quad (585)$$

$$\mathcal{L}_{cls_i}^r = -\mathbb{E}_{x, c_i} [\log \mathbf{D}_{cls_i}(c_i \mid x)] \quad (586)$$

$$\mathcal{L}_{cls_i}^f = -\mathbb{E}_{x, c_i} [\log \mathbf{D}_{cls_i}(c_i \mid \mathbf{R}(\mathbf{E}(x)))] \quad (587)$$

$$\mathcal{L}_{cyc}^{\mathbf{ER}} = \mathbb{E}_x [||\mathbf{R}(\mathbf{E}(x)) - x||_1] \quad (588)$$

$$\mathcal{L}_{cyc}^{\mathbf{T}_i} = \mathbb{E}_x [||\mathbf{T}_i(\mathbf{E}(x)) - \mathbf{E}(\mathbf{R}(\mathbf{T}_i(\mathbf{E}(x))))||_1] \quad (589)$$

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.

## Tacotron 2

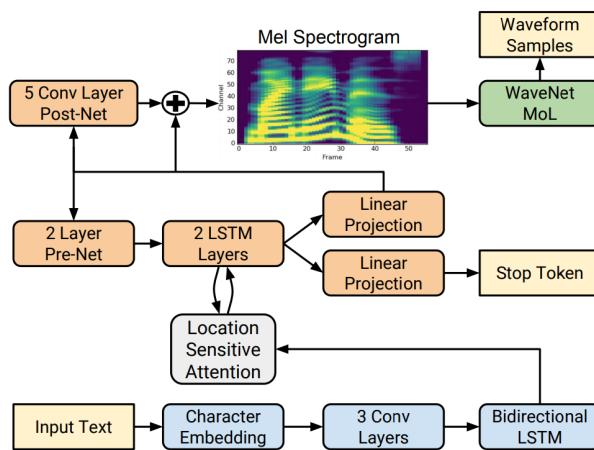
Table of Contents Local

Written by Brandon McKinzie

Shen et al., “NATURAL TTS SYNTHESIS BY CONDITIONING WAVENET ON MEL SPECTROGRAM PREDICTIONS” *Google, UCB* (Feb 2018).

**TL;DR:** Seq2seq network for mapping characters to mel-scale spectrograms<sup>115</sup>, followed by a modified WaveNet vocoder.

**Spectrogram Prediction Network** (2.2). Illustrated below.



- Encoder: character language model architecture. They use convolution layers in the middle under the hypothesis that these should learn ngram-like representations.
- Attention Decoder: **location-sensitive** encoder-decoder attention<sup>116</sup>. At each timestep, instead of incorporating the previous decoder prediction, they first feed it through a **pre-net** (2 FF layers). Furthermore, they also have a **post-net** which predicts a residual to add to the prediction “to improve the overall reconstruction”.

Finally, they also project the decoder/attention outputs to a scalar (sigmoid activation) for predicting probability that output sequence has completed.

<sup>115</sup>A **mel-frequency spectrogram** is obtained by applying a nonlinear transform to the frequency axis of the short-time Fourier transform (STFT).

<sup>116</sup>Basically additive attention with cumulative weights.

**WaveNet Vocoder** (2.3). Inverts the mel spectrogram feature representation into time-domain waveform samples. Specifically, I assume this means they take each time slice of the predicted spectrogram (a vector over frequencies?) and do **local conditioning** (as described in WaveNet paper) with  $\mathbf{h}_t$  the spectrogram at time slice  $t$ .

### Digression: Fourier Transforms

We measure sound by pressure/amplitudes as a function of time. What we actually measure is the *sum total* of all the individual sinusoidal waves each with their own frequency/amplitude. To reiterate: in the *time domain* our x-axis is time (duh) and our y-axis is amplitude. The **Fourier transform** maps our time-domain representation into a *frequency domain*. Now, the x-axis is frequency (duh) and the y-axis represents how much of our original signal consisted of waves with a given frequency. For example, if our original signal was literally a 3 Hz wave superimposed over some 5 Hz wave, our frequency-domain plot would show two spikes at x=3 and x=5, and be zero everywhere else.

Formally, let  $g(t)$  denote the amplitude of some sound wave at time  $t$ .

$$g(t)e^{-2\pi i f t} \quad (590)$$

where

- The negative sign in the exponent is a convention that we should think about *clockwise* rotations.
- $f$  denotes the “winding” frequency: the number of full cycles represented in our wound up graph.
- Multiplying by  $g(t)$  means that the distance-to-origin (magnitude) at time  $t$  will always equal  $g(t)$ .

Similarly, the **inverse Fourier transform** maps from frequency domain to time domain.

Remember that the trick for identifying the component frequencies is by computing the **center of mass** of this formula:

$$\frac{1}{t_2 - t_1} \int_{t_1}^{t_2} g(t)e^{-2\pi i f t} dt \quad (591)$$

where, to do this for one full circle, we’d set  $t_1=0$  and  $t_2=1/f$ . If our signal  $g(t)$  does contain the frequency  $f$ , this integral will be relatively large in comparison with other frequencies. The final formula for the Fourier transform is just removing the scale factor (i.e. the FT is just the CoM scaled by the time interval of our signal):

$$\hat{g}(\omega) \triangleq \Re \left\{ \int_{t_1}^{t_2} g(t)e^{-2\pi i \omega t} dt \right\} \quad (592)$$

More intuition regarding removal of the scale factor: if there is a component wave in  $g(t)$  that only exists for a small portion of time  $\delta t$ , it would have a smaller value of  $\hat{g}(\omega)$  than a component of some different frequency (but same amplitude) that persisted throughout the entirety of our time interval.

### Digression: Spectrograms

Now that we know about Fourier transforms (above example), we can define what a spectrogram is. Say that, instead of a single function  $g(t)$ , I split time into various windows  $\{t_0, t_1, t_2, t_T\}$  and store a separate function,  $g_\tau(t)$  with  $1 \leq \tau \leq T$ , for each window. Then, I apply the FT on each function, resulting in a set of  $\hat{g}_\tau(\omega)$ . If I plot a 2D heatmap with my x-axis denoting the window  $\tau$ , the y-axis denoting frequency  $\omega$ , and values (color) denoting  $\hat{g}_\tau(\omega)$ , I have a **spectrogram** (technically I need to transform my y-axis to  $\log \omega$  and my color/value axis to Decibels).

## Glow

Table of Contents Local

Written by Brandon McKinzie

Kingma et al., “Glow: Generative Flow with Invertible 1x1 Convolutions” *OpenAI*, (July 2018).

**Flow-based Generative Models (2).**

$$\mathbf{z} \sim p_{\theta}(\mathbf{z}) \quad (593)$$

$$\mathbf{x} = \mathbf{g}_{\theta}(\mathbf{z}) \quad (594)$$

where  $\mathbf{g}$  is invertible. *Inference* is done by  $\mathbf{z} = \mathbf{f}_{\theta}(\mathbf{x}) = \mathbf{g}_{\theta}^{-1}(\mathbf{x})$ .

$$\log p_{\theta}(\mathbf{x}) = \log p_{\theta}(\mathbf{z}) + \log |\det \frac{\partial \mathbf{z}}{\partial \mathbf{x}}| \quad (595)$$

$$= \log p_{\theta}(\mathbf{z}) + \sum_{i=1}^K \log |\det \frac{\partial \mathbf{h}_i}{\partial \mathbf{h}_{i-1}}| \quad (596)$$

The log-determinant gives the change in log-density when going  $\mathbf{h}_{i-1} \rightarrow \mathbf{h}_i$ . We can see that maximum likelihood will encourage  $\mathbf{f}_{\theta}(\mathbf{x})$  to increase volume.

## WaveGlow

Table of Contents Local

Written by Brandon McKinzie

Prenger et al., “WaveGlow: A Flow-Based Generative Network for Speech Synthesis”, NVIDIA, (Oct 2018).

**TL;DR:** flow-based **vocoder**<sup>117</sup>

**WaveGlow** (2). Model the distribution of audio samples *conditioned on* a mel-spectrogram. Note that the *forward pass* is defined as going from  $\mathbf{x}$  to  $\text{vec}z$ . Train by minimizing the negative log-likelihood:

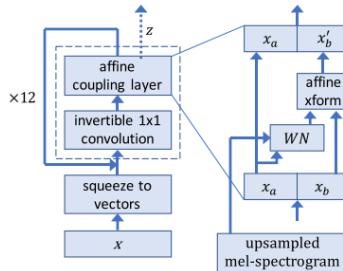
$$\mathbf{z} \sim \mathcal{N}(\mathbf{z}; \mathbf{0}, \mathbf{I}) \quad (597)$$

$$\mathbf{x} = \mathbf{f}_0 \circ \mathbf{f}_1 \circ \dots \circ \mathbf{f}_k(\mathbf{z}) \quad (598)$$

$$\log p_\theta(\mathbf{x}) = \log p_\theta(\mathbf{z}) + \sum_{i=1}^k \log |\det(\mathbf{J}(\mathbf{f}_i^{-1}(\mathbf{x})))| \quad (599)$$

$$\mathbf{z} = \mathbf{f}_k^{-1} \circ \mathbf{f}_{k-1}^{-1} \circ \dots \circ \mathbf{f}_0^{-1}(\mathbf{x}) \quad (600)$$

and  $\mathbf{x}$  denotes the output audio. The mel-spectrogram conditioning happens in the affine coupling layer(s). Training penalizes the norm of  $\mathbf{z}$  and encourages each layer  $f_i^{-1}$  to increase the log-density volume of the previous layer.



$$\mathbf{x}_a, \mathbf{x}_b = \text{split}(\mathbf{x}) \quad (601)$$

$$\log \mathbf{s}, \mathbf{t} = WN(\mathbf{x}_a, \text{mel-spectrogram}) \quad (602)$$

$$\mathbf{x}'_b = \mathbf{s} \odot \mathbf{x}_b + \mathbf{t} \quad (603)$$

$$\mathbf{f}_{coupling}^{-1}(\mathbf{x}) = concat(\mathbf{x}_a, \mathbf{x}'_b) \quad (604)$$

---

<sup>117</sup>Vocoder: network that transforms time-aligned features (e.g. spectrograms) into audio samples.

# META LEARNING

## CONTENTS

|     |                                |     |
|-----|--------------------------------|-----|
| 8.1 | Online Meta-Learning . . . . . | 211 |
|-----|--------------------------------|-----|

## Online Meta-Learning

Table of Contents Local

Written by Brandon McKinzie

Finn et al., “Online Meta-Learning” *Stanford University*, (Jul 2019).

**Online Learning** (2.3). The agent faces a *sequence* of loss functions  $\{\ell_t\}_{t=1}^\infty$ , one in each round  $t$ . At each such round  $t$ , the learner must decide how to update the model parameters  $\mathbf{w}_t$  that perform well on the loss *sequence*<sup>118</sup>. The most standard notion of **regret** is to compare our learner’s cumulative loss over time with the cumulative loss of the best *fixed* model in hindsight:

$$\text{Regret}_T = \sum_{t=1}^T \ell_t(\mathbf{w}_t) - \min_{\mathbf{w}} \sum_{t=1}^T \ell_t(\mathbf{w}) \quad (605)$$

A simple algorithm for this is **follow the leader (FTL)**<sup>119</sup>:

$$\mathbf{w}_{t+1} = \arg \min_{\mathbf{w}} \sum_{k=1}^t \ell_k(\mathbf{w}) \quad (606)$$

**The Online Meta-Learning Problem** (3). We now allow the agent to perform some local *task-specific* updates to the model before it’s deployed/evaluated at each round  $t$ . This is realized through an update procedure at each round  $t$  as a mapping  $\mathbf{U}_t : \mathbf{w} \mapsto \tilde{\mathbf{w}}$ . For example,  $\mathbf{U}_t$  could implement MAML:

$$\mathbf{U}_t(\mathbf{w}) = \mathbf{w} - \alpha \nabla \hat{\ell}_t(\mathbf{w}) \quad (607)$$

where the notation  $\nabla \hat{\ell}_t$  is meant to be interpreted as “approximate gradient of the true  $\ell_t$  obtained with e.g. a minibatch of data from the task at round  $t$ .” To compute regret, we just replace the inputs  $\mathbf{w}$  to the losses in equation 605 with  $\mathbf{U}_t(\mathbf{w})$ .

---

<sup>118</sup>Not just the current loss – we don’t want to forget what was previously learned.

<sup>119</sup>Not sure why they’re saying this is FTL...we aren’t selecting the current best leader at each step per-say. We are just training on everything we’ve seen. That’s it...

**Algorithm and Analysis** (4). Authors propose the **follow the meta leader (FTML)** algorithm for their update rule, which again is just equation 606 but replacing  $\mathbf{w}$  with  $\mathbf{U}_t(\mathbf{w})$  for the loss function inputs.

| <b>Algorithm 1</b> Online Meta-Learning with FTML                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                               | <b>Algorithm 2</b> FTML Subroutines                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                            |
|---------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|
| <pre> 1: <b>Input:</b> Performance threshold of proficiency, <math>\gamma</math> 2: randomly initialize <math>\mathbf{w}_1</math> 3: initialize the task buffer as empty, <math>\mathcal{B} \leftarrow []</math> 4: <b>for</b> <math>t = 1, \dots</math> <b>do</b> 5:   initialize <math>\mathcal{D}_t = \emptyset</math> 6:   Add <math>\mathcal{B} \leftarrow \mathcal{B} + [\mathcal{T}_t]</math> 7:   <b>while</b> <math> \mathcal{D}_{\mathcal{T}_t}  &lt; N</math> <b>do</b> 8:     Append batch of <math>n</math> new datapoints <math>\{(\mathbf{x}, \mathbf{y})\}</math> to <math>\mathcal{D}_t</math> 9:     <math>\mathbf{w}_t \leftarrow \text{Meta-Update}(\mathbf{w}_t, \mathcal{B}, t)</math> 10:    <math>\bar{\mathbf{w}}_t \leftarrow \text{Update-Procedure}(\mathbf{w}_t, \mathcal{D}_t)</math> 11:    <b>if</b> <math>\mathcal{L}(\mathcal{D}_t^{\text{test}}, \bar{\mathbf{w}}_t) &lt; \gamma</math> <b>then</b> 12:      Record efficiency for task <math>\mathcal{T}_t</math> as <math> \mathcal{D}_t </math> datapoints 13:    <b>end if</b> 14:  <b>end while</b> 15:  Record final performance of <math>\bar{\mathbf{w}}_t</math> on test set <math>\mathcal{D}_t^{\text{test}}</math> for task <math>t</math>. 16:  <math>\mathbf{w}_{t+1} \leftarrow \mathbf{w}_t</math> 17: <b>end for</b> </pre> | <pre> 1: <b>Input:</b> Hyperparameters parameters <math>\alpha, \eta</math> 2: <b>function</b> <math>\text{Meta-Update}(\mathbf{w}, \mathcal{B}, t)</math> 3:   <b>for</b> <math>n_m = 1, \dots, N_{\text{meta}}</math> steps <b>do</b> 4:     Sample task <math>\mathcal{T}_k</math>: <math>k \sim \nu^t(\cdot)</math> // (or a minibatch of tasks) 5:     Sample minibatches <math>\mathcal{D}_k^{\text{tr}}, \mathcal{D}_k^{\text{val}}</math> uniformly from <math>\mathcal{D}_k</math> 6:     Compute gradient <math>\mathbf{g}_t</math> using <math>\mathcal{D}_k^{\text{tr}}, \mathcal{D}_k^{\text{val}}</math>, and Eq. 5 7:     Update parameters <math>\mathbf{w} \leftarrow \mathbf{w} - \eta \mathbf{g}_t</math> // (or use Adam) 8:   <b>end for</b> 9:   Return <math>\mathbf{w}</math> 10: <b>end function</b> 11: <b>function</b> <math>\text{Update-Procedure}(\mathbf{w}, \mathcal{D})</math> 12:   Initialize <math>\bar{\mathbf{w}} \leftarrow \mathbf{w}</math> 13:   <b>for</b> <math>n_g = 1, \dots, N_{\text{grad}}</math> steps <b>do</b> 14:     <math>\bar{\mathbf{w}} \leftarrow \bar{\mathbf{w}} - \alpha \nabla \mathcal{L}(\mathcal{D}, \bar{\mathbf{w}})</math> 15:   <b>end for</b> 16:   Return <math>\bar{\mathbf{w}}</math> 17: <b>end function</b> </pre> |

# OPTIMIZATION

## CONTENTS

|       |                                                                                 |     |
|-------|---------------------------------------------------------------------------------|-----|
| 9.1   | Non-Convex Optimization for Machine Learning . . . . .                          | 214 |
| 9.1.1 | Non-Convex Projected Gradient Descent (3) . . . . .                             | 217 |
| 9.2   | Co-sampling: Training Robust Networks for Extremely Noisy Supervision . . . . . | 218 |
| 9.3   | Connectionist Temporal Classification . . . . .                                 | 219 |
| 9.3.1 | Sequence Modeling With CTC . . . . .                                            | 221 |
| 9.4   | Large-Scale Study of Curiosity Driven Learning . . . . .                        | 224 |
| 9.5   | The Marginal Value of Adaptive Gradient Methods in Machine Learning . . . . .   | 225 |
| 9.6   | Deep Double Descent . . . . .                                                   | 226 |
| 9.7   | A Tutorial on Bayesian Optimization of Expensive Cost Functions . . . . .       | 227 |
| 9.8   | Practical Bayesian Optimization of Machine Learning Algorithms . . . . .        | 229 |
| 9.8.1 | Review of Gaussian Processes . . . . .                                          | 231 |
| 9.9   | Dragonfly . . . . .                                                             | 233 |

## 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

$$f(\mathbf{y}) \geq f(\mathbf{x}) + \langle \nabla f(\mathbf{x}), \mathbf{y} - \mathbf{x} \rangle \quad (608)$$

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>120</sup>.

---

<sup>120</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>121</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

$$\frac{\alpha}{2} \|\mathbf{x} - \mathbf{y}\|_2^2 \leq f(\mathbf{y}) - f(\mathbf{x}) - \langle \nabla f(\mathbf{x}), \mathbf{y} - \mathbf{x} \rangle \leq \frac{\beta}{2} \|\mathbf{x} - \mathbf{y}\|_2^2 \quad (609)$$

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>122</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  $\forall \mathbf{x}, \mathbf{y} \in \mathbb{R}^p$ ,

$$|f(\mathbf{x}) - f(\mathbf{y})| \leq B \cdot \|\mathbf{x} - \mathbf{y}\|_2 \quad (610)$$

**Jensen's Inequality.** Generalizes behavior of convex functions on convex combinations<sup>123</sup>.

If  $X$  is a R.V. taking values in the domain of a convex function  $f$ , then

$$\mathbb{E}[f(X)] \geq f(\mathbb{E}[X]) \quad (611)$$

---

<sup>121</sup>Strong smoothness alone does not imply convexity.

<sup>122</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>123</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

$$\Pi_{\mathcal{C}}(\mathbf{z}) := \arg \min_{\mathbf{x} \in \mathcal{C}} \|\mathbf{x} - \mathbf{z}\|_2 \quad (612)$$

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>124</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  $\lambda \in [0, 1]$ ,

$$\mathbf{x}_\lambda := \lambda \mathbf{x} + (1 - \lambda) \hat{\mathbf{z}} \quad (613)$$

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

$$\min_{\mathbf{x} \in \mathbb{R}^p} f(\mathbf{x}) \quad \text{s.t.} \quad \mathbf{x} \in \mathcal{C} \quad (614)$$

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>124</sup>In this case only,  $\mathcal{C}$  need not be convex

---

### 9.1.1 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

$$\mathcal{B}_0(s) \triangleq \{\mathbf{x} \in \mathbb{R}^p \mid \|\mathbf{x}\|_0 \leq s\} \quad (615)$$

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

$$\mathcal{B}_{rank}(r) \triangleq \{A \in \mathbb{R}^{m \times n} \mid \text{rank}(A) \leq r\} \quad (616)$$

and we want to project some matrix  $A$  onto this set,

$$\Pi_{\mathcal{B}_{rank}(r)}(A) := \arg \min_{X \in \mathcal{B}_{rank}(r)} \|A - X\|_F \quad (617)$$

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

$$f(\mathbf{y}) \geq f(\mathbf{x}) + \langle \nabla f(\mathbf{x}), \mathbf{y} - \mathbf{x} \rangle \quad (618)$$

and a similar rephrasing for restricted strong convexity (RSC) and restricted strong smoothness (RSS).

## 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  $\forall i, j T_{ij} \triangleq \Pr[\tilde{y} = e^j | y = e^i]$ . The matrix is row-stochastic<sup>125</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:  $\mathfrak{R}_T = 1 - \min\left\{\frac{T}{T_k}, \mathfrak{D}\right\}$ 
end
Output:  $f_{w_1}$  and  $f_{w_2}$ 
```

---

<sup>125</sup>Each row sums to 1.

## 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  $\mathcal{Y} = L^*$  be the set of all sequences of a finite vocabulary of  $L$  labels.
- **Data distribution:** Denote by  $\mathcal{D}_{\mathcal{X} \times \mathcal{Y}}$  the probability distribution over samples  $(\mathbf{x}, \mathbf{y})$ . 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  $\pi$ , where

$$p(\pi \mid \mathbf{x}) = \prod_{t=1}^T y_{\pi_t}^t \quad (\forall \pi \in (L')^T) \quad (619)$$

and  $y_k^t$  denoting the probability of observing label  $k$  at time  $t$ . Now that we have paths  $\pi$ , 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  $\mathcal{B} : L'^T \mapsto L^{\leq T}$ . In other words,  $\mathcal{B}(\pi) \rightarrow \ell$ . We can then write the conditional posterior over possible output sequences  $\ell$ :

$$p(\ell \mid \mathbf{x}) = \sum_{\pi \in \mathcal{B}^{-1}(\ell)} p(\pi \mid \mathbf{x}) \quad (620)$$

**Constructing the Classifier** (3.2). There is no tractable algorithm for exact decoding, i.e. computing

$$h(\mathbf{x}) := \arg \max_{\ell \in L^{\leq T}} p(\ell \mid \mathbf{x}) \quad (621)$$

However, the following two approximate methods work well in practice:

1. **Best Path Decoding.**  $h(\mathbf{x}) \approx \mathcal{B}(\pi^*)$  where  $\pi^* = \arg \max_{\pi \in L'^T} p(\pi \mid \mathbf{x})$ .
2. **Prefix Search Decoding.**

The authors end up using neither of the above, but rather a heuristic approach:

We divide the output sequence into sections that are very likely to begin and end with a blank. We do this by choosing boundary points where the probability of observing a blank label is above a certain threshold. We then calculate the most probable labelling for each section individually and concatenate these to get the final classification.

**The CTC Forward-Backward Algorithm** (4.1). Define the total probability of the first  $s$  output labels,  $\ell_{\langle 1 \dots s \rangle}$ , given the first  $t$  inputs as

$$\alpha_t(s) \triangleq p(\ell_{\langle 1 \dots s \rangle} \mid \mathbf{x}_{\langle 1 \dots t \rangle}) \quad (622)$$

$$= \sum_{\substack{\pi \in L'^T \\ \mathcal{B}(\pi_{\langle 1 \dots t \rangle}) = \ell_{\langle 1 \dots s \rangle}}} \prod_{t'=1}^t y_{\pi_{t'}}^{t'} \quad (623)$$

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 input step  $t=1$ , followed by the recursion rule:

$$\alpha_1(s) = \begin{cases} y_\epsilon^1 & s=1 \\ y_{\ell_1}^1 & s=2 \\ 0 & s > 2 \end{cases} \quad (624)$$

$$\alpha_t(s) = \begin{cases} \bar{\alpha}_t(s)y_{\ell'_s}^t & \ell'_s = b \text{ or } \ell'_{s-2} \\ (\bar{\alpha}_t(s) + \alpha_{t-1}(s-2))y_{\ell'_s}^t & \text{otherwise} \end{cases} \quad (625)$$

$$\bar{\alpha}_t(s) \triangleq \alpha_{t-1}(s) + \alpha_{t-1}(s-1) \quad (626)$$

### Interpretation:

- *Initialization:* Given the first input  $x_1$ , the only valid paths<sup>126</sup> that could potentially result in the final labeling  $\ell$  are  $\{\epsilon\}$  and  $\{\ell_1\}$ . These respectively correspond to the augmented labelings  $\{\epsilon\}$  and  $\{\epsilon, \ell_1\}$ . This is what the authors are referring to when they say “We allow all prefixes to start with either a blank or the first symbol in  $\ell$ .”
- *Case 1.*  $\bar{\alpha}_t(s)$  corresponds to a right arrow ( $[s, t-1] \rightarrow [s, t]$ ) and a right-down-one arrow ( $[s-1, t-1] \rightarrow [s, t]$ ) (see explanations below).
- *Case 2.*  $\bar{\alpha}_t(s) + \alpha_{t-1}(s-2)$  corresponds to the 2 arrows from case 1 and additional a right-down-two arrow ( $[s-2, t-1] \rightarrow [s, t]$ ).

**Reading the lattice arrow diagrams.** The usage of the augmented label sequence  $\ell'$  can make reading these very confusing. Here’s how to read them:

- *Rows:* row  $s$  corresponds to  $\ell'_s$ .

---

<sup>126</sup>which remember are always the same length as  $x$

- *Columns*: column  $t$  corresponds to  $\mathbf{x}_t$ .
- *Arrows/Traversal*: this was the confusing part for me. The arrows correspond to a transition in a given *path*. Each arrow direction has a different interpretation:
  - right: the path either had a repeated character (which would've been merged into a single output label) or a repeated blank token (if the row is a blank token) which would've been removed entirely in the final output label sequence.
  - right-down-one: transition from either label-to-blank or blank-to-label.
  - right-down-two: direct transition between two unique characters.

Ok, I finally get it now. In my opinion, introducing  $\ell'$  as the “augmented label sequence” is confusing/misleading/unnecessary/etc. It is literally just introduced so we can meaningfully talk about transitions between elements of a given path  $\pi$ . The traversal being done in the  $\alpha$  formulas is referencing the valid *paths*, NOT the final labels (at least directly). Especially confusing is when the distill article says crap like “can’t jump over  $z_{s-1}$ ” which is just total nonsense – no one is jumping over anything! Literally all they mean is “transition between unique characters in a path”. Lost so many hours confused over the wording in the distill article, which only served to confuse me further (just read the paper).

The final probability of the label sequence  $\ell$  given input sequence  $\mathbf{x}$  is thus:

$$p(\ell \mid \mathbf{x}) = \alpha_T(|\ell'|) + \alpha_T(|\ell'| - 1) \quad (627)$$

which corresponds to “total probability of all paths ending in a blank token plus total probability of all paths ending in the final label of  $\ell$ .”

### 9.3.1 SEQUENCE MODELING WITH CTC

Notes on this distill article (note that I do NOT recommending reading this article – the wording is horrendously confusing compared to what’s actually going on).

**Introduction.** CTC is an 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 ( $T \neq U$ ).

**Alignment.** An **alignment** between input sequence  $X$  and label sequence  $Y$  is a function  $f : X \mapsto Y$ . 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:

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.

### 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.

**Loss Function.** 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  $X$ :

$$p(Y | X) = \sum_{\mathcal{A} \in \mathcal{A}_{X,Y}} \prod_{t=1}^T p_t(a_t | X) \quad (628)$$

where  $\mathcal{A}$  is one of the valid alignments from  $\mathcal{A}_{X,Y}$ . The value of  $a_t$  obeys the set of three constraints listed above.

How can we compute the loss efficiently? Let  $\mathbf{z} \triangleq [\epsilon, y_1, \epsilon, y_2, \dots, \epsilon, y_U, \epsilon]$ , and let  $\alpha$  denote the **score of the merged alignments** at a given node [in the **CTC lattice**]. We compute the forward probabilities  $\alpha_t(s)$ , defined as the probability of arriving at [prefix of] augmented label sequence  $\mathbf{z}_{\langle 1 \dots s \rangle}$  given unmerged alignments up to input step  $t$ :

$$\alpha_t(s) \triangleq p(\mathbf{z}_{\langle 1 \dots s \rangle} | \mathbf{x}_{\langle 1 \dots t \rangle})$$

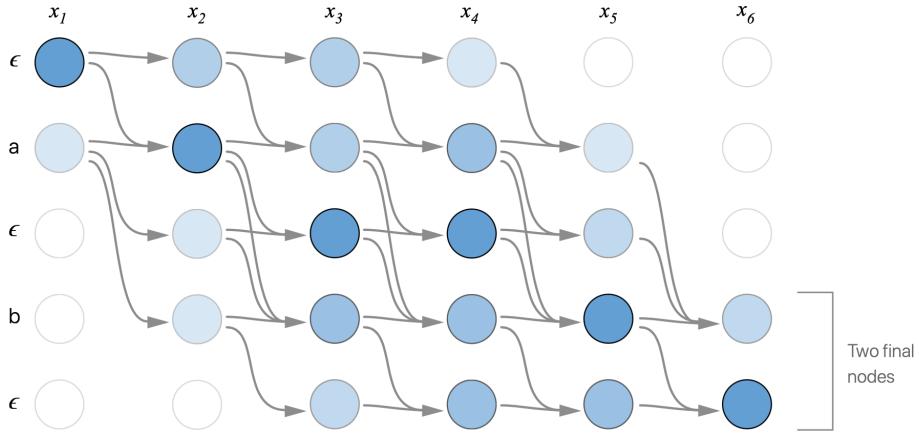
The key insight is that we can compute  $\alpha_t$  as long as we know  $\alpha_{t-1}$ . There are two cases to consider.

1. (1.1)  $z_s$  is the blank token  $\epsilon$ . 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,  $z_{s-1}$ . In other words,

$$\alpha_t(s) = p(z_s=\epsilon | x_t) \cdot (\alpha_{t-1}(s) + \alpha_{t-1}(s-1)) \quad (629)$$

2. (1.2)  $z_s$  is the same label as at step  $s - 2$ . This occurs when  $Y$  has repeated labels next to each other.

$$\alpha_t(s) = p(z_s=z_{s-2} | x_t) \cdot (\alpha_{t-1}(s) + \alpha_{t-1}(s-1)) \quad (630)$$



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)  **$z_{s-1}$  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  $z_{s-2}$  at the previous time ( $t - 1$ ) and then emitted  $z_s$  immediately after at time  $t$ .

$$\alpha_t(s) = p(z_s | x_t) \cdot (\alpha_{t-1}(s-2) + \alpha_{t-1}(s-1) + \alpha_{t-1}(s)) \quad (631)$$

## 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*):

$$r_t = -\log p(\phi(x_{t+1}) | x_t, a_t) \quad (632)$$

The authors choose to model the next state embedding with a Gaussian,

$$\phi(x_{t+1}) | x_t, a_t \sim \mathcal{N}(f(x_t, a_t), \epsilon) \quad (633)$$

$$r_t = \|f(x_t, a_t) - \phi(x_{t+1})\|_2^2 \quad (634)$$

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).

## 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>127</sup>

$$[\text{regular}] \quad \Delta w_{k+1} = -\alpha_k \tilde{\nabla} f(w_k) \quad (635)$$

$$[\text{momentum}] \quad \Delta w_{k+1} = -\alpha_k \tilde{\nabla} f(w_k + \gamma_k \Delta w_k) + \beta_k \Delta w_k \quad (636)$$

$$[\text{adaptive}] \quad \Delta w_{k+1} = -\alpha_k \mathbf{H}_k^{-1} \tilde{\nabla} f(w_k + \gamma_k \Delta w_k) + \beta_k \mathbf{H}_k^{-1} \mathbf{H}_{k-1} \Delta w_k \quad (637)$$

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:

$$H_k := \text{diag} \left[ \left( \sum_{i=1}^k \eta_i \mathbf{g}_i \odot \mathbf{g}_i \right)^{1/2} \right] \quad (638)$$

**The Potential Perils of Adaptivity (3).** Consider the binary least-squares classification problem, where we aim to minimize

$$R_s[w] := \frac{1}{2} \|Xw - y\|_2^2 \quad (639)$$

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>127</sup>I'm defining  $\Delta w_{k+1} \triangleq w_{k+1} - w_k$ .

## Deep Double Descent

Table of Contents Local

Written by Brandon McKinzie

Nakkiran et al., “Deep Double Descent: Where Bigger Models and More Data Hurt” *OpenAI*, (Dec 2019).

**TL;DR:** For a variety of DL tasks, increasing model size (or training epochs) first leads to *worse* performance, and then gets better. Define a new complexity measure called **effective model complexity**.

*Informally, our intuition is that for model-sizes at the interpolation threshold, there is effectively only one model that fits the train data and this interpolating model is very sensitive to noise in the train set and/or model mis-specification. That is, since the model is just barely able to fit the train data, forcing it to fit even slightly-noisy or mis-specified labels will destroy its global structure, and result in high test error.*

**Introduction** (1). Bias-variance tradeoff suggests that increasing model complexity results in lower bias and higher variance. After a certain threshold [of increasing model complexity], conventional wisdom says that the model will “overfit” with the variance term dominating the test error. Therefore, increasing the complexity beyond this threshold *should* merely result in a larger variance term and thus worse MSE (i.e. *larger models are worse* [after a certain point]). Conventioanal wisdom also tells us that *more data is always better* [for improving the test MSE].

**Results** (2). Define a **training procedure**  $\mathcal{T}$  to be any procedure that takes as input a training set  $S = \{(x_i, y_i)\}_{i=1}^n$  and outputs a classifier  $\mathcal{T}(S) : x \mapsto y$ .

### Effective Model Complexity

The **EMC** of  $\mathcal{T}$ , wrt distribution  $\mathcal{D}$  and  $\epsilon > 0$ , is defined as:

$$\text{EMC}_{\mathcal{D}, \epsilon}(\mathcal{T}) := \max\{n \mid \mathbb{E}_{S \sim \mathcal{D}^n} [\text{Error}_S(\mathcal{T}(S))] \leq \epsilon\} \quad (640)$$

where  $\text{Error}_S(M)$  is the mean error of model  $M$  on train samples  $S$ .

In other words, the EMC of  $\mathcal{T}$  is the maximum number of training samples for which  $\mathcal{T}$  gets [on average] zero *training error*. The authors then hypothesize the 3 following regimes (assuming  $n$  training examples):

- **Under-parameterized:**  $\text{EMC}_{\mathcal{D}, \epsilon}(\mathcal{T}) \ll n$ . Any  $\delta\mathcal{T}$  resulting in larger  $\text{EMC}_{\mathcal{D}, \epsilon}(\mathcal{T})$  will *decrease* test error.
- **Over-parameterized:**  $\text{EMC}_{\mathcal{D}, \epsilon}(\mathcal{T}) \gg n$ . (same as above).
- **Critically parameterized:**  $\text{EMC}_{\mathcal{D}, \epsilon}(\mathcal{T}) \approx n$ . Any  $\delta\mathcal{T}$  resulting in larger  $\text{EMC}_{\mathcal{D}, \epsilon}(\mathcal{T})$  may *decrease or increase* test error.

## A Tutorial on Bayesian Optimization of Expensive Cost Functions

Table of Contents   Local

*Written by Brandon McKinzie*

Brochu, Eric Cora, Vlad M. de Freitas, Nando, “ A Tutorial on Bayesian Optimization of Expensive Cost Functions, with Application to Active User Modeling and Hierarchical Reinforcement Learning ”, (Dec 2010).

**The Bayesian Optimization Approach** (2). Goal will to be *maximize* a real-valued function  $f$ :

$$\mathbf{x}^* = \arg \max_{\mathbf{x}} f(\mathbf{x}) \quad (641)$$

$$\|f(\mathbf{x}_1) - f(\mathbf{x}_2)\| \leq C\|\mathbf{x}_1 - \mathbf{x}_2\| \quad \forall \mathbf{x}_1, \mathbf{x}_2 \in \mathcal{A} \quad (642)$$

where the second line states our assumption that  $f$  is **Lipschitz-continuous**. Bayesian optimization uses a prior  $p(f)$  and evidence  $p(\mathcal{D})$  to define a posterior distribution over functions given evidence. The optimization follows the principle of *maximum expected utility*, with expectation taken wrt the posterior. The utility function will be referred to as the **acquisition function**.

**Priors over Functions** (2.1). By far the most common choice for the prior over functions is the **Gaussian process prior**. A GP is completely specified by its mean and covariance function.

*A GP is an extension of the multivariate Gaussian distribution to an infinite-dimension stochastic process for which any finite combination of dimensions will be a Gaussian distribution.*

$$f(\mathbf{x}) \sim \mathcal{GP}(m(\mathbf{x}), k(\mathbf{x}, \mathbf{x}')) \quad (643)$$

Here, we’ll use the zero mean function. A common choice for the kernel is the **squared exponential function**:

$$k(\mathbf{x}_i, \mathbf{x}_j) = \exp\left\{-\frac{1}{2}\|\mathbf{x}_i - \mathbf{x}_j\|^2\right\} \quad (644)$$

Now say we observed some data  $\mathcal{D}_{1:t}$  of queries/outputs to/from  $f$ . Given some next query  $\mathbf{x}_{t+1}$ , the predictive distribution (probs of output from the next query) is:

$$P(f_{t+1} | \mathcal{D}_{1:t}, \mathbf{x}_{t+1}) = \mathcal{N}\left(\mu_t(\mathbf{x}_{t+1}), \sigma_t^2(\mathbf{x}_{t+1})\right) \quad (645)$$

$$\mu_t(\mathbf{x}_{t+1}) = \mathbf{k}^T \mathbf{K}^{-1} \mathbf{f}_{1:t} \quad (646)$$

$$\sigma_t^2(\mathbf{x}_{t+1}) = k(\mathbf{x}_{t+1}, \mathbf{x}_{t+1}) - \mathbf{k}^T \mathbf{K}^{-1} \mathbf{k} \quad (647)$$

$$k_i \triangleq k(\mathbf{x}_{t+1}, \mathbf{x}_i) \quad (648)$$

**Choice of Covariance Functions** (2.2). The choice of  $k$  determines the smoothness properties of samples drawn from the GP. It is more common to incorporate additional *hyperparameters*  $\theta$  into the covariance function, e.g.:

$$k(\mathbf{x}_i, \mathbf{x}_j) = \exp\left(-\frac{1}{2}(\mathbf{x}_i - \mathbf{x}_j)^T \text{diag}(\boldsymbol{\theta})^{-2} (\mathbf{x}_i - \mathbf{x}_j)\right) \quad (649)$$

Typically the values of  $\theta$  are learned by “seeding” with a few random samples and maximizing the log-likelihood of the evidence  $\mathcal{D}$  given  $\theta$ . It is also common to have a scalar hyperparameter applied to  $k$  to control the magnitude of the variance.

**Acquisition Functions for Bayesian Optimization** (2.3). Typically, acquisition functions are defined such that high acquisition corresponds to potentially high values of the objective function, whether because the prediction is high, the uncertainty is great, or both. One approach, referred to as **probability of improvement** over the incumbent (current best)  $f(\mathbf{x}^+)$ , where  $\mathbf{x}^+ = \arg \max_{\mathbf{x}_i \in \mathbf{x}_{(1 \dots t)}} f(\mathbf{x}_i)$ , defines acquisition function

$$PI(\mathbf{x}) = \Pr[f(\mathbf{x}) \geq f(\mathbf{x}^+)] = \Phi\left(\frac{\mu(\mathbf{x}) - f(\mathbf{x}^+)}{\sigma(\mathbf{x})}\right) \quad (650)$$

where  $\mu$  and  $\sigma$  refer to the current predictive posterior defined earlier, and  $\Phi$  is the normal cumulative distribution. Note how this is a pure exploitation strategy. We can add a trade-off parameter  $\xi \geq 0$ :

$$PI(\mathbf{x}) = \Pr[f(\mathbf{x}) \geq f(\mathbf{x}^+) + \xi] = \Phi\left(\frac{\mu(\mathbf{x}) - f(\mathbf{x}^+) - \xi}{\sigma(\mathbf{x})}\right) \quad (651)$$

which has the effect of subtracting  $\xi/\sigma$  in the input of the cumulative distribution function – higher uncertain y points are now given relatively better chances of being selected.

## Practical Bayesian Optimization of Machine Learning Algorithms

Table of Contents   Local

*Written by Brandon McKinzie*

P. Snoek, H. Larochelle, and R. Adams, “ Practical Bayesian Optimization of Machine Learning Algorithms ” *University of Toronto*, (Aug 2012).

**Introduction** (1). For some continuous function  $f$ , BO assumes that  $f$  was sampled from a **Gaussian Process** (GP) and maintains a posterior distribution for  $f$  given observations. Here, the observations are the measure of generalization performance<sup>128</sup> for an associated setting of the input parameters to  $f$ . In order to choose/propose the next set of parameters to feed to  $f$ , we can do one of the following popular approaches:

- Optimize the **expected improvement** (EI) over the current best result **TODO**: formal definition
- Optimize the GP **upper confidence bound** (UCB) **TODO**: formal definition.

Main contributions of this paper:

- Identification of good practices for BO of machine learning algorithms.
- Argue that a fully Bayesian treatment of GP kernel params is of critical importance.
- Description of new algorithm that accounts for cost in experiments.
- Propose algorithm that can take advantage of multiple cores in parallel.

**Bayesian Optimization with GP Priors** (2). Two major choices when performing BO:

1. Select a prior over functions that will express assumptions about  $f$ .
2. Choose an **acquisition function**, which is used to construct a utility function from the model posterior, allowing us to determine the next point to evaluate.

For choice (1), we use the GP prior. For (2), first denote the acquisition function  $a : \mathcal{X} \rightarrow \mathbb{R}^+$  that basically gives a score for any input  $\mathbf{x}$ , such that we can get a proposal for the next input  $\mathbf{x}_{next}$  to feed to  $f$  via  $\mathbf{x}_{next} := \arg \max_{\mathbf{x}} a(\mathbf{x})$ . Note that  $a(\mathbf{x})$  is shorthand for  $a(\mathbf{x}; \{\mathbf{x}_n, y_n\}, \theta)$  where  $\theta$  are the GP hyperparameters (?) (same goes for predictive  $\mu$  and  $\sigma$  below).

---

<sup>128</sup>I believe the observations are always the output value of  $f$ , to be clear.

Acquisition functions<sup>129</sup>

$$[\text{Prob of Impr.}] \quad a_{PI}(\mathbf{x}) = \Phi(\gamma(\mathbf{x})) \quad \gamma(\mathbf{x}) = \frac{f(\mathbf{x}_{best}) - \mu(\mathbf{x})}{\sigma(\mathbf{x})} \quad (652)$$

Impr  $\equiv$  Improvement

$$[\text{Expected Impr.}] \quad a_{EI}(\mathbf{x}) = \sigma(\mathbf{x})(\gamma(\mathbf{x})\Phi(\gamma(\mathbf{x}))+\mathcal{N}(\gamma(\mathbf{x}); 0, 1)) \quad (653)$$

$$[\text{UCB}] \quad a_{LCB}(\mathbf{x}) = \mu(\mathbf{x}) - \kappa\sigma(\mathbf{x}) \quad (654)$$

with a tunable  $\kappa$  to balance exploitation and exploration<sup>130</sup>.

Covariance functions: Below are the **automatic relevance determination** (ARD) squared exponential and Matérn 5/2 kernels:

$$K_{SE}(\mathbf{x}, \mathbf{x}') = \theta_0 \exp\left\{-\frac{1}{2}r^2(\mathbf{x}, \mathbf{x}')\right\} \quad r^2(\mathbf{x}, \mathbf{x}') = \sum_{d=1}^D (x_d - x'_d)^2 / \theta_d^2 \quad (655)$$

$$K_{M52}(\mathbf{x}, \mathbf{x}') = \theta_0 \left(1 + \sqrt{5r^2(\mathbf{x}, \mathbf{x}')} + \frac{5}{3}r^2(\mathbf{x}, \mathbf{x}')\right) \exp\left\{\sqrt{5r^2(\mathbf{x}, \mathbf{x}')}\right\} \quad (656)$$

GP Hyperparameters:

- Length scales  $\theta_{1:D}$
- Covariance amplitude  $\theta_0$
- Observation noise  $\nu$
- Constant mean  $m$  (?)

We can either use a point estimate of these parameters by optimizing the marginal likelihood under the GP:

$$\arg \max_{\theta, \nu, m} p(\mathbf{y} \mid \{\mathbf{x}_n\}_{n=1}^N, \theta, \nu, m) = \arg \max_{\theta, \nu, m} \mathcal{N}(\mathbf{y} \mid m\mathbf{1}, \Sigma_\theta + \nu\mathbf{I}) \quad (657)$$

or the [fully-Bayesian] marginalization via the **integrated acquisition function**:

$$\hat{a}(\mathbf{x}) = \int a(\mathbf{x})p(\theta \mid \{\mathbf{x}_n, y_n\})d\theta \quad (658)$$

Monte Carlo Acquisition for parallelizing BO. Consider situation where  $N$  evaluations have completed, yielding data  $\{\mathbf{x}_n, y_n\}_{n=1}^N$ , and  $J$  evaluations are pending at locations  $\{\mathbf{x}_j\}_{j=1}^J$  (running on different threads). Ideally, we'd choose  $\mathbf{x}_{next}$  based on

$$\hat{a}(\mathbf{x}; \{\mathbf{x}_j\}) = \int_{\mathbb{R}^J} a(\mathbf{x}; \{\mathbf{x}_j, y_j\})p(y_{1:J} \mid \mathbf{x}_{\langle 1 \dots J \rangle})dy_1 \dots dy_J \quad (659)$$

Instead, since we don't have access to all the pending evaluations, we compute Monte Carlo estimates to approximate the above integral.

<sup>129</sup>Note that we are trying to minimize  $f$  in this context.

<sup>130</sup> $a_{LCB}$  exploits lower confidence bounds (upper, when considering maximization) to construct acquisition functions that minimize regret over the course of their optimization.

---

### 9.8.1 REVIEW OF GAUSSIAN PROCESSES

---

#### Gaussian Process Review

A GP is defined as a collection of random variables  $f(\mathbf{x})^a$ , any finite number of which have a joint Gaussian distribution. A GP is completely specified by mean function  $m(\mathbf{x})$  and covariance function  $k(\mathbf{x}, \mathbf{x}')$ .

$$m(\mathbf{x}) = \mathbb{E}[f(\mathbf{x})] \quad (660)$$

$$k(\mathbf{x}, \mathbf{x}') = \mathbb{E}[(f(\mathbf{x}) - m(\mathbf{x}))(f(\mathbf{x}') - m(\mathbf{x}'))] \quad (661)$$

Again, without getting too wrapped up in the “meaning”, we *define* the notation (which at this point is still just a placeholder basically) for representing a GP:

$$f(\mathbf{x}) \sim \mathcal{GP}(m(\mathbf{x}), k(\mathbf{x}, \mathbf{x}')) \quad (662)$$

**Key Insight:** if we don’t know the distribution that  $f$  is sampled from<sup>b</sup> we can instead *define* how our covariance function should look, which itself *implies* a distribution over functions  $f$  (!!).

Consider the widely popular choice: the **squared exponential** covariance function:

$$k(\mathbf{x}_p, \mathbf{x}_q) \triangleq \exp\left(-\frac{1}{2}|\mathbf{x}_p - \mathbf{x}_q|^2\right) \quad (663)$$

It can be shown that, by defining our covariance function as above, we are implicitly modeling  $f$  as being sampled from a Bayesian linear regression model with an infinite number of basis (feature) functions (see prev footnote).

<sup>a</sup>We also refer to  $f(\mathbf{x})$  as a real process. Don’t get too wrapped up in thinking about how one samples from a space of functions (yet) – we haven’t gotten to that part! For now, just interpret this as a regular collection of random variables.

<sup>b</sup>An example where we *do* know is if we let e.g.  $f(x) \triangleq \phi(x)^T w$ , with feature map  $\phi$  and weights  $w \sim \mathcal{N}(0, \sigma_w)$  – crucially that  $w$  is the only source of randomness and is drawn from a Gaussian. In this case, we can just directly compute the exact formula for  $k(\mathbf{x}, \mathbf{x}')$  and are basically done.

**Sampling Procedure.** Let’s assume we are using the squared exponential covariance function from the above review box. We’ll also assume  $m(\mathbf{x}) = \mathbf{0}$ . **Sampling from the GP prior** works as follows:

- Draw a set of samples over input space  $\{\mathbf{x}_i\}_{i=1}^N$  with each  $\mathbf{x}_i \in \mathbb{R}^d$ .
- Compute the  $N \times N$  covariance matrix  $\mathbf{K}$ , where  $K_{ij} := \exp\left(-\frac{1}{2}|\mathbf{x}_i - \mathbf{x}_j|^2\right)$ .
- Our GP prior for  $f$  is now well-defined as a multivariate Gaussian over  $N$  dimensions (NOT  $d$ )<sup>131</sup>.
- We can sample functions  $f_*$  from our prior – these will *each* be random Gaussian vectors of  $N$  dimensions with interpretation that the  $i$ th element equals the sampled  $f_*$  at input  $\mathbf{x}_i$ . For example, to sample a function  $f(\mathbf{x})$  from the prior:
  1. Compute the **Cholesky decomposition**  $\mathbf{L}$  of the symmetric psd covariance matrix  $\mathbf{K} = \mathbf{L}\mathbf{L}^T$ .
  2. Sample  $N$  times from the standard normal to get  $\mathbf{u} \sim \mathcal{N}(0, I)$ .
  3. Compute  $f(\cdot) = \mathbf{m} + \mathbf{L}\mathbf{u}$ <sup>132</sup>.

Now, how do incorporate observations and sample from the posterior? Let's pretend that we live in infinite-computation land for a moment, and that we've collected all the possible training inputs  $\mathbf{X}_{tr} = \{\mathbf{x}_i^{tr}\}_{i=1}^{N_{tr}}$  and [a disjoint set of] test inputs  $\mathbf{X}_{test} = \{\mathbf{x}_i^{test}\}_{i=1}^{N_{test}}$ .

1. Compute each element of the *combined*  $N \times N$  (where  $N = N_{tr} + N_{test}$ ) squared-exponential covariance matrix to obtain a full (massive) GP prior<sup>133</sup>.
2. Evaluate  $f(\mathbf{x}_i^{tr})$  for all  $\mathbf{x}_i \in \mathbf{X}_{tr}$ , resulting in a vector  $\mathbf{f} \in \mathbb{R}^{N_{tr}}$  of training outputs.
3. Condition our GP prior to obtain the corresponding posterior  $f_* | \mathbf{X}_{test}, \mathbf{X}_{tr}, \mathbf{f}$ , over function evaluations  $f_*$  (outputs) on the test set  $\mathbf{X}_{test}$ .

**Math/other TODO** reorganize section. The **Choleskey decomposition** of a symmetric psd matrix  $\mathbf{A}$  decomposes  $\mathbf{A}$  into a product of a lower triangular matrix  $\mathbf{L}$  and its transpose:

$$\mathbf{L}\mathbf{L}^T = \mathbf{A} \quad (664)$$

where  $\mathbf{L}$  is called the **Cholesky factor**. This decomposition is useful for solving linear systems ( $\mathbf{Ax} = \mathbf{b}$  for vector  $\mathbf{x}$ ) with symmetric psd coefficient matrix  $\mathbf{A}$ .

---

<sup>131</sup>Notice that the functions we sample from this prior are only defined at each of the  $N$  query inputs, so we are really only getting pointwise defined functions...

<sup>132</sup>Our prior mean function is  $\mathbf{m} = \mathbf{0}$ .

<sup>133</sup>Same as before – nothing new here except how I've organized/defined the sampled input sets)

## Dragonfly

Table of Contents Local

Written by Brandon McKinzie

Kandasamy et al., “ Tuning Hyperparameters without Grad Students: Scalable and Robust Bayesian Optimisation with Dragonfly ” *CMU*, (Apr 2020).

**Introduction** (1). Given black-box function  $f : \mathcal{X} \mapsto \mathbb{R}$ , we wish to find its optimum using [as few as possible] repeated evaluations to  $f$ . Methods aim to determine the next point  $x_t$  for evaluation, given knowledge of  $f$  acquired via previous query-observation pairs  $\{(x_i, y_i)\}_{i=1}^{t-1}$ . They do this by taking the argmax of an **acquisition function**  $\varphi_t : \mathcal{X} \mapsto \mathbb{R}$ :

$$x_t \leftarrow \arg \max_{x \in \mathcal{X}} \varphi_t(x) \quad (665)$$

**Review of Gaussian Processes and Bayesian Optimization** (2). A **Gaussian Process** (GP) is characterized by a mean function  $\mu : \mathcal{X} \mapsto \mathbb{R}$  and a kernel (covariance function)  $\kappa : \mathcal{X}^2 \mapsto \mathbb{R}$ . Suppose we are given  $n$  observations  $\mathcal{D}_n = \{(x_i, y_i)\}_{i=1}^n$  from this GP, where  $y_i = f(x_i) + \epsilon_i \in \mathbb{R}$ . Then the *posterior*  $f|\mathcal{D}_n$  is also a GP with mean  $\mu_n$  and covariance  $\kappa_n$ :

$$\mu_n(x) = k^T \left( K + \eta^2 I_n \right)^{-1} Y \quad (666)$$

$$\kappa_n(x, x') = \kappa(x, x') - k^T \left( K + \eta^2 I_n \right)^{-1} k' \quad (667)$$

- $Y \in \mathbb{R}^n$  is a vector with  $Y_i = y_i$
- $k_i = \kappa(x, x_i)$
- $k'_i = \kappa(x', x_i)$
- Gram matrix  $K \in \mathbb{R}^{n \times n}$  is given by  $K_{ij} = \kappa(x_i, x_j)$ .

**Bayesian Optimization**

*A suite of methods for black box optimisation in the Bayesian paradigm.*

- Use a prior belief distribution for  $f$ .
- To determine the next point  $x_t$  for evaluation, construct an acquisition function  $\varphi_t : \mathcal{X} \mapsto \mathbb{R}$  which captures the utility of performing an evaluation at  $x$  according to the posterior  $f$  conditioned on past evaluations.

One popular acquisition function is the **Gaussian process upper confidence bound** (GP-UCB)<sup>134</sup>:

$$\varphi_t(x) = \mu_{t-1}(x) + \sqrt{\beta_t} \sigma_{t-1}(x) \quad (668)$$

---

<sup>134</sup>Upper confidence bound for  $f$ .

where the posterior standard deviation  $\sigma_{t-1}$  is the uncertainty associated with our current estimate of  $f$  (after t-1 observations).

- The  $\mu_{t-1}$  term encourages an *exploitative strategy*: we want to query regions where we already believe  $f$  is high.
- $\sigma_{t-1}$  encourages an *exploratory strategy*: we want to query where we are uncertain about  $f$ lest we miss high valued regions which have not been queried yet.
- $\beta_t$  controls the trade-off between exploration and exploitation.

**Choice of Acquisition** (4.1). Authors employ an adaptive sampling strategy which chooses different acquisition functions at different iterations instead of attempting to pick a single best one.

**GP Hyperparameters** (4.2). Authors either employ maximum likelihood or use posterior sampling at each iteration to obtain a set of GP hyperparameters.

## BO Implementation in Dragonfly (5).

- **Domains** of variables. Support Euclidean, integral, discrete and neural network variables.
- **Kernels**. Squared exponential and Matérn kernels.
- **Optimising the Acquisition**. To maximize  $\varphi_t$  in purely Euclidean spaces, they use **DiRect** or **PDOO**
- **Initialisation**. Bootstrap the BO routine with  $n_{init}$  evaluations. For Euclidean/integral vars, these points are chosen via **latin hypercube sampling**. Discrete vars are sampled uniformly at random.

### BO in Dragonfly with M asynchronous workers

Let  $n_{init}$  be the number of random initial query points to evaluate on.

- $\mathcal{D}_0 \leftarrow$  evaluate  $f$  at  $n_{init}$  points.
- $y_{max} \leftarrow$  maximum  $y$  value in  $\mathcal{D}_0$ .
- Repeat the following for some number of iterations (ok how many tho?)
  - Wait for a worker to finish.
  - $\mathcal{D}_t \leftarrow \mathcal{D}_{t-1} \cup \{(q, y)\}$  where  $(q, y)$  are the worker's previous query and observation.
  - $\mu_{t-1} \leftarrow$  compute posterior GP mean given  $\mathcal{D}_j$  using  $\theta$ .
  - Compute hallucinated posterior (necessary for simple first impl?)
  - $q' \leftarrow$  determine next query for evaluation using acquisition  $\alpha$  and GP  $\mathcal{GP}_{t-1}$ .
  - Re-deploy worker with an evaluation at  $q'$ .

# MISCELLANEOUS

## CONTENTS

|        |                                                                                                |     |
|--------|------------------------------------------------------------------------------------------------|-----|
| 10.1   | Neural Style . . . . .                                                                         | 236 |
| 10.2   | TextRank . . . . .                                                                             | 238 |
| 10.2.1 | Keyword Extraction . . . . .                                                                   | 240 |
| 10.2.2 | Sentence Extraction . . . . .                                                                  | 241 |
| 10.3   | Survey of Text Clustering Algorithms . . . . .                                                 | 242 |
| 10.3.1 | Distance-based Clustering Algorithms . . . . .                                                 | 245 |
| 10.3.2 | Probabilistic Document Clustering and Topic Models . . . . .                                   | 246 |
| 10.3.3 | Online Clustering with Text Streams . . . . .                                                  | 248 |
| 10.4   | Clustering Massive Text Streams . . . . .                                                      | 250 |
| 10.5   | Hierarchical Attention Networks . . . . .                                                      | 252 |
| 10.6   | Neural Architecture Search with Reinforcement Learning . . . . .                               | 255 |
| 10.7   | Neural Ordinary Differential Equations . . . . .                                               | 257 |
| 10.8   | A Framework for Intelligence and Cortical Function . . . . .                                   | 259 |
| 10.9   | A Theoretically Grounded Application of Dropout in Recurrent Neural Networks . . . . .         | 260 |
| 10.10  | Protection Against Reconstruction and Its Applications in Private Federated Learning . . . . . | 261 |
| 10.11  | Product quantization for nearest neighbor search . . . . .                                     | 263 |
| 10.12  | Did the Model Understand the Question? . . . . .                                               | 264 |
| 10.13  | Solving Rubik's Cube with a Robot Hand . . . . .                                               | 265 |
| 10.14  | Decoupled Weight Decay Regularization . . . . .                                                | 267 |
| 10.15  | SmartChoices . . . . .                                                                         | 268 |
| 10.16  | The Contextual Bandits Problem . . . . .                                                       | 269 |
| 10.17  | The principles of adaptation in organisms and machines . . . . .                               | 272 |

## 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>135</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>135</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  $1 \times 10^{-4}$ .

## 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:

$$S(V_i) = (1 - d) + d * \sum_{V_j \in \text{In}(V_i)} \frac{1}{|\text{Out}(V_i)|} S(V_j) \quad \text{where } d \in \mathbb{R}[0, 1] \quad (669)$$

and the damping factor  $d$  is interpreted as the probability of jumping from a given vertex<sup>136</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>137</sup>
- (2) Iterate over vertices, computing equation 10.2 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:

$$\text{Error}^k(V_i) \approx S^k(V_i) - S^{k-1}(V_i) \quad (670)$$

<sup>136</sup>Note that  $d$  is a single parameter for the graph, i.e. it is the same for all vertices.

<sup>137</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 10.2 to incorporate *weighted* connections may be useful for NLP applications.

$$WS(V_i) = (1 - d) + d * \sum_{j \in \text{In}(V_i)} \frac{w_{ji}}{\sum_{V_k \in \text{Out}(V_j)} w_{jk}} WS(V_j) \quad (671)$$

$w_{ij}$  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.

---

### 10.2.1 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>138</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>139</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>140</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>138</sup>That is . . . simpler than expected. *Can we do better?*

<sup>139</sup>Another approach is to have  $T$  be a fixed value, where typically  $5 < T < 20$ .

<sup>140</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.

---

### 10.2.2 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  $S_j$ :

$$\text{Similarity}(S_i, S_j) = \frac{|S_i \cap S_j|}{\log(|S_i|) + \log(|S_j|)} \quad (672)$$

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).

## 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:

$$s(t) = \Pr [t \in y \mid t \in x] \quad (673)$$

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

$$s(t) = \frac{\text{Num pairs in which } t \text{ occurs in both}}{\text{Num pairs in which } t \text{ occurs in the first of the pair}} \quad (674)$$

*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:

$$E(t) = - \sum_{i=1}^n \sum_{j=1}^n (S_{ij} \cdot \log(S_{ij}) + (1 - S_{ij}) \cdot \log(1 - S_{ij})) \quad (675)$$

$$S_{ij} = 2^{-d_{ij}/\bar{d}} \quad (676)$$

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>141</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>142</sup>, assuming the data is mean-centered. Quick check/reminder:

$$(\mathbf{A}^T \mathbf{A})_{ij} = \mathbf{A}_{:,i}^T \mathbf{A}_{:,j} \triangleq \mathbf{a}_i^T \mathbf{a}_j \quad (677)$$

$$\approx n \cdot \mathbb{E} [\mathbf{a}_i \mathbf{a}_j] \quad (678)$$

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>143</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>141</sup>The difference between LSI and PCA is that PCA subtracts out the means, which destroys the sparseness of the design matrix.

<sup>142</sup>Approximation because it is based on our training data, not on true expectations over the underlying data-distribution.

<sup>143</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:

$$J = \frac{1}{2} \|\mathbf{A} - \mathbf{UV}^T\|_F^2 \quad (679)$$

$$= \frac{1}{2} \left( \text{tr}(\mathbf{AA}^T) - 2\text{tr}(\mathbf{AVU}^T) + \text{tr}(\mathbf{UV}^T \mathbf{VU}^T) \right) \quad (680)$$

$$\begin{aligned} u_{ij} &\geq 0 \\ v_{ij} &\geq 0 \end{aligned}$$

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

$$\mathbf{a}_i \approx \mathbf{u}_i \cdot \mathbf{V}^T \quad (681)$$

$$= \sum_{i=1}^k \mathbf{u}_i \mathbf{V}_i^T \quad (682)$$

Therefore, the document vector  $\mathbf{a}$  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>144</sup>:

$$\mathcal{L} = J + \text{tr}(\alpha \cdot \mathbf{U}^T) + \text{tr}(\beta \cdot \mathbf{V}^T) \quad (683)$$

$$\text{where } \text{tr}(\alpha \cdot \mathbf{U}^T) = \sum_{i=1}^n \alpha_i \cdot \mathbf{u}_i = \sum_{i=1}^n \sum_{j=1}^n \alpha_{ij} u_{ij} \quad (684)$$

Any matrix multiplication with 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>144</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>145</sup>:

$$\frac{\partial \mathcal{L}}{\partial U} = -\mathbf{A} \cdot \mathbf{V} + \mathbf{U} \cdot \mathbf{V}^T \cdot \mathbf{V} + \boldsymbol{\alpha} = 0 \quad (685)$$

$$\frac{\partial \mathcal{L}}{\partial V} = -\mathbf{A}^T \cdot \mathbf{U} + \mathbf{V} \cdot \mathbf{U}^T \cdot \mathbf{U} + \boldsymbol{\beta} = 0 \quad (686)$$

Since, ultimately, these just say [some matrix] = 0, we can multiply both sides (element-wise) by a constant ( $x \times 0 = 0$ ). Using<sup>146</sup> the **Kuhn-Tucker conditions**  $\alpha_{ij} \cdot u_{ij} = 0$  and  $\beta_{ij} \cdot v_{ij} = 0$ , we get:

$$(\mathbf{A} \cdot \mathbf{V})_{ij} \cdot u_{ij} - (\mathbf{U} \cdot \mathbf{V}^T \cdot \mathbf{V})_{ij} \cdot u_{ij} = 0 \quad (687)$$

$$(\mathbf{A}^T \cdot \mathbf{U})_{ij} \cdot v_{ij} - (\mathbf{V} \cdot \mathbf{U}^T \cdot \mathbf{U})_{ij} \cdot v_{ij} = 0 \quad (688)$$

- Update rules:

$$u_{ij} = \frac{(\mathbf{A} \cdot \mathbf{V})_{ij} \cdot u_{ij}}{(\mathbf{U} \cdot \mathbf{V}^T \cdot \mathbf{V})_{ij}} \quad (689)$$

$$v_{ij} = \frac{(\mathbf{A}^T \cdot \mathbf{U})_{ij} \cdot v_{ij}}{(\mathbf{V} \cdot \mathbf{U}^T \cdot \mathbf{U})_{ij}} \quad (690)$$

### 10.3.1 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>147</sup>.

<sup>145</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>146</sup>i.e. the equations that follow are *not* the KT conditions, they just *use/exploit* them...

<sup>147</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

### 10.3.2 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  $D_i$ :

$$\Pr [t_l | D_i] = \sum_{j=1}^k \Pr [t_l | T_j] \cdot \Pr [T_j | D_i] \quad (691)$$

which should be interpreted as a weighted average<sup>148</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:

$$J = -\log (\Pr [\mathbf{X}]) \quad (692)$$

$$= -\log \left( \prod_{i,l} \Pr [t_l | D_i]^{\mathbf{X}_{i,l}} \right) \quad (693)$$

$$= -\sum_{i,l} \mathbf{X}_{i,l} \cdot \log (\Pr [t_l | D_i]) \quad (694)$$

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 691 to for evaluating  $\Pr [t_l | D_i]$ . We want to optimize the value of  $J$ , subject to the constraints:

$$(\forall T_j) : \sum_l \Pr [t_l | T_j] = 1 \quad (\forall D_i) : \sum_j \Pr [T_j | D_i] = 1 \quad (695)$$

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>149</sup>. Typically preferred over PLSI because PLSI more prone to overfitting.

---

<sup>148</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>149</sup>LDA is the Bayesian version of PLSI

---

### 10.3.3 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>150</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>151</sup>, insert  $X$  to the cluster with maximum similarity.
  - Otherwise, a new cluster is created<sup>152</sup> containing the solitary data point  $X$ .

---

<sup>150</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>151</sup>We specify some max allowed number of clusters  $k$ .

<sup>152</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*.

## 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:

$$t_0 \text{ s.t. } f(t_0) = \frac{1}{2}f(0) \quad (696)$$

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:

$$\mathcal{D}(t, \mathcal{C}) \triangleq (\bar{DF}_2, \bar{DF}_1, n, w(t), l) \quad (697)$$

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>153</sup>.  $\sum_{i=1}^d \sum_{j \neq i}^d v_i v_j$  entries total<sup>154</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>153</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>154</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:

$$\mathcal{D}(t, \mathcal{C}) \triangleq (\bar{DF}2, \bar{DF}1, n, w(t), l) \quad (698)$$

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>155</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>155</sup>In other words, the statistics for a cluster do not decay, until a new point is added to it.

## 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.

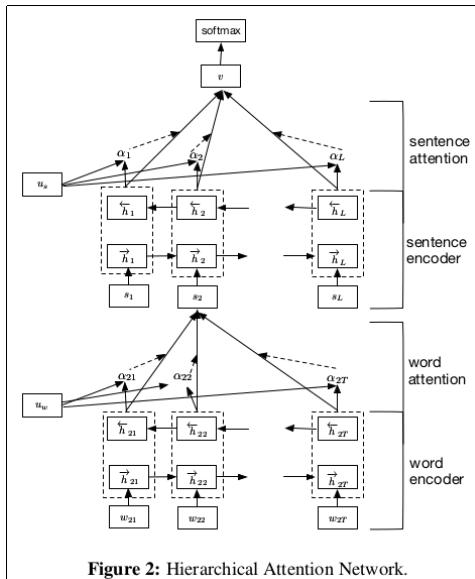


Figure 2: Hierarchical Attention Network.

The authors choose the GRU as their underlying RNN. For ease of reference, the defining equations of the GRU are shown below:

$$h_t = (1 - z_t) \odot h_{t-1} + z_t \odot \tilde{h}_t \quad (699)$$

$$z_t = \sigma(W_z x_t + U_z h_{t-1} + b_z) \quad (700)$$

$$\tilde{h}_t = \tanh(W_h x_t + r_t \odot (U_h h_{t-1}) + b_h) \quad (701)$$

$$r_t = \sigma(W_r x_t + U_r h_{t-1} + b_r) \quad (702)$$

**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.

$$u_{it} = \tanh(W_w h_{it} + b_w) \quad (703)$$

$$\alpha_{it} = \frac{\exp(u_{it}^T u_w)}{\sum_t \exp(u_{it}^T u_w)} \quad (704)$$

$$s_i = \sum_t \alpha_{it} h_{it} \quad (705)$$

Note the *context vector*  $u_w$ , which is shared for all words<sup>156</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.

$$u_i = \tanh(W_s h_i + b_s) \quad (706)$$

$$\alpha_i = \frac{\exp(u_i^T u_s)}{\sum_i \exp(u_i^T u_s)} \quad (707)$$

$$v = \sum_t \alpha_i h_i \quad (708)$$

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:

$$p = \text{softmax}(W_c v + b_c) \quad (709)$$

---

<sup>156</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:

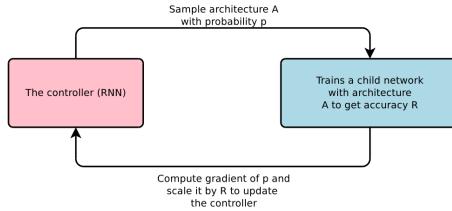
- **Tokenization:** Stanford CoreNLP. Vocabulary consists of words occurring more than 5 times, all others are replaced with UNK token.
- **Word Embeddings:** train word2vec on the training and validation splits. Dimension of 200.
- **GRU.** Dimension of 50 (so 100 because bidirectional).
- **Context 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.

## 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>157</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,  $J(\theta_c)$ :

$$J(\theta_c) = \mathbb{E}_{P(a_{1:T}; \theta_c)} [R] \quad (710)$$

Since the quantity  $\nabla_{\theta_c} R$  is non-differentiable<sup>158</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.

$$\nabla_{\theta_c} J(\theta_c) = R \sum_{t=1}^T \mathbb{E}_{P(a_{1:T}; \theta_c)} [\nabla_{\theta_c} \log P(a_t | a_{1:t-1}; \theta_c)] \quad (711)$$

$$\approx \frac{1}{m} \sum_{k=1}^m R_k \sum_{t=1}^T \nabla_{\theta_c} \log P(a_t | a_{1:t-1}; \theta_c) \quad (712)$$

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>159</sup>. Also note that we do

<sup>157</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>158</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>159</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 711 is an unbiased estimator for 712, 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:

$$\frac{1}{m} \sum_{k=1}^m \sum_{t=1}^T \nabla_{\theta_c} \log P(a_t | a_{1:t-1}; \theta_c) (R_k - b) \quad (713)$$

## 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:

$$\frac{d\mathbf{h}(t)}{dt} = f(\mathbf{h}(t), t, \theta) \quad (714)$$

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  $\mathbf{h}(t)$ .

$$d\mathbf{h}(t) = f(\mathbf{h}(t), t, \theta) dt \quad (715)$$

$$\int d\mathbf{h}(t) = \int f(\mathbf{h}(t), t, \theta) dt \quad (716)$$

$$\mathbf{h}(t) + c_1 = \int f(\mathbf{h}(t), t, \theta) dt \quad (717)$$

$$(718)$$

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

$$L(\mathbf{z}(t_1)) = L\left(\int_{t_0}^{t_1} f(\mathbf{z}(t), t, \theta) dt\right) \quad (719)$$

Given our starting definition (eq 714), we can say

$$\mathbf{z}(t + \epsilon) = \mathbf{z}(t) + \int_t^{t+\epsilon} f(\mathbf{z}(t), t, \theta) dt := T_\epsilon(\mathbf{z}(t), t) \quad (720)$$

which we can use to define the **adjoint**  $a(t)$ :

$$a(t) \triangleq -\frac{\partial L}{\partial z(t)} = -\frac{\partial L}{\partial z(t+\epsilon)} \frac{dz(t+\epsilon)}{dz(t)} \quad (721)$$

$$= a(t+\epsilon) \frac{\partial T_\epsilon(z(t), t)}{\partial z(t)} \quad (722)$$

$$\frac{da(t)}{dt} = -a(t)^T \frac{\partial f(z(t), t, \theta)}{\partial z} \quad (723)$$

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  $z(t_1)$ , parameters  $\theta$ , and gradient  $\frac{\partial L}{\partial z(t_1)}$  compute all gradients of  $L$ .*

1. Compute  $t_1$  gradient:

$$\frac{\partial L}{\partial t_1} = \frac{\partial L}{\partial z(t_1)}^T \frac{\partial z(t_1)}{\partial t_1} = \frac{\partial L}{\partial z(t_1)}^T f(z(t_1), t_1, \theta) \quad (724)$$

2. Initialize the augmented state:

$$s_0 := \left[ z(t_1), a(t_1), \mathbf{0}, -\frac{\partial L}{\partial t_1} \right] \quad (725)$$

3. Define augmented state dynamics:

$$\frac{ds}{dt} \triangleq \left[ f(z(t), t, \theta), -a(t)^T \frac{\partial f}{\partial z}, -a(t)^T \frac{\partial f}{\partial \theta}, -a(t)^T \frac{\partial f}{\partial t} \right] \quad (726)$$

4. Solve **reverse-time**<sup>a</sup> ODE:

$$\left[ z(t_0), \frac{\partial L}{\partial z(t_0)}, \frac{\partial L}{\partial \theta}, \frac{\partial L}{\partial t_0} \right] = \text{ODESolve}\left(s_0, \frac{ds}{dt}, t_1, t_0, \theta\right) \quad (727)$$

5. Return  $\frac{\partial L}{\partial z(t_0)}$ ,  $\frac{\partial L}{\partial \theta}$ ,  $\frac{\partial L}{\partial t_0}$ ,  $\frac{\partial L}{\partial t_1}$

---

<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.

## 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>160</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>161</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>162</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>160</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>161</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>162</sup>A single module alone cannot, because it repeats periodically. In other words, it can only provide relative location information.

## 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,

$$p(y=d | x, \omega) = \text{Cat}_d(\text{softmax}(f^\omega(x))) \quad (728)$$

for a classification setting. Given a dataset  $X, Y$ , and some new point  $x^*$ , we can predict its output via

$$p(y^* | x^*, X, Y) = \int p(y^* | x^*, \omega) p(\omega | X, Y) d\omega \quad (729)$$

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:

$$\text{KL}(q(\omega) || p(\omega | X, Y)) \propto - \int q(\omega) \log p(Y | X, \omega) d\omega + \text{KL}(q(\omega) || p(\omega)) \quad (730)$$

$$= \sum_{i=1}^N \int q(\omega) \log p(y_i | f^\omega(x_i)) d\omega + \text{KL}(q(\omega) || p(\omega)) \quad (731)$$

**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:

$$\mathcal{L} \approx - \sum_{i=1}^N \log p(y_i | f_y^{\hat{\omega}_i} (f_h^{\hat{\omega}_i}(x_{i,T}, h_{T-1}))) + \text{KL}(q(\omega) || p(\omega)) \quad (732)$$

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

$$q(w_k) \triangleq p\mathcal{N}(w_k; \mathbf{0}, \sigma^2 I) + (1-p)\mathcal{N}(w_k; m_k, \sigma^2 I) \quad (733)$$

with  $m_k$  **variational parameter** (row vector).

## 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  $S \subset \mathcal{Z}$ :

$$\frac{\Pr[M(x) \in S]}{\Pr[M(x') \in S]} \leq e^{\varepsilon} \quad (734)$$

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

$$U = \frac{\Delta W}{\|\Delta W\|_2} \quad (735)$$

$$R = \|\Delta W\|_2 \quad (736)$$

**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>163</sup>, we want to generate an  $\varepsilon$ -differentially private vector  $Z$  such that

$$\mathbb{E}[Z | u] = u \quad \forall u \in \mathbb{S}^{d-1} \quad (737)$$

### Privatized Unit Vector: $\text{PrivUnit}_2$

Sample random vector  $V$ :

$$V \sim \begin{cases} U(\{v \in \mathbb{S}^{d-1} | \langle v, u \rangle \geq \gamma\}) & \text{with probability } p \\ U(\{v \in \mathbb{S}^{d-1} | \langle v, u \rangle < \gamma\}) & \text{otherwise} \end{cases} \quad (738)$$

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

$$m = \frac{(1-\gamma^2)\alpha}{2^{d-2}(d-1)} \left[ \frac{p}{B(\alpha, \alpha) - B(\tau; \alpha, \alpha)} - \frac{1-p}{B(\tau; \alpha, \alpha)} \right] \quad (739)$$

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>163</sup>Here, this denotes an n-sphere:

$$\mathbb{S}^n \triangleq \{x \in \mathbb{R}^{n+1} : \|x\| = r\}$$

## 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:

$$MSE(q) = \mathbb{E}_{\mathbf{x} \sim p(\mathbf{x})} \left[ \|\mathbf{x} - \mathbf{q}(\mathbf{x})\|_2^2 \right] \quad (740)$$

In order for an optimizer to be optimal, it must satisfy the **Lloyd optimality conditions**:

$$(1) \quad \mathbf{q}(\mathbf{x}) = \arg \min_{c_i \in \mathcal{C}} \|\mathbf{x} - \mathbf{c}_i\|_2 \quad (741)$$

$$(2) \quad \mathbf{c}_i = \mathbb{E}_{\mathbf{x}} [\mathbf{x} \mid i] \triangleq \int_{\mathcal{V}_i} \mathbf{x} p(\mathbf{x}) d\mathbf{x} \quad (742)$$

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  $a \in \mathbb{Z}$ ).

$$\mathbf{x} = \underbrace{\mathbf{x}_{\langle 1 \dots D^* \rangle}}_{\mathbf{u}_1(\mathbf{x})}, \dots, \underbrace{\mathbf{x}_{\langle D-D^*+1 \dots D \rangle}}_{\mathbf{u}_m(\mathbf{x})} \rightarrow \mathbf{q}_1(\mathbf{u}_1(\mathbf{x})), \dots, \mathbf{q}_m(\mathbf{u}_m(\mathbf{x})) \quad (743)$$

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  $\mathcal{C} = \mathcal{C}_1 \times \dots \times \mathcal{C}_m$ .

## 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.”

$$F(\mathbf{x} = \{x_1, \dots, x_n\}) \in [0, 1] \quad (744)$$

$$A_F(\mathbf{x}, \mathbf{x}') = \{a_1, \dots, a_n\} \in \mathbb{R}^n \quad (745)$$

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>164</sup>.

*Given an input  $x$  and baseline  $x'$ , the **integrated gradient** along the  $i$ th dimension is as follows.*

$$IG_i(x, x') \triangleq (x_i - x'_i) \times \int_{\alpha=0}^1 \frac{\partial F(x' + \alpha \times (x - x'))}{\partial x_i} d\alpha \quad (746)$$

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>164</sup>The use the same context though (e.g. the associated image for VQA). Only the question is changed.

## Solving Rubik's Cube with a Robot Hand

Table of Contents Local

Written by Brandon McKinzie

Akkaya et al., "Solving Rubik's Cube with a Robot Hand" *OpenAI*, (Oct 2019).

**TL;DR: Automatic Domain Randomization** (ADR) + robot platform = solving rubik's cube from simulation alone. The two main structural components are (1) the hand and (2) a few cameras that observe the pose/state of the hand/cube.

**ADR** (5). They denote an “environment” as  $e_\lambda$ , parameterized by  $\lambda \in \mathbb{R}^d$ . In other words, there are  $d$  parameters/knobs they can fiddle with in simulation to change various characteristics of the given environment<sup>165</sup>. Here, the authors choose  $d' = 2d$ . Let  $\phi^L, \phi^H \in \mathbb{R}^d$  be some partition of  $\phi$ .

$$\begin{aligned}\lambda &\in \mathbb{R}^d \\ \phi &\in \mathbb{R}^{d'}\end{aligned}$$

$$P_\phi(\lambda) = \prod_{i=1}^d P_\phi(\lambda_i) = \prod_{i=1}^d U(\phi_i^L, \phi_i^H) \quad (747)$$

$$\mathcal{H}(P_\phi) = -\frac{1}{d} \int P_\phi(\lambda) \ln P_\phi(\lambda) d\lambda = \frac{1}{d} \sum_{i=1}^d \ln (\phi_i^H - \phi_i^L) \quad (748)$$

where they define a normalized entropy  $\mathcal{H}$  in nats/dimension. The pairs  $(\phi_i^L, \phi_i^H)$  are referred to as **boundary values**.

---

<sup>165</sup>I guess it comes down to semantics whether we interpret a  $\Delta\lambda$  as a new environment or just the same environment with different characteristics.

## Automatic Domain Randomization

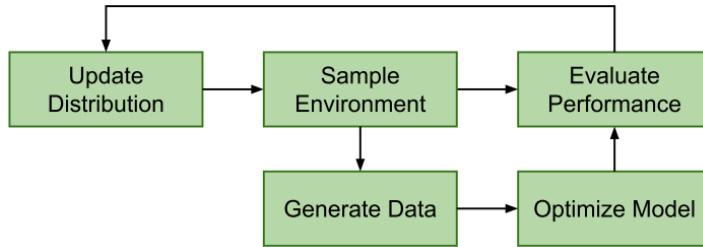
TL;DR: algorithm for updating parameter ranges  $\phi_i$  in distribution  $P_\phi$ .

Init with **performance buffers**  $\{D_i^L, D_i^H\}_{i=1}^d$  (presumably empty at start), **thresholds**  $m, t_L, t_H$  ( $t_L < t_H$ ). Repeat the following until “training is complete”:

1. Sample  $\lambda \sim P_\phi$ .
2. Randomly sample index  $1 \leq i \leq d$ , and number  $x \sim U(0, 1)$ .
3. If  $x < 0.5$ , set  $D_i = D_i^L$ ,  $\lambda_i = \phi_i^L$ , else  $D_i = D_i^H$ ,  $\lambda_i = \phi_i^H$ .
4. Collect model performance  $p$  on environment parameterized by  $\lambda$ .
5.  $D_i \leftarrow D_i \cup \{p\}$ .
6. If  $\text{Length}(D_i) \geq m$ :
  - (a)  $\bar{p} \leftarrow \text{AVERAGE}(D_i)$
  - (b)  $\text{CLEAR}(D_i)$
  - (c) if  $\bar{p} \geq t_H$  increase  $\phi_i$  by  $\Delta$ . If  $\bar{p} \leq t_L$ , decrease  $\phi_i$  by  $\Delta$ . (not sure if this means increase range scale, or shift range left/right (seems like the former))

In English: we sample environment params  $\lambda$  at each iteration, but then set a random dimension of  $\lambda$  to a boundary value. We measure model performance and append to performance buffer index associated with the boundary value. After we have enough performance measurements for a given boundary value index  $i$ , we increase the associated pair  $(\phi_i^L, \phi_i^H)$ .

To generate training data, they use the ADR algorithm above in conjunction with sampling from the [changing each time ADR is run] distribution  $P_\phi$  and running the model (running the model results in new training data<sup>166</sup>).




---

<sup>166</sup> **TODO:** how? how does running a model generate data? what is the data anyway?

## Decoupled Weight Decay Regularization

Table of Contents Local

Written by Brandon McKinzie

I. Loshchilov and F. Hutter, “Decoupled Weight Decay Regularization” *University of Freiburg*, (Jan 2019).

**Decoupling the Weight Decay from the Gradient-Based Update** (2). Denote learning rate as  $\alpha$ , rate of weight decay as  $\lambda$ , and batch-averaged loss at step  $t$  as  $f_t$ . Let  $f_t^{reg}(\theta) = f_t(\theta) + \frac{\lambda'}{2} \|\theta\|_2^2$  denote the loss function with L2 regularization, where  $\lambda' = \frac{\lambda}{\alpha}$ . The formal distinction between **weight decay** and **L2 regularization** is as follows:

$$[\text{weight decay}] \quad \theta_{t+1} = (1 - \lambda)\theta_t - \alpha \nabla f_t(\theta_t) \quad (749)$$

$$[\text{L2 reg}] \quad \theta_{t+1} = \theta_t - \alpha \nabla f_t^{reg}(\theta_t) \quad (750)$$

**Proposition 1**

*Weight decay (with  $\lambda$ ) is identical to L2 regularization (with  $\lambda' = \frac{\lambda}{\alpha}$ ) for standard SGD.*

The key realization here is that the optimal value of  $\alpha$  may change differently than the optimal weight decay value  $\lambda$ . The L2 regularization formula, however, has the weight decay strongly coupled with the learning rate, since  $\lambda' = \frac{\lambda}{\alpha}$  (a requirement for Proposition 1 to hold). The formula for weight decay above was proposed for decoupling this effect.

## SmartChoices

Table of Contents Local

*Written by Brandon McKinzie*

V. Carbune, T. Coppey, A. Daryin, T. Deselaers, N. Sarda, J. Yagnik, “SmartChoices: Hybridizing Programming and Machine Learning” *Google Research*, (June 2019).

**Software Development with SmartChoices** (2). Below is an example for creating a SmartChoice object, along with a common sequence of method invocations:

---

```
choice = SmartChoice(
    output_def=(float, shape=[1], range=[0, 1]),
    observation_defs={k: (float, [1], [0, 10] for k in ['low', 'high', 'target'])},
    initial_function=lambda observations: 0.5)
choice.Observe({'high': 0.56, 'target': 0.43})
value = choice.Predict()
# ...Do something with value that allows us to compute reward...
choice.Feedback(reward=10)
```

---

SmartChoice aims to maximize the sum of reward values received over time (possibly discounted).

**SmartChoices in Algorithms** (4). To evaluate the impact of SmartChoices, the authors measure **cumulative regret** over training episodes<sup>167</sup>

---

<sup>167</sup>Recall that an episode in RL is defined as a sequence of the form  $(s_0; a_1, r_1, s_1; a_2, r_2, s_2; \dots s_F)$ , starting with a start state and terminating with an end state.

## The Contextual Bandits Problem

Table of Contents   Local

*Written by Brandon McKinzie*

Robert Schapire, “The Contextual Bandits Problem” *Lecture at Microsoft Research*, (March 2017).

(Link to talk)

### Formal Model

Repeat the following for  $t \in [1..T]$ :

1. Learner observes **context**  $x_t$
2. **Reward** vector  $r_t \in [0, 1]^K$  chosen (but **not** observed)
3. Learner selects **action**  $a_t \in \{1, \dots, K\}$
4. Learner receives observed **reward**  $r_t(a_t)$

**Goal:** maximize the total reward:  $\sum_{t=1}^T r_t(a_t)$  relative to the best policy  $\pi \in \Pi$ , i.e. want small **regret**:

$$\max_{\pi \in \Pi} \frac{1}{T} \sum_{t=1}^T r_t(\pi(x_t)) - \frac{1}{T} \sum_{t=1}^T r_t(a_t) \quad (751)$$

For now, we assume pairs  $(x_t, r_t)$  are chosen at random iid.

Some miscellaneous additional points:

- Want to learn a good **policy**  $\pi : x \mapsto a$ .
- Before learning, we still have to decide the general form/space of possible policies that can be learned (by deciding general model architecture, etc).
- A **no-regret** learning algorithm is one that satisfies  $\text{regret} \rightarrow 0$  as  $T \rightarrow \infty$ .

**Starting Point: Full-Information Setting.** Same as bandit, but learner can see rewards for *all* actions (for the given context). Note that the learner does not see the rewards for all actions until it *after* it chooses an action – it first takes an action, then receives a reward and also is shown the rewards for all the other actions it could've taken. Similarly, after our learner makes a series of actions for various contexts, it can look back and compute the rewards it *would've* received for any other policy<sup>168</sup>.

### Follow-the-Leader algorithm

Optimal-regret algorithm for the full-information setting (and iid  $(x_t r_t)$ ). At each round  $t$ , do:

1. Find the *empirically best*  $\pi \in \Pi$ , the policy with the maximum average reward over all previously-observed contexts.
2. Use this policy to choose the next action:  $a_t = \pi(x_t)$ .

Has optimal regret  $\mathcal{O}\left(\sqrt{\ln |\Pi| / T}\right)$ .

---

<sup>168</sup> **Question:** is it more common to reflect after a series of N actions, or to do so after every action?

The main challenge of FTL is determining how to find the empirical best policy, which is the task of finding some single function  $\pi : x \mapsto a$  that would've resulted in maximum average reward over all previous observations<sup>169</sup>. Considering that the context and action spaces can be large, this is not trivial.

**Non-Stochastic (Adversarial) Setting.** If the  $(x_t, r_t)$  are *not* iid, then FTL will not work, especially if the contexts are provided by an adversary (e.g. in game playing). The **hedge algorithm** (below) has optimal regret even in adversarial setting, but time/space are linear  $|\Pi|$ .

### Hedge Algorithm

Maintain a real-valued *weight* for every  $\pi \in \Pi$ . On each round  $t$ :

1. Choose *random* policy  $\pi$  with probability proportional to weights
2. Use action chosen by  $\pi$
3. *Increase* weight of each policy according the reward it would've received

**Epoch-Greedy Algorithm.** Back to the bandit setting (and stochastic/iid) where we only observe rewards for the actions we take. Although we *could* attempt to just use the same algorithms as before – but only on actions/rewards we took of course – we end up with a lot of missing/biased information (since we can only evaluate other policies on some  $x_t$  when they take the same action as the policy we used). One simple improvement is to use the same algorithms as before (FTL) with probability  $1 - \epsilon$ , and to choose an action uniformly at random with probability  $\epsilon$ .

**De-biasing Biased Estimates.** Need to address the problem of selection bias. A simple trick is **inverse-propensity weighting**:

### Inverse-Propensity Weighting

**Goal:** estimate  $\mathbb{E}[X]$  for some RV  $X$ . With probability  $p$ , we observe  $X$  once, and with probability  $(1 - p)$ , we don't get to observe  $X$  at all (on each trial).

**Trick:** Define new RV  $\hat{X} = X/p$  if observed, else 0. Then  $\mathbb{E}[\hat{X}] = p\mathbb{E}[X/p] + (1-p)\mathbb{E}[0] = X$  – unbiased!

We can use this technique to get unbiased estimates for the rewards of *all* actions (not just observed). **Problem:** variance can be extremely large!

---

<sup>169</sup>Note that this is identical to the general task of finding a classifier that performs optimally on some training data

**Mini-Monster Algorithm**<sup>170</sup>. Each round, finds a weighted combination of policies satisfying:

1. Low (estimated) regret
2. Low (estimated) variance

---

<sup>170</sup>A.K.A. ILOVETOCONBANDITS

## The principles of adaptation in organisms and machines

Table of Contents   Local

*Written by Brandon McKinzie*

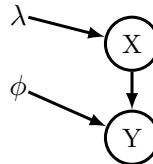
Hideaki Shimazaki, “The principles of adaptation in organisms and machines I: machine learning, information theory, and thermodynamics” *Graduate School of Informatics, Kyoto University*, (Feb 2019).

**Hierarchical Models** (2.1). The author defines the following [hierarchical] generative model for data  $\mathbf{Y}$  from the outside world:

$$p(\mathbf{y}, \mathbf{x}, \mathbf{w}) = p(\mathbf{y} | \mathbf{x}, \mathbf{w})p(\mathbf{x} | \mathbf{w})p(\mathbf{w}) \quad (752)$$

$$\text{[generative model]} \approx p(\mathbf{y}, \mathbf{x}; \mathbf{w}) = p(\mathbf{y} | \mathbf{x}; \mathbf{w})p(\mathbf{x}; \mathbf{w}) \quad (753)$$

factorized using the **observation model** and **prior** model, where the second line is the approximation the author makes henceforth: treat  $\mathbf{w}$  as a fixed hyperparameter instead of a value sampled from a distribution<sup>171</sup>.  $\mathbf{X}$  represents the underlying causative factors of the data  $\mathbf{Y}$ , while  $\mathbf{W}$  is basically a blanket set of factors containing any environmental factors/contexts underlying either/both of  $\mathbf{X}$  and  $\mathbf{Y}$ . We can further simplify by letting  $\mathbf{w} = (\lambda, \phi)$ , defined such that we have the following graphical model:



In generative modeling, the typical goal is to learn the marginalized likelihood, which we can then use to estimate underlying causes from data<sup>172</sup>:

$$\text{[marginal likelihood]} \quad p(\mathbf{y} | \mathbf{w}) = \int p(\mathbf{y}, \mathbf{x} | \mathbf{w}) d\mathbf{x} \quad (754)$$

$$\text{[posterior distribution]} \quad p(\mathbf{x} | \mathbf{Y}, \mathbf{w}) = \frac{p(\mathbf{Y} | \mathbf{x}, \mathbf{w})p(\mathbf{x} | \mathbf{w})}{p(\mathbf{Y} | \mathbf{w})} \quad (755)$$

<sup>171</sup>Note also that in sec 2.2. he ends up treating  $\mathbf{w}$  as trainable model parameters.

<sup>172</sup>I'll use the notation  $p(\mathbf{y} | \mathbf{w})$  interchangeably with  $p(\mathbf{y}; \mathbf{w})$ , etc.

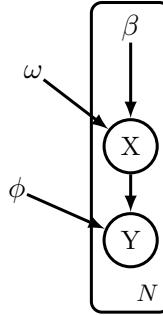
**Learning** (2.2). Author defines **Type II maximum likelihood estimation**: the process of determining parameters that maximize the marginal likelihood function  $p(\mathbf{Y} | \mathbf{w})$ :

$$\mathbf{W}^* = \arg \max_{\mathbf{w}} \log p(\mathbf{Y} | \mathbf{w}) \quad (756)$$

$$= \arg \max_{\mathbf{w}} \log \int p(\mathbf{Y}, \mathbf{x} | \mathbf{w}) d\mathbf{x} \quad (757)$$

(758)

The method for constructing  $p(\mathbf{x} | \mathbf{Y}, \mathbf{w})$  using the generative model with  $\mathbf{w} := \mathbf{W}^*$  is called the **empirical Bayes method**. As usual, the author points out the difficulty of obtaining the exact form of  $p(\mathbf{Y} | \mathbf{x}, \mathbf{w})$  and  $p(\mathbf{x} | \mathbf{w})$  (and thus also  $p(\mathbf{y} | \mathbf{w})$  and  $p(\mathbf{x} | \mathbf{Y}, \mathbf{w})$ ), and proceeds to introduce approximate methods. We'll use a slightly modified version of the graphical model from before:



We'll also define some **crucial assumptions** that will hold henceforth:

- Assumption 1: Observation model parameters  $\phi$  are learned from an entire set of observed data  $\{\mathbf{Y}_1, \dots, \mathbf{Y}_n\}$ .
- Assumption 2: Using a factorized set of prior distribution parameters  $\lambda \rightarrow \{\beta, \omega\}$ , assume ...
  - $\omega$  are optimized using all samples.
  - $\beta$  are optimized for each sample.

**Approximate Inference** (3). Recall that variational inference is the process of learning a parameterized function  $q(\mathbf{x} | \mathbf{Y}) \approx p(\mathbf{x} | \mathbf{Y}, \mathbf{w})$  to approximate the posterior. Quick review of the main equations:

$$\log p(\mathbf{Y} | \mathbf{w}) = \mathcal{L}[q, p] + D_{KL}(q || p) \quad (759)$$

$$\mathcal{L}[q, p] \triangleq \int q(\mathbf{x} | \mathbf{Y}) \log \frac{p(\mathbf{Y}, \mathbf{x} | \mathbf{w})}{q(\mathbf{x} | \mathbf{Y})} d\mathbf{x} \quad (760)$$

$$= \mathcal{Q}[q, p] + H[q] \quad (761)$$

$$\mathcal{Q}[q, p] \triangleq \int q(\mathbf{x} | \mathbf{Y}) \log p(\mathbf{Y}, \mathbf{x} | \mathbf{w}) d\mathbf{x} \quad (762)$$

with an illustration of the EM-algorithm below:

|                    |                       |                                                                                                                                                                              |                                                                                   |                                                |
|--------------------|-----------------------|------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|-----------------------------------------------------------------------------------|------------------------------------------------|
|                    | Inference<br>(E-step) | $\mathcal{L}[q, p] = \log p(\mathbf{Y}   \mathbf{w}) - KL[q(\mathbf{x}   \mathbf{Y})    p(\mathbf{x}   \mathbf{Y}, \mathbf{w})]$                                             |                                                                                   |                                                |
|                    |                       | <small>fixed</small><br><small>marginal log-likelihood</small>                                                                                                               | <small>optimized</small><br><small>KL between approx. and exact posterior</small> |                                                |
| <b>Lower bound</b> |                       |                                                                                                                                                                              |                                                                                   |                                                |
|                    | Learning<br>(M-step)  | $\mathcal{L}[q, p] = \int q(\mathbf{x}   \mathbf{Y}) \log p(\mathbf{Y}, \mathbf{x}   \mathbf{w}) d\mathbf{z} + H[q(\mathbf{x}   \mathbf{Y})]$                                | <small>optimized</small><br><small>expected complete data log-likelihood</small>  | <small>fixed</small><br><small>entropy</small> |
|                    |                       | <small>accuracy</small>                                                                                                                                                      |                                                                                   |                                                |
|                    |                       | $\mathcal{L}[q, p] = \int q(\mathbf{x}   \mathbf{Y}) \log p(\mathbf{Y}   \mathbf{x}, \mathbf{w}) d\mathbf{z} - KL[q(\mathbf{x}   \mathbf{Y})    p(\mathbf{x}   \mathbf{w})]$ | <small>optimized</small><br><small>complexity</small>                             |                                                |
|                    |                       |                                                                                                                                                                              |                                                                                   |                                                |

**Summary of Interpretations.** The author provides a lot of interpretations of formulas scattered throughout – below I've summarized them:

- **Posterior**  $p(\mathbf{x} | \mathbf{Y}, \mathbf{w})$ 
  - Neural activity  $\mathbf{X}$  caused by stimulus  $\mathbf{Y}$ .
  - Constructed by modulating the internal neural activity represented as the prior distribution by observing data.
- **Observation Model**  $p(\mathbf{y} | \mathbf{x}, \phi)$ :
  - Represents how data is expressed by combinations of neuronal activity.
- **Prior**  $p(\mathbf{x} | \lambda)$ :
  - Constraints on neural activity.
  - Spontaneous activity of neurons.

**Further Reading.** Some quotes and their citations from the paper that I want to investigate further.

- “It is likely that the brain implements functions similar to Bayesian inference” – K. Doya, Bayesian brain: Probabilistic approaches to neural coding. MIT press, 2007.
- “Algorithms for learning approximate inference model in the statistics and machine learning is instructive when we consider models of the brain.” K. Friston, “Learning and inference in the brain,” Neural Networks.
- “Parameters that represent spontaneous activity of neurons may be learned from many samples.” P. Berkes, G. Orban, M. Lengyel, and J. Fiser, “Spontaneous cortical activity reveals hallmarks of an optimal internal model of the environment,” Science, vol. 331, no. 6013, pp. 83-87, 2011.
- “Early sensory neurons receive feedback/recurrent modulation via top-down or lateral connections, which plays an essential role in perceptual experiences, attention, and reward modulation.” V. A. Lamme and P. R. Roelfsema, “The distinct modes of vision offered by feedforward and recurrent processing,” Trends in neurosciences, vol. 23, no. 11, pp. 571–579, 2000.