

CONTENTS

1	Math and Machine Learning Basics	4
1.1	Linear Algebra (Quick Review)	5
1.1.1	Example: Principal Component Analysis	7
1.2	Probability & Information Theory (Quick Review)	9
1.3	Numerical Computation	11
1.3.1	Gradient-Based Optimization (4.3)	11
1.4	Machine Learning Basics	13
1.4.1	Capacity, Overfitting, and Underfitting (5.2)	13
1.4.2	Estimators, Bias and Variance (5.4)	13
1.4.3	Bayesian Statistics (5.6)	15
1.4.4	Maximum Likelihood Estimation (5.5)	15
1.4.5	Supervised Learning Algorithms (5.7)	16
2	Deep Networks: Modern Practices	18
2.1	Deep Feedforward Networks	19
2.1.1	Back-Propagation (6.5)	20
2.2	Regularization for Deep Learning	21
2.2.1	Parameter Norm Penalties (7.1)	21
2.2.2	Sparse Representations (7.10)	22
2.2.3	Adversarial Training (7.13)	22
2.3	Convolutional Neural Networks	23
2.4	Sequence Modeling (RNNs)	24
2.4.1	Review: The Basics of RNNs	24
2.4.2	Encoder-Decoder Seq2Seq Architectures (10.4)	30
2.4.3	Challenge of Long-Term Deps. (10.7)	31
2.4.4	LSTMs and Other Gated RNNs (10.10)	32
2.5	Applications (Ch. 12)	33

2.5.1	Natural Language Processing (12.4)	33
2.5.2	Neural Language Models (12.4.2)	35
3	Deep Learning Research	36
3.1	Linear Factor Models (Ch. 13)	37
3.2	Autoencoders (Ch. 14)	39
3.3	Representation Learning (Ch. 15)	40
3.4	Monte Carlo Methods (Ch. 17)	41
4	Papers and Tutorials	44
4.1	Conv Nets: A Modular Perspective	45
4.2	Understanding Convolutions	46
4.3	Deep Reinforcement Learning	48
4.4	Deep Learning for Chatbots (WildML)	51
4.5	WaveNet (Paper)	53
4.6	Neural Style	57
4.7	Neural Conversation Model	59
4.8	NMT By Jointly Learning to Align & Translate	61
4.8.1	Detailed Model Architecture	62
4.9	Effective Approaches to Attention-Based NMT	64
4.10	Using Large Vocabularies for NMT	66
4.11	Candidate Sampling – TensorFlow	70
4.12	Attentional Interfaces – Neural Perspective	72
4.13	Attention Terminology	73
4.14	TextRank	76
4.14.1	Keyword Extraction	78
4.14.2	Sentence Extraction	79
4.15	Simple Baseline for Sentence Embeddings	80
4.16	Deep Models for Semantic Compositionality	82

4.17	Survey of Text Clustering Algorithms	83
4.17.1	Feature Selection and Transformation Methods	83
4.17.2	Distance-based Clustering Algorithms	87
5	NLP with Deep Learning	89
5.1	Word Vector Representations (Lec 2)	90
5.2	GloVe (Lec 3)	93
6	Speech and Language Processing	94
6.1	Naive Bayes and Sentiment (Ch. 6)	95
6.2	Vector Semantics (Ch. 15)	97
6.3	Semantics with Dense Vectors (Ch. 16)	100

MATH AND MACHINE LEARNING

BASICS

CONTENTS

1.1	Linear Algebra (Quick Review)	5
1.1.1	Example: Principal Component Analysis	7
1.2	Probability & Information Theory (Quick Review)	9
1.3	Numerical Computation	11
1.3.1	Gradient-Based Optimization (4.3)	11
1.4	Machine Learning Basics	13
1.4.1	Capacity, Overfitting, and Underfitting (5.2)	13
1.4.2	Estimators, Bias and Variance (5.4)	13
1.4.3	Bayesian Statistics (5.6)	15
1.4.4	Maximum Likelihood Estimation (5.5)	15
1.4.5	Supervised Learning Algorithms (5.7)	16

Linear Algebra (Quick Review): January 23

Table of Contents Local

Written by Brandon McKinzie

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

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

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

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

$$\|x\|_\infty = \max_i |x_i|$$

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

$$f(x + y) \leq f(x) + f(y) \quad (2)$$

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

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

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

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

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

- Suppose square matrix $A \in \mathbb{R}^{n \times n}$ has n linearly independent eigenvectors $\{v^{(1)}, \dots, v^{(n)}\}$. The **eigendecomposition** of A is then given by¹

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

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

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

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

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

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

- **Positive definite:** all λ are positive; **positive semidefinite:** all λ are positive or zero.

→ PSD: $\forall \mathbf{x}, \mathbf{x}^T \mathbf{A} \mathbf{x} \geq 0$

→ PD: $\mathbf{x}^T \mathbf{A} \mathbf{x} = 0 \Rightarrow \mathbf{x} = \mathbf{0}$.²

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

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

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

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

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

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

- The **singular values** are the diagonal entries \mathbf{D}_{ii} .
- The **left(right)-singular vectors** are the columns of $\mathbf{U}(\mathbf{V})$.
- Eigenvectors of $\mathbf{A}\mathbf{A}^T$ are the L-S vectors. Eigenvectors of $\mathbf{A}^T\mathbf{A}$ are the R-S vectors. The eigenvalues of both $\mathbf{A}\mathbf{A}^T$ and $\mathbf{A}^T\mathbf{A}$ are given by the singular values squared.
- The Moore-Penrose **pseudoinverse**, denoted \mathbf{A}^+ , enables us to find an “inverse” of sorts for a (possibly) non-square matrix \mathbf{A} . Most algorithms compute \mathbf{A}^+ via

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

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

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

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

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

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

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

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

EXAMPLE: PRINCIPAL COMPONENT ANALYSIS

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

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

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

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

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

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

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

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

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

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

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

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

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

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

Probability & Information Theory (Quick Review): January 24

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

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

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

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

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

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

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

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

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

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

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

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

where m is the number of data points.

Measure Theory.

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

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

Functions of RVs.

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

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

Numerical Computation: January 24

Table of Contents Local

Written by Brandon McKinzie

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

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

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

GRADIENT-BASED OPTIMIZATION (4.3)

Optimization algorithms that use only the gradient (e.g. SGD) are called 1st-order optimization algorithms. Likewise, ones using the Hessian matrix are 2nd-order.

Jacobian and Hessian Matrices. For when we want partial derivatives of some function f whose input and output are both vectors. The **Jacobian matrix** contains all such partial derivatives. Sometimes we want to know about second derivatives too, since this tells us whether a gradient step will cause as much of an improvement as we would expect based on the gradient alone. The **Hessian matrix** $\mathbf{H}(f)(\mathbf{x})$ is defined such that³

$$\begin{aligned} f &: \mathbb{R}^m \rightarrow \mathbb{R}^n \\ \mathbf{J} &\in \mathbb{R}^{n \times m} \text{ where} \\ J_{i,j} &= \frac{\partial}{\partial x_j} f(\mathbf{x})_i \end{aligned}$$

The Hessian is the Jacobian of the gradient.

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

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

³Recall that the directional derivative of $f(\mathbf{x})$ in direction \mathbf{u} is $\mathbf{u}^T \nabla_x f(\mathbf{x})$

shows us how \mathbf{H} can give information related to how “good” of a step that really was. Mathematically,

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

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

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

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

The best (and perhaps only) way to take what we learned about the “second derivative test” in single-variable calculus and apply it to the multidimensional case with \mathbf{H} is by using the *eigendecomposition of \mathbf{H}* . Why? Because we can examine the eigenvalues of the Hessian to determine whether the critical point $\mathbf{x}^{(0)}$ is a local maximum, local minimum, or saddle point⁴. If all eigenvalues are positive (remember that this is equivalent to saying that the Hessian is **positive definite!**), the point is a local minimum.

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

⁴Emphasis on “values” in “eigenvalues” because it’s important not to get tripped up here about what the eigenvectors of the Hessian mean. The reason for the decomposition is that it gives us an orthonormal basis (out of which we can get any direction) and therefore the magnitude of the second derivative along each of these directions as the eigenvalues.

Machine Learning Basics: January 25

Table of Contents Local

Written by Brandon McKinzie

CAPACITY, OVERFITTING, AND UNDERFITTING (5.2)

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

ESTIMATORS, BIAS AND VARIANCE (5.4)

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

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

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

Bias. Defined below, where the expectation is taken over the data-generating distribution⁵. Bias measures the expected deviation from the true value of the func/param.

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

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

Bias-Variance Tradeoff.

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

→ **Summary of main formulas.**

$$\text{bias} \left[\hat{\theta}_m \right] = \mathbb{E} \left[\hat{\theta}_m \right] - \theta \quad (20)$$

$$\text{Var} \left[\hat{\theta}_m \right] = \mathbb{E} \left[\left(\hat{\theta}_m - \mathbb{E} \left[\hat{\theta}_m \right] \right)^2 \right] \quad (21)$$

→ **MSE decomposition.** The MSE of the estimates is given by⁶

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

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

and desirable estimators are those with low MSE.

⁵May want to double-check this, but I'm fairly certain this is what the book meant when it said "data," based on later examples.

⁶Derivation:

$$\text{MSE} = \mathbb{E} \left[\hat{\theta}^2 + \theta^2 - 2\theta\hat{\theta} \right] \quad (22)$$

$$= \mathbb{E} \left[\hat{\theta}^2 \right] + \theta^2 - 2\theta\mathbb{E} \left[\hat{\theta} \right] \quad (23)$$

$$= (\mathbb{E} \left[\hat{\theta} \right]^2 - \mathbb{E} \left[\hat{\theta} \right]^2) + \mathbb{E} \left[\hat{\theta}^2 \right] + \theta^2 - 2\theta\mathbb{E} \left[\hat{\theta} \right] \quad (24)$$

$$= \left(\mathbb{E} \left[\hat{\theta} \right]^2 + \theta^2 - 2\theta\mathbb{E} \left[\hat{\theta} \right] \right) + \left(\mathbb{E} \left[\hat{\theta}^2 \right] - \mathbb{E} \left[\hat{\theta} \right]^2 \right) \quad (25)$$

$$= \text{Bias}(\hat{\theta})^2 + \text{Var} \left[\hat{\theta}_m \right] \quad (26)$$

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

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

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

BAYESIAN STATISTICS (5.6)

Distinction between frequentist and bayesian approach:

- **Frequentist:** Estimate $\theta \rightarrow$ make predictions thereafter based on this estimate.
- **Bayesian:** Consider all possible values of θ when making predictions.

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

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

MAXIMUM LIKELIHOOD ESTIMATION (5.5)

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

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

where we've chosen to express with \log for underflow/gradient reasons. One interpretation of ML is to view it as minimizing the dissimilarity, as measured by the KL divergence⁷, between \hat{p}_{data} and p_{model} .

⁷ The KL divergence is given by

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

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

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

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

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

Consider linear regression as viewed with ML: Instead of a single prediction value \hat{y} given input \mathbf{x} , think of the model as producing a *conditional distribution* $p(y \mid \mathbf{x})$ ⁸. To derive the standard linear regression algorithm, we *define*

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

Assume σ^2 is some fixed constant chosen by the user.

We can use this (and the i.i.d. assumption) to evaluate the conditional log-likelihood as

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

and we see that finding the \mathbf{w} that maximizes the conditional log-likelihood is equivalent to finding the \mathbf{w} that minimizes the training MSE.

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

SUPERVISED LEARNING ALGORITHMS (5.7)

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

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

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

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

Equation 5.81 is the definition of logistic regression

Unfortunately, there is no closed-form solution for θ , so we must search via maximizing the log-likelihood.

⁸After all, we might have several training points with the same value of \mathbf{x} but different labels y .

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

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

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

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

DEEP NETWORKS: MODERN PRACTICES

CONTENTS

2.1	Deep Feedforward Networks	19
2.1.1	Back-Propagation (6.5)	20
2.2	Regularization for Deep Learning	21
2.2.1	Parameter Norm Penalties (7.1)	21
2.2.2	Sparse Representations (7.10)	22
2.2.3	Adversarial Training (7.13)	22
2.3	Convolutional Neural Networks	23
2.4	Sequence Modeling (RNNs)	24
2.4.1	Review: The Basics of RNNs	24
2.4.2	Encoder-Decoder Seq2Seq Architectures (10.4)	30
2.4.3	Challenge of Long-Term Deps. (10.7)	31
2.4.4	LSTMs and Other Gated RNNs (10.10)	32
2.5	Applications (Ch. 12)	33
2.5.1	Natural Language Processing (12.4)	33
2.5.2	Neural Language Models (12.4.2)	35

Deep Feedforward Networks: January 26

Table of Contents Local

Written by Brandon McKinzie

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

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

with ϕ defining a hidden layer.

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

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

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

a.e. is short for “almost everywhere”

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

$$h_i = g(\mathbf{z}, \boldsymbol{\alpha})_i = \max(0, z_i) + \alpha_i \min(0, z_i) \quad (27)$$

→ Absolute value rectification: fix $\alpha_i = -1$ to obtain $g(z) = |z|$.

→ Leaky ReLU: fix α_i to a small value like 0.01.

→ Parametric ReLU (PReLU): treats α_i like a learnable parameter.

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

BACK-PROPAGATION (6.5)

The chain rule. Suppose $z = f(\mathbf{y})$ where $\mathbf{y} = g(\mathbf{x})$ (see margin for dimensions). Then⁹,

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

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

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

$$\mathbf{x} \in \mathbb{R}^m$$

$$\mathbf{y} \in \mathbb{R}^n$$

$$z : \mathbb{R}^n \rightarrow \mathbb{R}$$

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

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

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

Regularization for Deep Learning: January 12

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

PARAMETER NORM PENALTIES (7.1)

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

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

where we typically choose Ω to only penalize the *weights* and leave biases unregularized.

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

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

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

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

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

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

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

Sparse Representations (7.10)

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

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

Adversarial Training (7.13)

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

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

$$\mathbf{x}' \longleftarrow \mathbf{x} + \epsilon \cdot \text{sign}(\nabla_{\mathbf{x}} J(\boldsymbol{\theta}; \mathbf{x}, \mathbf{y})) \quad (28)$$

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

Convolutional Neural Networks: January 24

Table of Contents Local

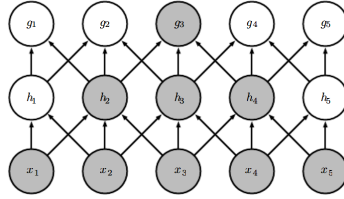
Written by Brandon McKinzie

We use a 2-D image I as our input (and therefore require a 2-D kernel K). Note that most neural networks do not technically implement convolution¹⁰, but instead implement a related function called the *cross-correlation*, defined as

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

Convolution leverages the following three important ideas:

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



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

¹⁰Technically the convolution output is defined as

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

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

where 9.5 can be asserted due to commutativity of convolution.

Sequence Modeling (RNNs): January 15

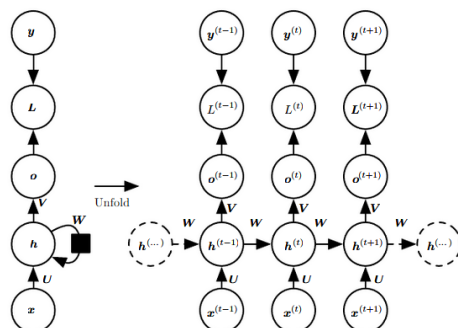
Table of Contents Local

Written by Brandon McKinzie

REVIEW: THE BASICS OF RNNs

Notation/Architecture Used.

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

Shape of $\mathbf{x}^{(t)}$ fixed, e.g. vocab length.Black square on recurrent connection \equiv interaction w/delay of a single time step.

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

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

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

Note that this is an example of an RNN that maps input seqs to output seqs of the same length¹¹. We can then compute, e.g., the log-likelihood loss $L = \sum_t L^{(t)}$ over all time

¹¹Where “same length” is related to the number of timesteps (i.e. τ input steps means τ output steps), not anything about the actual shapes/sizes of each individual input/output.

steps as:

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

Convince yourself this is identical to cross-entropy.

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

Back-Propagation Through Time.

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

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

$$\begin{aligned} (\nabla_{\mathbf{o}^{(t)}} L)_i &= \frac{\partial L}{\partial \mathbf{o}_i^{(t)}} \\ &= \frac{\partial L}{\partial L^{(t)}} \cdot \frac{\partial L^{(t)}}{\partial \mathbf{o}_i^{(t)}} \\ &= (1) \cdot \frac{\partial L^{(t)}}{\partial \mathbf{o}_i^{(t)}} \\ &= \frac{\partial}{\partial \mathbf{o}_i^{(t)}} \left\{ -\log \left(\hat{\mathbf{y}}_{y^{(t)}}^{(t)} \right) \right\} \\ &= -\frac{\partial}{\partial \mathbf{o}_i^{(t)}} \left\{ \log \left(\frac{e^{\mathbf{o}_{y^{(t)}}^{(t)}}}{\sum_j e^{\mathbf{o}_j^{(t)}}} \right) \right\} \\ &= -\frac{\partial}{\partial \mathbf{o}_i^{(t)}} \left\{ \mathbf{o}_{y^{(t)}}^{(t)} - \log \left(\sum_j e^{\mathbf{o}_j^{(t)}} \right) \right\} \\ &= -\left\{ \mathbf{1}_{i,y^{(t)}} - \frac{\partial}{\partial \mathbf{o}_i^{(t)}} \log \left(\sum_j e^{\mathbf{o}_j^{(t)}} \right) \right\} \end{aligned} \quad (29)$$

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

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

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

which leaves all entries of $\mathbf{o}^{(t)}$ unchanged *except* for the entry corresponding to the true label, which will become negative in the gradient. All this means is, since we want to increase the probability of this entry, driving this value up will *decrease* the loss (hence negative) and driving any other entries up will *increase* the loss proportional to its current estimated probability (driving up an [incorrect] entry that is already high is “worse” than driving up a small [incorrect entry]).

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

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

Before proceeding, **notice the following useful pattern:** If two nodes a and b , each containing n_a and n_b neurons, are fully connected by parameter matrix $W_{n_b \times n_a}$ and directed like $a \rightarrow b \rightarrow L$, then¹² $\nabla_a L = W^T \nabla_b L$. Using

¹²More generally,

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

this result, we can then iterate and take gradients back in time from $t = \tau - 1$ to $t = 1$ as follows:

$$\nabla_{\mathbf{h}^{(t)}} L = \left(\frac{\partial \mathbf{h}^{(t+1)}}{\partial \mathbf{h}^{(t)}} \right)^T (\nabla_{\mathbf{h}^{(t+1)}} L) + \left(\frac{\partial \mathbf{o}^{(t)}}{\partial \mathbf{h}^{(t)}} \right)^T (\nabla_{\mathbf{o}^{(t)}} L) \quad (10.20)$$

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

$$\frac{d}{dx} \tanh(x) = 1 - \tanh^2(x)$$

$$(\text{diag}(\mathbf{a}))_{ii} \triangleq a_i$$

2. Parameter Gradients. Now we can compute the gradients for the parameter matrices/vectors, where it is crucial to remember that a given parameter matrix (e.g. U) is shared across *all* time steps t . We can treat tensor derivatives in the same form as previously done with vectors after a quick abstraction: For any tensor \mathbf{X} of arbitrary rank (e.g. if rank-4 then index like \mathbf{X}_{ijkl}), use single variable (e.g. i) to represent the complete tuple of indices¹³.

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

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

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

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

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

¹³More details on tensor derivatives: Consider the chain defined by $\mathbf{Y} = g(\mathbf{X})$, and $z = f(\mathbf{Y})$, where z is some vector. Then

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

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

$$\nabla_{\mathbf{V}} L = \sum_t^{\tau} \nabla_{\mathbf{V}} L^{(t)} \quad (30a)$$

$$= \sum_t^{\tau} \nabla_{\mathbf{V}} L^{(t)}(\mathbf{o}_1^{(t)}, \dots, \mathbf{o}_{n_{out}}^{(t)}) \quad (30b)$$

$$= \sum_t^{\tau} \sum_i^{n_{out}} (\nabla_{\mathbf{o}^{(t)}} L)_i \nabla_{\mathbf{V}} \mathbf{o}_i^{(t)} \quad (30c)$$

$$= \sum_t^{\tau} \sum_i^{n_{out}} (\nabla_{\mathbf{o}^{(t)}} L)_i \nabla_{\mathbf{V}} \left\{ c_i + \sum_{j=1}^{n_{hid}} V_{ij} \mathbf{h}_j^{(t)} \right\} \quad (30d)$$

$$= \sum_t^{\tau} \sum_i^{n_{out}} (\nabla_{\mathbf{o}^{(t)}} L)_i \begin{bmatrix} 0 & 0 & \dots & 0 \\ \vdots & \vdots & \dots & \vdots \\ \mathbf{h}_1^{(t)} & \mathbf{h}_2^{(t)} & \dots & \mathbf{h}_{n_{hid}}^{(t)} \\ \vdots & \vdots & \dots & \vdots \\ 0 & 0 & \dots & 0 \end{bmatrix} \quad (30e)$$

NOTE: In equation 30e, all rows are 0 except the i th row, which is $(\mathbf{h}^{(t)})^T$. (See footnote for more)

$$= \sum_t^{\tau} (\nabla_{\mathbf{o}^{(t)}} L) (\mathbf{h}^{(t)})^T \quad (30f)$$

where if 30e confuses you, see the footnote¹⁴.

- \mathbf{W} ($n_{hid} \times n_{hid}$). This one is a bit odd, since \mathbf{W} is, in a sense, even more “shared” across time steps than \mathbf{V} ¹⁵. The authors here define/choose, when evaluating $\nabla_{\mathbf{W}} h_i^{(t)}$ to only concern themselves with $\mathbf{W} := \mathbf{W}^{(t)}$, i.e. the direct

¹⁴ The general lesson learned here is that, for some matrix $\mathbf{W} \in \mathbb{R}^{a \times b}$ and vector $\mathbf{x} \in \mathbb{R}^b$,

$$\sum_i \nabla_{\mathbf{W}} [(\mathbf{W}\mathbf{x})_i] = \begin{bmatrix} \mathbf{x}^T \\ \mathbf{x}^T \\ \vdots \\ \mathbf{x}^T \end{bmatrix} \quad (31)$$

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

¹⁵Specifically, $\mathbf{h}^{(t)}$ is both

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

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

connections to \mathbf{h} at time t .

$$\nabla_{\mathbf{w}} L = \sum_t^\tau \nabla_{\mathbf{w}} L^{(t)} \quad (32a)$$

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

$$= \sum_t^\tau \sum_i^{n_{hid}} (\nabla_{\mathbf{h}^{(t)}} L)_i \left(\text{diag} (1 - (\mathbf{h}^{(t)})^2) \begin{bmatrix} 0 & 0 & \dots & 0 \\ \vdots & \vdots & \dots & \vdots \\ \mathbf{h}_1^{(t-1)} & \mathbf{h}_2^{(t-1)} & \dots & \mathbf{h}_{n_{hid}}^{(t-1)} \\ \vdots & \vdots & \dots & \vdots \\ 0 & 0 & \dots & 0 \end{bmatrix} \right) \quad (32b)$$

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

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

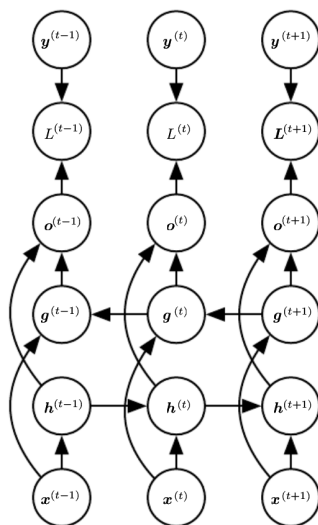
$$\nabla_{\mathbf{U}} L = \sum_t^\tau \nabla_{\mathbf{U}} L^{(t)} \quad (33a)$$

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

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

RNNs as Directed Graphical Models. The advantage of RNNs is their efficient parameterization of the joint distribution over $\mathbf{y}^{(i)}$ via parameter sharing. This introduces a built-in assumption that we can model the effect of $y^{(i)}$ in the distant past on the current $y^{(t)}$ *via its effect on \mathbf{h}* . We are also assuming that the conditional probability distribution over the variables at $t+1$ given the variables at time t is **stationary**. Next, we want to know how to draw *samples* from such a model. Specifically, how to sample from the conditional distribution ($y^{(t)}$ given $y^{(t-1)}$) at each time step. The author then only discusses methods of determining stopping conditions for sampling (pg. 380), and not really “how to draw samples,” which I suppose is self-explanatory.

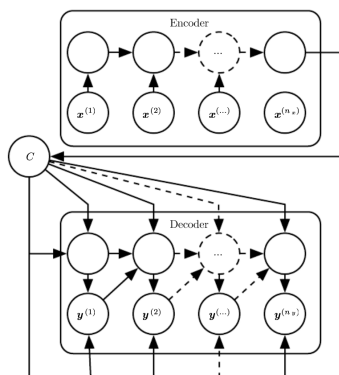
Bidirectional RNNs. In many applications, it is desirable to output a prediction of $\mathbf{y}^{(t)}$ that may depend on *the whole sequence*. For example, in speech recognition, the interpretation of words/sentences can also depend on what is *about* to be said. Below is a typical bidirectional RNN,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



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

ENCODER-DECODER SEQ2SEQ ARCHITECTURES (10.4)

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



CHALLENGE OF LONG-TERM DEPS. (10.7)

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

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

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

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

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

Any component of $\mathbf{h}^{(0)}$ that isn't aligned with the largest eigenvector will eventually be discarded.¹⁶

If, however, we have a non-recurrent network such that the state elements are repeatedly multiplied by different $w^{(t)}$ at each time step, the situation is different. Suppose the different $w^{(t)}$ are i.i.d. with mean 0 and variance v . The variance of the product is easily seen to be $\mathcal{O}(v^n)$ ¹⁷. *To obtain some desired variance v^* we may choose the individual weights with variance $v = \sqrt[n]{v^*}$.*

¹⁶Make sure to think about this from the right perspective. The largest value of $t = \tau$ in the RNNs we've seen would correspond with either (1) the largest output sequence or (2) the largest input sequence (if fixed-vector output). After we extract the output from a given forward pass, we reset the clock and either back-propagate errors (if training) or get ready to feed another sequence.

¹⁷Quick sketch of (my) proof:

$$\text{Var} [w^{(i)}] = v = \mathbb{E} [(w^{(i)})^2] - \cancel{\mathbb{E} [w^{(i)}]^2} \quad (34)$$

$$\text{Var} \left[\prod_t^n w^{(t)} \right] = \mathbb{E} \left[\left(\prod_t^n w^{(t)} \right)^2 \right] = \prod_t^n \mathbb{E} [(w^{(t)})^2] = v^n \quad (35)$$

LSTMS AND OTHER GATED RNNs (10.10)

While leaky units have connection weights that are either manually chosen constants or are trainable parameters, gated RNNs generalize this to connection weights that may change *at each time step*. Furthermore, gated RNNs can learn to both accumulate and *forget*, while leaky units are designed for just accumulation¹⁸

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

- **Forget gate** $f_i^{(t)}$.
- **Internal state** $s_i^{(t)}$.
- **External input gate** $g_i^{(t)}$.
- **Output gate** $q_i^{(t)}$.

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

¹⁸Q: Isn't choosing to update with higher relative weight on the present the same as forgetting? A: Sort of. It's like "soft forgetting" and will inevitably erase more/less than desired (smeary). In this context, "forget" means to set the weight of a specific past cell to zero.

Applications (Ch. 12): February 14

Table of Contents Local

Written by Brandon McKinzie

NATURAL LANGUAGE PROCESSING (12.4)

Begins on pg. 448

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

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

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

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

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

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

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

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

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

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

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

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

In addition, n -gram models are vulnerable to the curse of dimensionality, since most n -grams won't occur in the training set¹⁹, even for modest n .

¹⁹For a given vocabulary, which usually has much more than n possible words, consider how many possible sequences of length n .

NEURAL LANGUAGE MODELS (12.4.2)

Designed to overcome curse of dimensionality by using a distributed representation of words. Recognize that any model trained on sentences of length n and then told to generalize to new sentences [also of length n] must deal with a space²⁰ of possible sentences that is exponential in n . Such word representations (i.e. viewing words as existing in some high-dimensional space) are often called **word embeddings**. The idea is to map the words (or sentences) from the raw high-dimensional [vocab sized] space to a smaller feature space, where similar words are closer to one another. Using distributed representations may also be used with graphical models (think Bayes' nets) in the form multiple *latent variables*.

²⁰Ok I tried re-wording that from the book's confusing wording but that was also a bit confusing. Let me break it down. Say you train on a thousand sentences each of length 5. For a given vocabulary of size VOCAB.SIZE, the number of possible sequences of length 5 is $(\text{VOCAB.SIZE})^5$, which can be quite a lot more than a thousand (not to mention the possibility of duplicate training examples). To the naive model, all points in this high-dimensional space are basically the same. A neural language model, however, tries to arrange the space of possibilities in a meaningful way, so that an unforeseen sample at test time can be said "similar" as some previously seen training example. It does this by *embedding* words/sentences in a lower-dimensional feature space.

DEEP LEARNING RESEARCH

CONTENTS

3.1	Linear Factor Models (Ch. 13)	37
3.2	Autoencoders (Ch. 14)	39
3.3	Representation Learning (Ch. 15)	40
3.4	Monte Carlo Methods (Ch. 17)	41

Linear Factor Models (Ch. 13): January 12

Table of Contents Local

Written by Brandon McKinzie

Overview. Much research is in building a *probabilistic model*²¹ of the input, $p_{\text{model}}(x)$. Why? Because then we can perform *inference* to predict stuff about our environment given any of the other variables. We call the other variables **latent variables**, h , with

$$p_{\text{model}}(x) = \sum_h \Pr(h) \Pr(x \mid h) = \mathbb{E}_h [p_{\text{model}}(x \mid h)] \quad (36)$$

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

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

Note that \mathbf{h} is a *vector* of arbitrary size, where we assume $p(\mathbf{h})$ is a **factorial distribution**: $p(\mathbf{h}) = \prod_i p(h_i)$. This just means that the elements of \mathbf{h} are mutually independent²². The LFM describes the data-generation process as follows:

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

$$\mathbf{x} = \mathbf{W}\mathbf{h} + \mathbf{b} + \text{noise} \quad (37)$$

²¹Whereas, before, we've been building *functions* of the input (deterministic).

²²Note that, technically, this assumption isn't strictly the definition of mutual independence, which requires that every *subset* (i.e. not just the full set) of $\{h_i \in \mathbf{h}\}$ follow this factorial property.

Probabilistic PCA and Factor Analysis.

- Factor analysis:

$$\mathbf{h} \sim \mathcal{N}(\mathbf{h}; \mathbf{0}, \mathbf{I}) \quad (38)$$

$$\text{noise} \sim \mathcal{N}(\mathbf{0}, \boldsymbol{\psi} \equiv \text{diag}(\boldsymbol{\sigma}^2)) \quad (39)$$

$$\mathbf{x} \sim \mathcal{N}(\mathbf{x}; \mathbf{b}, \mathbf{W}\mathbf{W}^T + \boldsymbol{\psi}) \quad (40)$$

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

It is worth emphasizing the interpretation of $\boldsymbol{\psi}$ as the matrix of **conditional variances** σ_i^2 . *Huh?* Let's take a step back. The fact that we were able to separate the distributions in the above relations for \mathbf{h} and noise is from a built-in assumption that $\Pr(x_i | \mathbf{h}, x_{j \neq i}) = \Pr(x_i | \mathbf{h})$ ²³.

The Big Idea

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

- **Probabilistic PCA:** Just factor analysis with $\boldsymbol{\psi} = \sigma^2 \mathbf{I}$. So zero-mean spherical Gaussian noise. It becomes regular PCA as $\sigma \rightarrow 0$. Here we can use an iterative EM algorithm for estimating the parameters \mathbf{W} .

²³Due to $\langle \text{MATH} \rangle$, this introduces a constraint that knowing the value of some element x_j doesn't alter the probability $\Pr(x_i = W_i \cdot h + b_i + \text{noise})$. Given how we've defined the variable \mathbf{h} , this means that knowing noise_j provides no clues about noise_i . Mathematically, the noise must have a diagonal covariance matrix.

Autoencoders (Ch. 14): May 07

Table of Contents Local

Written by Brandon McKinzie

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

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

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

Representation Learning (Ch. 15): May 07

Table of Contents Local

Written by Brandon McKinzie

Greedy Layer-Wise Unsupervised Pretraining. Given:

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

The procedure is as follows:

1. Initialize.

$$f(\cdot) \leftarrow I(\cdot) \quad (41)$$

$$\tilde{\mathbf{X}} = \mathbf{X} \quad (42)$$

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

$$f^{(k)} = \mathcal{L}(\tilde{\mathbf{X}}) \quad (43)$$

$$f(\cdot) \leftarrow f^{(k)}(f(\cdot)) \quad (44)$$

$$\tilde{\mathbf{X}} \leftarrow f^{(k)}(\tilde{\mathbf{X}}) \quad (45)$$

In English: just apply the regular learning/training process for each layer/stage **sequentially and individually**²⁴.

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

²⁴In other words, you proceed one layer at a time *in order*. You don't touch layer i until the weights in layer $i - 1$ have been learned.

Monte Carlo Methods (Ch. 17): May 09

Table of Contents Local

Written by Brandon McKinzie

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

$$s = \int p(\mathbf{x})f(\mathbf{x})d\mathbf{x} = \mathbb{E}_p[f(\mathbf{x})] \quad \rightarrow \quad \hat{s}_n = \frac{1}{n} \sum_{i=1, \mathbf{x}^{(i)} \sim p}^n f(\mathbf{x}^{(i)}) \quad (46)$$

As we've seen before, the empirical average is an unbiased²⁵ estimator. Furthermore, the central limit theorem tells us that the distribution of \hat{s}_n converges to a normal distribution with mean s and variance $\text{Var}[f(\mathbf{x})]/n$.

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

$$p(\mathbf{x})f(\mathbf{x}) = q(\mathbf{x})\frac{p(\mathbf{x})f(\mathbf{x})}{q(\mathbf{x})} \quad (50)$$

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

²⁵ Recall that expectations on such an average are still taken over the underlying (assumed) probability distribution:

$$\mathbb{E}_p[\hat{s}_n] = \frac{1}{n} \sum_{i=1}^n \mathbb{E}_p[f(\mathbf{x}^{(i)})] \quad (47)$$

$$= \frac{1}{n} \sum_{i=1}^n s \quad (48)$$

$$= s \quad (49)$$

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

to denote the sampling distribution) as \hat{s}_q :

$$\hat{s}_q = \frac{1}{n} \sum_{i=1, \mathbf{x}^{(i)} \sim q}^n \frac{p(\mathbf{x}^{(i)})f(\mathbf{x}^{(i)})}{q(\mathbf{x}^{(i)})} \quad (51)$$

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

$$\text{Var}[\hat{s}_q] = \text{Var}\left[\frac{p(\mathbf{x})f(\mathbf{x})}{q(\mathbf{x})}\right] / n \quad (52)$$

with the optimal (minimum) value of q at:

$$q^* = \frac{p(\mathbf{x}) | f(\mathbf{x}) |}{Z} \quad (53)$$

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

$$\hat{s}_{BIS} = \frac{\sum_{i=1}^n \frac{\tilde{p}(\mathbf{x}^{(i)})}{\tilde{q}(\mathbf{x}^{(i)})} f(\mathbf{x}^{(i)})}{\sum_{i=1}^n \frac{\tilde{p}(\mathbf{x}^{(i)})}{\tilde{q}(\mathbf{x}^{(i)})}} \quad (54)$$

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

PAPERS AND TUTORIALS

CONTENTS

4.1	Conv Nets: A Modular Perspective	45
4.2	Understanding Convolutions	46
4.3	Deep Reinforcement Learning	48
4.4	Deep Learning for Chatbots (WildML)	51
4.5	WaveNet (Paper)	53
4.6	Neural Style	57
4.7	Neural Conversation Model	59
4.8	NMT By Jointly Learning to Align & Translate	61
4.8.1	Detailed Model Architecture	62
4.9	Effective Approaches to Attention-Based NMT	64
4.10	Using Large Vocabularies for NMT	66
4.11	Candidate Sampling – TensorFlow	70
4.12	Attentional Interfaces – Neural Perspective	72
4.13	Attention Terminology	73
4.14	TextRank	76
4.14.1	Keyword Extraction	78
4.14.2	Sentence Extraction	79
4.15	Simple Baseline for Sentence Embeddings	80
4.16	Deep Models for Semantic Compositionality	82
4.17	Survey of Text Clustering Algorithms	83
4.17.1	Feature Selection and Transformation Methods	83
4.17.2	Distance-based Clustering Algorithms	87

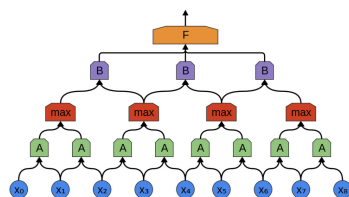
Conv Nets: A Modular Perspective: December 21

Table of Contents Local

Written by Brandon McKinzie

From this post on Colah's Blog.

The title is inspired by the following figure. Colah mentions how groups of neurons, like A , that appear in multiple places are sometimes called **modules**, and networks that use them are sometimes called modular neural networks. You can feed the output of one convolutional layer into another. With each layer, the network can detect higher-level, more abstract features.



- Function of the A neurons: compute certain *features*.
- Max pooling layers: kind of “zoom out”. They allow later convolutional layers to work on larger sections of the data. They also make us invariant to some very small transformations of the data.

Understanding Convolutions: December 21

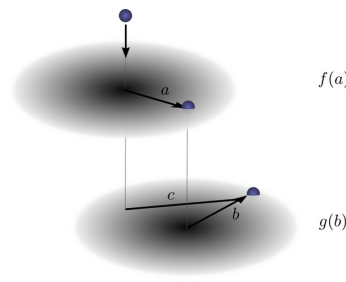
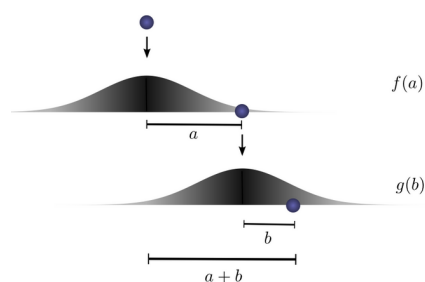
Table of Contents Local

Written by Brandon McKinzie

From Colah's Blog.

Ball-Dropping Example. The posed problem:

Imagine we drop a ball from some height onto the ground, where it only has one dimension of motion. How likely is it that a ball will go a distance c if you drop it and then drop it again from above the point at which it landed?



From basic probability, we know the result is a sum over possible outcomes, constrained by $a + b = c$. It turns out this is actually the definition of the convolution of f and g .

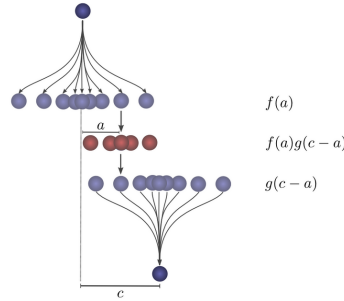
$$\Pr(a + b = c) = \sum_{a+b=c} f(a) \cdot g(b) \quad (55)$$

$$(f * g)(c) = \sum_{a+b=c} f(a) \cdot g(b) \quad (56)$$

$$= \sum_a f(a) \cdot g(c - a) \quad (57)$$

Visualizing Convolutions. Keeping the same example in the back of our heads, consider a few interesting facts.

- **Flipping directions.** If $f(x)$ yields the probability of landing a distance x away from where it was dropped, what about the probability that it was dropped a distance x from where it *landed*? Apparently²⁶ it is $f(-x)$.



- Above is a visualization of one term in the summation of $(f * g)(c)$. It is meant to show how we can move the bottom around to think about evaluating the convolution for different c values.

We can relate these ideas to image recognition. Below are two common kernels used to convolve images with.

0	0	0	0	0
0	1/9	1/9	1/9	0
0	1/9	1/9	1/9	0
0	1/9	1/9	1/9	0
0	0	0	0	0

0	0	0	0	0
0	0	0	0	0
0	-1	1	0	0
0	0	0	0	0
0	0	0	0	0

On the left is a kernel for *blurring* images, accomplished by taking simple averages. On the right is a kernel for *edge detection*, accomplished by taking the difference between two pixels, which will be largest at edges, and essentially zero for similar pixels.

²⁶Not entirely sold on the generalization of this, or even how true it is here.

Deep Reinforcement Learning: December 23

Table of Contents Local

Written by Brandon McKinzie

[Link to tutorial](#) – Part I of “Demystifying deep reinforcement learning.”

Reinforcement Learning. Vulnerable to the *credit assignment problem* - i.e. unsure which of the preceding actions was responsible for getting some reward and to what extent. Also need to address the famous *explore-exploit dilemma* when deciding what strategies to use.

Markov Decision Process. Most common method for representing a reinforcement problem. MDPs consist of states, actions, and rewards. Total reward is sum of current (includes previous) and *discounted* future rewards:

$$R_t = r_t \gamma (r_{t+1} + \gamma (r_{t+2} + \dots)) = r_t + \gamma R_{t+1} \quad (58)$$

Q - learning. Define function $Q(s, a)$ to be best possible score at end of game after performing action a in state s ; the “quality” of an action from a given state. The recursive definition of Q (for one transition) is given below in the *Bellman equation*.

$$Q(s, a) = r + \gamma \max_{a'} Q(s', a')$$

and updates are computed with a learning rate α as

$$Q(s_t, a_t) = (1 - \alpha) \cdot Q(s_{t-1}, a_{t-1}) + \alpha \cdot (r + \gamma \max_{a'} Q(s'_{t+1}, a'_{t+1}))$$

Deep Q Network. Deep learning can take deal with issues related to prohibitively large state spaces. The implementation chosen by DeepMind was to represent the Q-function with a neural network, with the states (pixels) as the input and Q-values as output, where the number of output neurons is the number of possible actions from the input state. We can optimize with simple squared loss:

$$L = \frac{1}{2} [\underbrace{r + \max_{a'} Q(s', a')}_{\text{target}} - \underbrace{Q(s, a)}_{\text{prediction}}]^2$$

and our algorithm from some state s becomes

1. **First forward pass** from s to get all predicted Q-values for each possible action. Choose action corresponding to max output, leading to next s' .
2. **Second forward pass** from s' and again compute $\max_{a'} Q(s', a')$.
3. **Set target output** for each action a' from s' . For the action corresponding to max (from step 2) set its target as $r + \gamma \max_{a'} Q(s', a')$, and for all other actions set target to same as originally returned from step 1, making the error 0 for those outputs. (Interpret as update to our guess for the best Q-value, and keep the others the same.)
4. **Update weights** using backprop.

Experience Replay. This the most important trick for helping convergence of Q-values when approximating with non-linear functions. During gameplay all the experience $\langle s, a, r, s' \rangle$ are stored in a replay memory. When training the network, random minibatches from the replay memory are used instead of the most recent transition.

Exploration. One could say that initializing the Q-values randomly and then picking the max is essentially a form of exploitation. However, this type of exploration is *greedy*, which can be tamed/fixed with **ϵ -greedy exploration**. This incorporates a degree of randomness when choosing next action at *all* time-steps, determined by probability ϵ that we choose the next action randomly. For example, DeepMind decreases ϵ over time from 1 to 0.1.

Deep Q-Learning Algorithm.

```
initialize replay memory  $D$ 
initialize action-value function  $Q$  with random weights
observe initial state  $s$ 
repeat
    select an action  $a$ 
        with probability  $\epsilon$  select a random action
        otherwise select  $a = \operatorname{argmax}_{a'} Q(s, a')$ 
    carry out action  $a$ 
    observe reward  $r$  and new state  $s'$ 
    store experience  $\langle s, a, r, s' \rangle$  in replay memory  $D$ 

    sample random transitions  $\langle ss, aa, rr, ss' \rangle$  from replay memory  $D$ 
    calculate target for each minibatch transition
        if  $ss'$  is terminal state then  $tt = rr$ 
        otherwise  $tt = rr + \gamma \max_{a'} Q(ss', aa')$ 
    train the  $Q$  network using  $(tt - Q(ss, aa))^2$  as loss

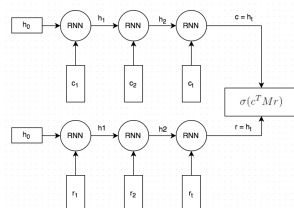
     $s = s'$ 
until terminated
```

Deep Learning for Chatbots (WildML): January 15

Table of Contents Local

*Written by Brandon McKinzie***Overview.**

- **Model.** Implementing a retrieval-based model. Input: conversation/context c . Output: response r .
 - **Data.** Ubuntu Dialog Corpus (UDC). 1 million examples of form (context, utterance, label). The label can be 1 (utterance was actual response to the context) or a 0 (utterance chosen randomly). Using NLTK, the data has been . . .
 - **Tokenized.** dividing strings into lists of substrings.
 - **Stemmed.** **IDK**
 - **Lemmatized.** **IDK**
- The test/validation set consists (context, ground-truth utterance, [9 distractors (incorrect utterances)]). The distractors are picked at random²⁷

Dual-Encoder LSTM.

1. **Inputs.** Both the context and the response text are split by words, and each word is embedded into a vector and fed into the same RNN.
2. **Prediction.** Multiply the [vector representation ("meaning")] c with param matrix M to predict some response r' .
3. **Evaluation.** Measure similarity of predicted r' to actual r via simple dot product. Feed this into sigmoid to obtain a probability [of r' being the correct response]. Use (binary) cross-entropy for loss function:

$$L = -y \cdot \ln(y') - (1 - y) \cdot \ln(1 - y') \quad (59)$$

²⁷Better example/approach: Google's Smart Reply uses clustering techniques to come up with a set of possible responses.

where y' is the predicted probability that r' is correct response r , and $y \in \{0, 1\}$ is the true label for the context-response pair (c, r) .

Data Pre-Processing. Courtesy of WildML, we are given 3 files after preprocessing: `train.tfrecords`, `validation.tfrecords`, and `test.tfrecords`, which use TensorFlow's 'Example' format. Each Example consists of . . .

- `context`: Sequence of word ids.
- `context.len`: length of the aforementioned sequence.
- `utterance`: seq of word ids representing utterance (response).
- `utterance.len`.
- `label`: only in training data. 0 or 1.
- `distractor_[N]`: Only in test/validation. N ranges from 0 to 8. Seq of word ids reppin the distractor utterance.
- `distractor_[N].len`.

WaveNet (Paper): January 15

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

- **What's WaveNet?** A deep neural network for generating raw audio waveforms.
- **How does it generate them?** **IDK**
- **What's it good for?** Text-to-speech, generating music, any waveform really.

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,

$$p(x) = \prod_{t=1}^T p(x_t \mid x_1, \dots, x_{t-1}) \quad (60)$$

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

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

Main ingredient of WaveNet is *dilated* causal convolutions, illustrated below. Note the absence of recurrent connections, which makes them faster to train than RNNs, but at the cost of requiring many layers, or large filters to increase the receptive field²⁸.

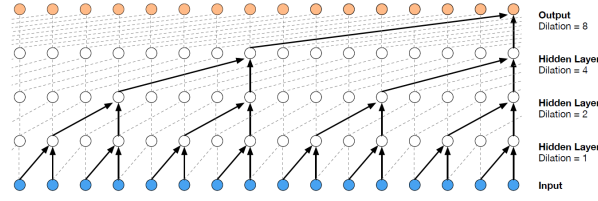


Figure 3: Visualization of a stack of *dilated* causal convolutional layers.

Excellent concise definition from paper:

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

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

$$f(x_t) = \text{sign}(x_t) \frac{\ln(1 + \mu|x_t|)}{\ln(1 + \mu)} \quad (61)$$

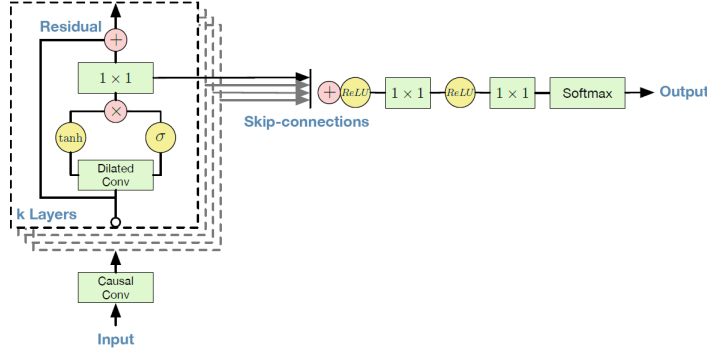
which (after plotting in Wolfram) looks identical to the sigmoid function.

²⁸Loose 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”

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

$$z = \tanh(W_{f,k} * x) \odot \sigma(W_{g,k} * x) \quad (62)$$

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



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

$$p(x \mid h) = \prod_{t=1}^T p(x_t \mid x_1, \dots, x_{t-1}, h) \quad (63)$$

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

$$z = \tanh(W_{f,k} * x + V_{f,k}^T h) \odot \sigma(W_{g,k} * x + V_{g,k}^T h) \quad (64)$$

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

Experiments.

- **Multi-Speaker Speech Generation.** Dataset: multi-speaker corpus of 44 hours of data from 109 different speakers²⁹. Receptive field of 300 milliseconds.
- **Text-to-Speech.** Single-speaker datasets of 24.6 hours (English) and 34.8 hours (Chinese) speech. Locally conditioned on *linguistic features*. Receptive field of 240 milliseconds. Outperformed both LSTM-RNN and HMM.
- **Music.** Trained the WaveNets to model two music datasets: (1) 200 hours of annotated music audio, and (2) 60 hours of solo piano music from youtube. Larger receptive fields sounded more musical.
- **Speech Recognition.** “With WaveNets we have shown that layers of dilated convolutions allow the receptive field to grow longer in a much cheaper way than using LSTM units.”

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

²⁹Speakers encoded as ID in form of a one-hot vector

Neural Style: January 22

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)³⁰.
- **Reconstructed image:** \mathbf{x} (initially random noise). Denote its corresponding filter response matrix at layer l as P^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} .

³⁰If 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

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

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

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

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 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 α/β , they explored 1×10^{-3} and 1×10^{-4} .

Neural Conversation Model: February 8

Table of Contents Local

Written by Brandon McKinzie

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

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

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

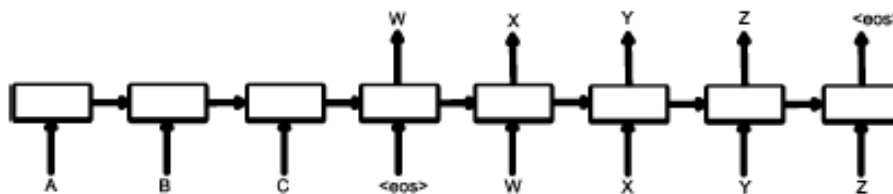
Related Work. The authors' approach is based on the following (linked and saved) papers on seq2seq:

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

Model. Succinctly described by the authors:

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

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



The **thought vector** is the hidden state of the model when it receives [as input] the end of sequence symbol $\langle eos \rangle$, because it stores the info of the sentence, or *thought*, “ABC”. The authors acknowledge that this model will *not* be able to “solve” the problem of modeling dialogue due to the objective function not capturing the actual objective achieved through human communication, which is typically longer term and based on exchange of information [rather than next step prediction]³¹.

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.

³¹I’d imagine that, in order to model human conversation, one obvious element needed would be a *memory*. Reminds me of DeepMind’s DNC. There would need to be some online filtering & output process to capture the crucial aspects/info to store in memory for later, and also some method of retrieving them when needed later. The method for retrieval would likely be some inference process where, given a sequence of inputs, the probability of them being related to some portion of memory could be trained. This would allow for conversations that stretch arbitrarily back in the past. Also, when storing the memories, I’d imagine a reasonable architecture would be some encoder-decoder for a sparse distributed representation of memory.

NMT By Jointly Learning to Align & Translate: February 27

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 Align³² and translate.

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

$$p(y_i \mid y_1, \dots, y_{i-1}, \mathbf{x}) = g(y_{i-1}, s_i, c_i) \quad (65)$$

$$s_i = f(s_{i-1}, y_{i-1}, c_i) \quad (66)$$

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:

$$c_i = \sum_{j=1}^{T_x} \alpha_{ij} h_j \quad (67)$$

$$\alpha_{ij} = \frac{\exp(e_{ij})}{\sum_{k=1}^{T_x} \exp(e_{ik})} \quad (68)$$

$$e_{ij} = a(s_{i-1}, h_j) \quad (69)$$

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.

³²By “align” the authors are referring to aligning the source-search to the relevant parts for prediction.

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

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_{i-1} , we now pass a *concatenated* state, $[s_{i-1}, c_i]$, 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:**

$$z_t = \sigma(W_z y_{t-1} + U_z[s_{t-1}, c_t]) \quad \text{[update gate]} \quad (70)$$

$$r_t = \sigma(W_r y_{t-1} + U_r[s_{t-1}, c_t]) \quad \text{[reset gate]} \quad (71)$$

$$(72)$$

3. **Proposal state:**

$$\tilde{s}_t = \tanh(W y_{t-1} + U[r_t \circ s_{t-1}, c_t]) \quad (73)$$

4. **Hidden state:**

$$s_t = (1 - z_t) \circ s_{t-1} + z_t \circ \tilde{s}_t \quad (74)$$

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:

$$a_j = a(s_{t-1}, h_j) = v_a^T \tanh(W_a s_{t-1} + U_a h_j) \quad (75)$$

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

$$a_j \leftarrow \frac{e^{a_j}}{\sum_{k=0}^{L_{enc}-1} e^{a_k}} \quad (76)$$

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

$$c = \sum_{j=1}^{L_{enc}} a_j h_j \quad (77)$$

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]:

$$\Pr(y \mid s, c) \propto e^{y^T W_o u} \quad (78)$$

$$u = [\max\{\tilde{u}_{2j-1}, \tilde{u}_{2j}\}]_{j=1, \dots, \ell}^T \quad (79)$$

$$\tilde{u} = U_o[s_{t-1}, c] + V_o y_{t-1} \quad (80)$$

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

Effective Approaches to Attention-Based NMT: May 11

[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:

$$\tilde{\mathbf{h}}_t = \tanh(\mathbf{W}_c[\mathbf{c}_t; \mathbf{h}_t]) \quad (81)$$

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

$$p(y_t \mid y_{<t}, x) = \text{softmax}(\mathbf{W}_s \tilde{\mathbf{h}}_t) \quad (82)$$

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:

$$\text{score}(\mathbf{h}_t, \bar{\mathbf{h}}_s) = \begin{cases} \mathbf{h}_t^T \bar{\mathbf{h}}_s & \text{dot} \\ \mathbf{h}_t^T \mathbf{W}_a \bar{\mathbf{h}}_s & \text{general} \\ \mathbf{v}_a^T \tanh(\mathbf{W}_a[\mathbf{h}_t; \bar{\mathbf{h}}_s]) & \text{concat} \end{cases} \quad (83)$$

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

$$\mathbf{a}_t(s) = \text{align}(\mathbf{h}_t, \bar{\mathbf{h}}_s) \quad (84)$$

$$= \frac{\exp(\text{score}(\mathbf{h}_t, \bar{\mathbf{h}}_s))}{\sum_{s'} \exp(\text{score}(\mathbf{h}_t, \bar{\mathbf{h}}_{s'}))} \quad (85)$$

3. \mathbf{c}_t : The weighted average over all source (encoder) hidden states³³:

$$\mathbf{c}_t = \sum_{i=1}^{L_{enc}} \mathbf{a}_t(i) \bar{\mathbf{h}}_i \quad (86)$$

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

$$\tilde{\mathbf{h}}_t = \tanh(\mathbf{W}_c[\mathbf{c}_t; \mathbf{h}_t]) \quad (87)$$

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 $[\tilde{\mathbf{h}}_t; \mathbf{x}_{t+1}]$.

³³NOTE: 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 \mathbf{a}_t and for long sentences, we ignore words near the end.*

Using Large Vocabularies for NMT: March 11

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:

$$x = (x_1, \dots, x_T) \quad [\text{source sentence}] \quad (88)$$

$$h = (h_1, \dots, h_T) \quad [\text{encoder state seq}] \quad (89)$$

$$h_t = f(x_t, h_{t-1}) \quad (90)$$

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:

$$y = (y_1, \dots, y'_T) \quad [y_i \in \mathbb{Z}] \quad (91)$$

$$\Pr[y_t \mid y_{<t}, x] \propto e^{q(y_{t-1}, z_t, c_t)} \quad (92)$$

$$z_t = g(y_{t-1}, z_{t-1}, c_t) \quad [\text{decoder hidden?}] \quad (93)$$

$$c_t = r(z_{t-1}, h_1, \dots, h_T) \quad [\text{decoder inp?}] \quad (94)$$

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,

$$\theta^* = \arg \max_{\theta} \sum_{n=1}^N \sum_{t=1}^{T_n} \log (\Pr[y_t^n \mid y_{<t}^n, x^n]) \quad (95)$$

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 .

$$c_t = \sum_{i=1}^T \alpha_i h_i \quad (96)$$

$$\alpha_t = \frac{e^{a(h_t, z_{t-1})}}{\sum_k e^{a(h_k, z_{t-1})}} \quad (97)$$

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 92 as³⁴,

$$\Pr[y_t \mid y_{<t}, x] = \frac{1}{Z} e^{\mathbf{w}_t^T \phi(y_{t-1}, z_t, c_t) + b_t} \quad (98)$$

$$Z = \sum_{k: y_k \in V} e^{\mathbf{w}_k^T \phi(y_{t-1}, z_t, c_t) + b_k} \quad (99)$$

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

where ϕ 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-likelihood³⁵, written in terms of the energy \mathcal{E} .

$$\nabla \log (\Pr[y_t \mid y_{<t}, x]) = \nabla \mathcal{E}(y_t) - \sum_{k: y_k \in V} \Pr[y_k \mid y_{<t}, x] \nabla \mathcal{E}(y_k) \quad (100)$$

$$\mathcal{E}(y_j) = \mathbf{w}_j^T \phi(y_{t-1}, z_t, c_t) + b_j \quad (101)$$

³⁴Note: 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.

³⁵**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 $Pr(y \mid y_{<t}, x)$. 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' .

$$\mathbb{E}_P [\nabla \mathcal{E}(y)] \approx \sum_{k: y_k \in V'} \frac{\omega_k}{\sum_{k': y_{k'} \in V'} \omega_{k'}} \nabla \mathcal{E}(y_k) \quad (102)$$

$$\omega_k = e^{\mathcal{E}(y_k) - \log Q(y_k)} \quad (103)$$

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

$$\omega_k = \frac{e^{\mathcal{E}(y_k)}}{Q(y_k)} \quad \text{thus} \quad p(y_k \mid y_{<t}, x) = \omega_k \frac{Q(y_k)}{Z} \quad (104)$$

$$\rightarrow e^{\mathcal{E}(y_k)} = Z \cdot p(y_k \mid y_{<t}, x) = Q(y_k) \cdot \omega_k \quad (105)$$

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.

$$\mathbb{E}_P [\nabla \mathcal{E}(y)] = \sum_{k: y_k \in V} \frac{\omega_k \cdot Q(y_k)}{\sum_{k': y_{k'} \in V} \omega_{k'} \cdot Q(y_{k'})} \nabla \mathcal{E}(y_k) \quad (106)$$

$$\mathbb{E}_P [\nabla \mathcal{E}(y)] = \sum_{k: y_k \in V'} \frac{\omega_k}{\sum_{k': y_{k'} \in V'} \omega_{k'}} \nabla \mathcal{E}(y_k) \quad (107)$$

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 τ . 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: March 19

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:

$$\Pr(y_i^{true} = y \mid C_i, x_i) = \frac{\Pr(y_i^{true} = y \mid x_i) \cdot \Pr(C_i \mid y_i^{true} = y, x_i)}{\Pr(C_i \mid x_i)} \quad (108)$$

$$= \frac{\Pr(y \mid x_i)}{Q(y \mid x_i)} \cdot \frac{1}{K(x_i, C_i)} \quad (109)$$

where they've just defined

$$K(x_i, C_i) \triangleq \frac{\Pr(C_i \mid x_i)}{\prod_{y' \in C_i} Q(y' \mid x_i) \prod_{y' \in (L - C_i)} (1 - Q(y' \mid x_i))} \quad (110)$$

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

$$F(x, y) = \log(\Pr(y \mid x)) + K(x) \quad (111)$$

which is the general form because

$$\text{Softmax}(\log(\Pr(y \mid x)) + K(x)) = \frac{e^{\log(\Pr(y \mid x)) + K(x)}}{\sum_{y'} e^{\log(\Pr(y' \mid x)) + K(x)}} \quad (112)$$

$$= \Pr(y \mid x) \quad (113)$$

Note that I've been a little sloppy here, since $\Pr(y \mid x)$ 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 111 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 109. We can easily satisfy that by subtracting $\log(Q(y \mid x))$ 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

$$\text{Softmax Input: } F(x, y) - \log(Q(y \mid x)) \quad (114)$$

instead. That's it.

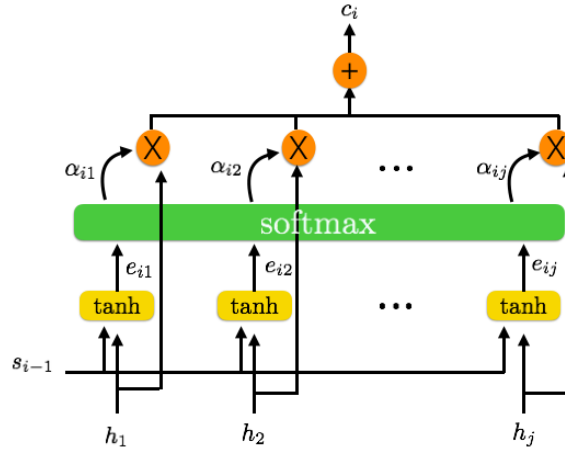
Attentional Interfaces – Neural Perspective: April 04

Table of Contents Local

Written by Brandon McKinzie

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

Attention Mechanism. Below is a close-up view/diagram of an attention layer. Technically, it only corresponds to a single time step i ; we are using the previous decoder state s_{i-1} to compute the i th context vector c_i which will be fed as an input to the decoder for step i .



For convenience, I'll rewrite the familiar equations for computing quantities at some step i .

$$\text{[decoder state]} \quad s_i = f(s_{i-1}, y_{i-1}, c_i) \quad (115)$$

$$\text{[context vect]} \quad c_i = \sum_{j=1}^{T_x} \alpha_{ij} h_j \quad (116)$$

$$\text{[annotation weights]} \quad \alpha_{ij} = \frac{\exp(e_{ij})}{\sum_{k=1}^{T_x} \exp(e_{ik})} \quad (117)$$

$$e_{ij} = a(s_{i-1}, h_j) \quad (118)$$

Now we can see just how simple this really is. Recall that Bahdanau *et al.*, 2015 use the wording: “ e_{ij} is an alignment model which scores how well the inputs around position j and the output at position i match.” But we can see an example implementation of an alignment model above: the \tanh function (that's it).

Attention Terminology: April 04

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

$$e_i = v^T \tanh(\text{FC}(s_{i-1}) + \text{FC}(h)) \quad (119)$$

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 α 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 deserves³⁶

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.
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:

$$\text{score}(\mathbf{h}_t, \bar{\mathbf{h}}_s) = \begin{cases} \mathbf{h}_t^T \bar{\mathbf{h}}_s & \text{dot} \\ \mathbf{h}_t^T \mathbf{W}_a \bar{\mathbf{h}}_s & \text{general} \\ \mathbf{v}_a^T \tanh(\mathbf{W}_a [\mathbf{h}_t; \bar{\mathbf{h}}_s]) & \text{concat} \end{cases} \quad (120)$$

Synonyms:
- scores
- unnormalized alignments

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

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

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

Synonyms:
- softmax outputs
- attention dist.
- alignments

$$\mathbf{a}_t(s) = \text{align}(\mathbf{h}_t, \bar{\mathbf{h}}_s) \quad (121)$$

$$= \frac{\exp(\text{score}(\mathbf{h}_t, \bar{\mathbf{h}}_s))}{\sum_{s'} \exp(\text{score}(\mathbf{h}_t, \bar{\mathbf{h}}_{s'}))} \quad (122)$$

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

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

Synonyms:
- context vector
- attention

³⁶For 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.

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:

$$\tilde{\mathbf{h}}_t = \tanh(\mathbf{W}_c [\mathbf{c}_t, \mathbf{h}_t]) \quad (123)$$

$$p(y_t \mid y_{<t}, x) = \text{softmax}(\mathbf{W}_s \tilde{\mathbf{h}}_t) \quad (124)$$

TextRank: May 03

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.

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

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:

$$S(V_i) = (1 - d) + d * \sum_{V_j \in \text{In}(V_i)} \frac{1}{|\text{Out}(V_j)|} S(V_j) \quad \text{where } d \in \mathbb{R}[0, 1] \quad (125)$$

The factor d is usually set to 0.85.

and the damping factor d is interpreted as the probability of jumping from a given vertex³⁷ to another random vertex in the graph. In practice, the algorithm is implemented through the following steps:

- (1) Initialize all vertices with arbitrary values.³⁸
- (2) Iterate over vertices, computing equation 4.14 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:

$$\text{Error}^k(V_i) \approx S^k(V_i) - S^{k-1}(V_i) \quad (126)$$

³⁷Note that d is a single parameter for the graph, i.e. it is the same for all vertices.

³⁸The 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!

Weighted Graphs. In contrast with the PageRank model, here we are concerned with natural language texts, which may include multiple or partial links between the units (vertices). The authors hypothesize that modifying equation 4.14 to incorporate *weighted* connections may be useful for NLP applications.

$$WS(V_i) = (1 - d) + d * \sum_{j \in \text{In}(V_i)} \frac{w_{ji}}{\sum_{V_k \in \text{Out}(V_j)} w_{jk}} WS(V_j) \quad (127)$$

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.

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 words³⁹. 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 $T = |V|/3$.⁴⁰ 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**⁴¹. 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.

³⁹That is ... simpler than expected. *Can we do better?*

⁴⁰Another approach is to have T be a fixed value, where typically $5 < T < 20$.

⁴¹ 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.

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 :

$$\text{Similarity}(S_i, S_j) = \frac{|S_i \cap S_j|}{\log(|S_i|) + \log(|S_j|)} \quad (128)$$

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

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

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

Simple Baseline for Sentence Embeddings: June 12, 2017

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:

$$\frac{1}{N} \sum_i^N \frac{a}{a + p(\mathbf{w}_i)} \mathbf{w}_i \quad (129)$$

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** $c_t \in \mathfrak{R}^d$ (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:

$$\Pr[w_t | c_t] \propto e^{c_t \cdot v_w} \quad (130)$$

- **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 MAP⁴² 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 130.**
 1. Additive term $\alpha p(w)$ where α is a scalar. Allows words to occur even if $c_t \cdot v_w$ is very small.
 2. Common discourse vector $c_0 \in \mathfrak{R}^d$. 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 (?)):

$$\Pr[w | c_s] = \alpha p(w) + (1 - \alpha) \frac{e^{\tilde{c}_s \cdot v_w}}{Z_{\tilde{c}_s}} \quad (131)$$

$$\tilde{c}_s = \beta c_0 + (1 - \beta) c_s \quad (132)$$

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

Embedding Calculations. Here I'll go into more depth on how the algorithm described earlier is derived. **TODO**

⁴²Review of MAP:

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

Deep Models for Semantic Compositionality: June 22, 2017

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

Overview. Full title: “Recursive Deep Models for Semantic Compositionality Over a Sentiment Treebank,” by Socher et al., 2014. Authors introduce the Sentiment Treebank, a dataset containing sentiment labels for 215,154 phrases in the parse trees of 11,855 sentences. They also introduce a model called the **Recursive Neural Tensor Network**. **TODO** Finish note on this . . .

Survey of Text Clustering Algorithms: July 03, 2017

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 AND TRANSFORMATION METHODS

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:

See ref 94 of the paper for more.

$$s(t) = \Pr[t \in y \mid t \in x] \quad (133)$$

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

$$s(t) = \frac{\text{Num pairs in which } t \text{ occurs in both}}{\text{Num pairs in which } t \text{ occurs in the first of the pair}} \quad (134)$$

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:

$$E(t) = - \sum_{i=1}^n \sum_{j=1}^n (S_{ij} \cdot \log(S_{ij}) + (1 - S_{ij}) \cdot \log(1 - S_{ij})) \quad (135)$$

$$S_{ij} = 2^{-d_{ij}/\bar{d}} \quad (136)$$

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, except that we use an approximation of the covariance matrix C which is appropriate for the sparse and high-dimensional nature of text data.

See ref 28. of paper for more on LSI.

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 matrix⁴³, assuming the data is mean-centered. Quick check/reminder:

$$(\mathbf{A}^T \mathbf{A})_{ij} = \mathbf{A}_{:,i}^T \mathbf{A}_{:,j} \triangleq \mathbf{a}_i^T \mathbf{a}_j \quad (137)$$

$$\approx n \cdot \mathbb{E} [\mathbf{a}_i \mathbf{a}_j] \quad (138)$$

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?**).

⁴³Approximation because it is based on our training data, not on true expectations over the underlying data-distribution.

As usual, we using the eigenvectors of $\mathbf{A}^T \mathbf{A}$ with the largest variance in order to represent the text⁴⁴. 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.

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:

$$J = \frac{1}{2} \|\mathbf{A} - \mathbf{UV}^T\|_F^2 \quad (139)$$

$$= \frac{1}{2} (tr(\mathbf{AA}^T) - 2tr(\mathbf{AVU}^T) + tr(\mathbf{UV}^T \mathbf{VU}^T)) \quad (140)$$

$$u_{ij} \geq 0$$

$$v_{ij} \geq 0$$

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} , of \mathbf{A} (document vector), this is

$$\mathbf{a} \approx \mathbf{u} \cdot \mathbf{V}^T \quad (141)$$

$$= \sum_{i=1}^k \mathbf{u}_i \mathbf{V}_i^T \quad (142)$$

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 .

⁴⁴In typical collections, only about 300 to 400 eigenvectors are required for the representation.

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 α and β , with same dimensions as \mathbf{U} and \mathbf{V} , respectively.
- Lagrangian: I'll just show it here first, and then explain in this footnote⁴⁵:

$$\mathcal{L} = J + \text{tr}(\alpha \cdot \mathbf{U}^T) + \text{tr}(\beta \cdot \mathbf{V}^T) \quad (143)$$

$$\text{where } \text{tr}(\alpha \cdot \mathbf{U}^T) = \sum_{i=1}^n \alpha_i \cdot \mathbf{u}_i = \sum_{i=1}^n \sum_{j=1}^n \alpha_{ij} u_{ij} \quad (144)$$

Any matrix multiplication with a \cdot is just a reminder to think of the matrices as column vectors.

You should think of α 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.

- Optimization: Set the partials of \mathcal{L} w.r.t both U and V (separately) to zero⁴⁶:

$$\frac{\partial \mathcal{L}}{\partial \mathbf{U}} = -\mathbf{A} \cdot \mathbf{V} + \mathbf{U} \cdot \mathbf{V}^T \cdot \mathbf{V} + \alpha = 0 \quad (145)$$

$$\frac{\partial \mathcal{L}}{\partial \mathbf{V}} = -\mathbf{A}^T \cdot \mathbf{U} + \mathbf{V} \cdot \mathbf{U}^T \cdot \mathbf{U} + \beta = 0 \quad (146)$$

Since, ultimately, these just say [some matrix] = 0, we can multiply both sides (element-wise) by a constant ($x \times 0 = 0$). Using⁴⁷ the **Kuhn-Tucker conditions** $\alpha_{ij} \cdot u_{ij} = 0$ and $\beta_{ij} \cdot v_{ij} = 0$, we get:

$$(\mathbf{A} \cdot \mathbf{V})_{ij} \cdot u_{ij} - (\mathbf{U} \cdot \mathbf{V}^T \cdot \mathbf{V})_{ij} \cdot u_{ij} = 0 \quad (147)$$

$$(\mathbf{A}^T \cdot \mathbf{U})_{ij} \cdot v_{ij} - (\mathbf{V} \cdot \mathbf{U}^T \cdot \mathbf{U})_{ij} \cdot v_{ij} = 0 \quad (148)$$

- Update rules:

$$u_{ij} = \frac{(\mathbf{A} \cdot \mathbf{V})_{ij} \cdot u_{ij}}{(\mathbf{U} \cdot \mathbf{V}^T \cdot \mathbf{V})_{ij}} \quad (149)$$

$$v_{ij} = \frac{(\mathbf{A}^T \cdot \mathbf{U})_{ij} \cdot v_{ij}}{(\mathbf{V} \cdot \mathbf{U}^T \cdot \mathbf{U})_{ij}} \quad (150)$$

⁴⁵ Recall that in Lagrangian minimization, \mathcal{L} takes the form of [the-function-to-be-minimized] + λ ([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].

⁴⁶ 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.

⁴⁷ i.e. the equations that follow are *not* the KT conditions, they just *use/exploit* them...

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 similarity⁴⁸.
- **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

K-Medoid isn't great for clustering text, especially short texts.

⁴⁸Here, 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.

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).

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

NLP WITH DEEP LEARNING

CONTENTS

5.1	Word Vector Representations (Lec 2)	90
5.2	GloVe (Lec 3)	93

Word Vector Representations (Lec 2): May 08

Table of Contents Local

Written by Brandon McKinzie

Meaning of a word. Common answer is to use a *taxonomy* like WordNet that has hypernyms (is-a) relationships. Problems with this discrete representation: misses nuances, e.g. the words in a set of synonyms can actually have quite different meanings/connotations. Furthermore, viewing words as atomic symbols is a bit like using one-hot vectors of words in a vocabulary space (inefficient).

Distributed representations. Want a way to encode word vectors such that two similar words have a similar structure/representation. The **distributional similarity-based**⁴⁹ approach represents words by means of its *neighbors* in the sentences in which it appears. You end up with a dense “vector for each word type, chosen so that it is good at predicting other words appearing in its context.”

Skip-gram prediction. Given a word w_t at position t in a sentence, learn to predict [probability of] some number of surrounding words, given w_t . Standard minimization with negative log-likelihood:

$$J(\theta) = -\frac{1}{T} \sum_{t=1}^T \sum_{-m \leq j \leq m} \log \Pr(w_{t+j} | w_t) \quad (151)$$

$$\Pr(o \mid c) = \frac{e^{\mathbf{u}_o^T \mathbf{v}_c}}{\sum_{w=1}^{\text{vocab size}} e^{\mathbf{u}_w^T \mathbf{v}_c}} \quad (152)$$

I cannot believe this actually works.

where

- The params θ are the vector representation of the words (they are the *only* learnable parameters here).
- m is our radius/window size.
- o and c are indices into our vocabulary (somewhat inconsistent notation).
- Yes, they are using different vector representations for \mathbf{u} (context words) and \mathbf{v} (center words). I’m assuming one reason this is done is because it makes the model architecture simpler/easier to build.

⁴⁹Note that this is distinct from the way “distributed” is meant in “distributed representation.” In contrast, distributional similarity-based representations refers to the notion that you can describe the meaning of words by understanding the context in which they appear.

Some subtleties:

- Looks like e.g. $Pr(w_{t+j} \mid w_t)$ doesn't really care what the value of j is, it is just modeling the probability that it is *somewhere* in the context window. The w_t are one-hot vectors into the vocabulary. Standard tricks for simplifying the cross-entropy loss apply.
- Equation 152 should be interpreted as the probability that the o^{th} word in our vocabulary occurs in the context window of the c^{th} word in our vocabulary.
- The model architecture is *identical* to an autoencoder. However, the (big) difference is the training procedure and interpretation of the model parameters going “in” versus the parameters going “out”.

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

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

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

$$\frac{1}{N} \sum_i^N \frac{a}{a + p(\mathbf{w}_i)} \mathbf{w}_i \quad (153)$$

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

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

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

1

Algorithm 1 Sentence Embedding

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

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

```

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

Further Reading.

- Learning representations by back-propagating errors (Rumelhard et al., 1986)
- A Neural Probabilistic Language Model (Bengio et al., 2003)
- NLP (almost) from Scratch (Collobert & Watson, 2008)
- Word2Vec (Miklov et al. 2013)

GloVe (Lec 3): May 08

Table of Contents Local

Written by Brandon McKinzie

Skip-gram and negative sampling. Main idea:

- Split the loss function from last lecture into two (additive) terms corresponding to the numerator and denominator respectively (you’ve done this a trillion times).
- The second term is an expectation over all the words in your vocab space. That is huge, so instead we only use a subsample of size k (the negative samples).
- Interpretation: First term is maximizing $\Pr(o | c)$, the probability that the true outside word (given by index o) occurs given context (index) c . Second term is minimizing the probability of random words (the negative samples) occurring around the center (context) word given by c .

To sample the negative samples, draw from $P(w) = U(w)^{3/4} / Z$, where U is the **unigram distribution**.

GloVe (Global Vectors). Given some co-occurrence matrix we computed with previous methods, we can use the following GloVe loss function over all pairs of co-occurring words in our matrix:

$$J(\theta) = \sum_{i,j=1}^W f(P_{ij})(u_i^T v_j - \log P_{ij})^2 \quad (154)$$

where P_{ij} is computed simply from the counts of words i and j co-occurring (empirical probability) and f is some squashing function that really isn’t discussed in this lecture (**TODO**).

Evaluating word vectors.

- **Word Vector Analogies:** Basically, determining if we can do standard analogy fill-in-the-blank problems: “*man [a] is to woman [b] as king [c] is to <blank>*” (if you answered “queen”, you’d make a good AI). We can determine this using a standard cosine distance measure:

$$d = \arg \max_i \frac{(x_b - x_a + x_c)^T x_i}{||x_b - x_a + x_c||} \quad (155)$$

Woah that is pretty neat. The solution is $x_i = \text{queen}$. $x_b - x_a$ is the vector pointing from man to woman, which encodes the type of similarity we are looking for with the other pair. Therefore, we take the vector to “king” and *add* the aforementioned difference vector – the resultant vector should point to “queen”. Neat!

SPEECH AND LANGUAGE PROCESSING

CONTENTS

6.1	Naive Bayes and Sentiment (Ch. 6)	95
6.2	Vector Semantics (Ch. 15)	97
6.3	Semantics with Dense Vectors (Ch. 16)	100

Naive Bayes and Sentiment (Ch. 6): June 21, 2017

Table of Contents Local

Written by Brandon McKinzie

[Quick Review]

Overview. This chapter is concerned with **text categorization**, the task of classifying an entire text by assigning it a label drawn from some set of labels. *Generative classifiers* like naive Bayes build a model of each class. Given an observation, they return the class most likely to have generated the observation. *Discriminative classifiers* like logistic regression instead learn what features from the input are most useful to discriminate between the different possible classes. Notation: we will assume we have a training set of N documents, each hand-labeled with some class: $\{(d_1, c_1), \dots, (d_N, c_N)\}$.

Discriminative systems are often more accurate and hence more commonly used.

Naive Bayes. A multinomial⁵⁰ classifier with a naive assumption about how the features interact. We model a text document as a **bag of words**, meaning we store (1) the words that occurred and (2) their frequencies. It's a probabilistic classifier, meaning it estimates the label/class of a document d as

$$\hat{c} = \arg \max_{c \in C} \Pr [c \mid d] \quad (157)$$

$$= \arg \max_{c \in C} \frac{\Pr [d \mid c] \Pr [c]}{\Pr [d]} \quad (158)$$

$$= \arg \max_{c \in C} \Pr [d \mid c] \Pr [c] \quad (159)$$

Computing the class-conditional distribution (the likelihood) $\Pr [d \mid c]$ over all possible $d \in D$ is typically intractable; we must introduce some simplifying assumptions and use an approximation of it. Our assumptions in this section will be:

50

Rapid review: **multinomial distribution.** Let $\mathbf{x} = (x_1, \dots, x_k)$ denote the result of an experiment with n independent trials ($n = \sum_i^k x_i$) and k possible outcomes for any given trial. i.e. x_i is the number of trials that had outcome i ($1 \leq i \leq k$). The pmf of this multinomial distribution, over all possible \mathbf{x} constrained by $n = \sum_i^k x_i$, is:

$$\Pr [\mathbf{x} = (x_1, \dots, x_k); n] = \frac{n!}{x_1! \dots x_k!} p_1^{x_1} \times \dots \times p_k^{x_k} \quad (156)$$

where p_i is the probability of outcome i for any single trial.

- **Bag-of-Words:** Assume that word position within a document is irrelevant. (Counts still matter)
- **Naive Bayes Assumption:** First, recall that d is typically modeled as a (random) vector consisting of features f_1, \dots, f_n , each of which has an associated probability distribution $\Pr[f_i | c]$. The NB assumption is that the features are mutually independent given the class c :

$$\Pr[f_1, \dots, f_n | c] = \Pr[f_1 | c] \cdots \Pr[f_n | c] \quad (160)$$

$$c_{NB} = \arg \max_{c \in C} \Pr[c] \prod_{f \in F} \Pr[f | c] \quad (161)$$

where 161 is the final equation for the class chosen by the naive Bayes classifier.

In text classification we typically use the word at position i in the document as f_i , and move to log space to avoid underflow/increase speed:

$$c_{NB} = \arg \max_{c \in C} \log \Pr[c] + \sum_i^{\text{len}(d)} \log \Pr[w_i | c] \quad (162)$$

Classifiers that use a linear combination of the inputs to make a classification decision (e.g. NB, logistic regression) are called linear classifiers.

Training the NB Classifier. No real "training" in my opinion, just simple counting from the data:

$$\hat{P}[c] = \frac{N_c}{N_{docs}} \quad (163)$$

$$\hat{P}[w_i | c] = \frac{\text{count}(w_i, c) + 1}{(\sum_{w \in V} \text{count}(w, c)) + |V|} \quad (164)$$

The Laplace smoothing is added to avoid the occurrence of zero-probabilities in equation 161.

Optimizing for Sentiment Analysis.

- It often improves performance [for sentiment] to clip word counts in each document to 1 ("binary NB").
- deal with *negations* in some way.
- Use sentiment lexicons, lists of words that are pre-annotated with positive or negative sentiment.

Vector Semantics (Ch. 15): June 21, 2017

Table of Contents Local

Written by Brandon McKinzie

Words and Vectors. Vector models are generally based on a **co-occurrence**, an example of which is a **term-document matrix**: Each row is identified by a word, and each column a document. A given cell value is the number of times the assoc. word occurred in the assoc. document. Can also view each column as a document vector.

Information Retrieval:
task of finding document
 d , from D docs total,
that best matches a
query q .

For individual word vectors, however, it is most common to instead use a **term-term matrix**⁵¹, in which columns are also identified by individual words. Now, cell values are the number of times the row (target) word and the column (context) word co-occur in some context in some training corpus. The context is most commonly a window around the row/target word, meaning a cell gives the number of times the column word occurs in a window of $\pm N$ words from the row word.

- **Q**: What about the co-occurrence of a word with itself (row i , col i)?
 - **A**: It is included, yes. Source: “The size of the window . . . is generally between 1 and 8 words on each side of the target word (*for a total context of 3-17 words*).”
- **Q**: Why is the size of each vector generally $|V|$ (vocab size)? Shouldn’t this vary substantially with window and corpus size?
 - **A**: idk

(**TODO**: revisit end of page 5 in my pdf of this)

Pointwise Mutual Information (PMI). Motivation: raw frequencies in a co-occurrence matrix aren’t that great, and words like “the” (which aren’t useful and occur everywhere) can really skew things. *The best weighting or measure of association between words should tell us how much more often than chance the two words co-occur.* PMI is such a measure.

$$[\text{Mutual Information}] \quad I(X, Y) = \sum_x \sum_y P(X = x, Y = y) \text{PMI}(x, y) \quad (165)$$

$$[\text{PMI}] \quad \text{PMI}(x, y) = \ln \frac{P(x, y)}{P(x)P(y)} \quad (166)$$

⁵¹Also called the word-word or term-context matrix

which can be applied for our specific use case as $\text{PMI}(w, c) = \ln \frac{P(w, c)}{P(w)P(c)}$. The interpretation is simple: the denominator tells the joint probability of the given target word w occurring with context word c if they were independent of each other, while the numerator tells us how often we observed the two words together (assuming we compute probability by using the MLE). Therefore, the ratio gives us how an estimate of how much more the target and feature co-occur than we expect by chance⁵². Most people use **Positive PMI**, which is just $\max(0, \text{PMI})$. We can compute a **PPMI matrix** (to replace our co-occurrence matrix), where PPMI_{ij} gives the PPMI value of word w_i with context c_j . The authors show a few formulas which is really distracting since all we actually need is the counts $f_{ij} = \text{counts}(w_i, c_j)$, and from there we can use basic probability and Bayes rule to get the PPMI formula.

- **Q**: Explain why the following is true: very rare words tend to have very high PMI values.
 - **A**: hi
- **Q**: What is the range of α used in PPMI_α ? What is the intuition behind doing this?
 - **A**: For reference:

$$\text{PPMI}_\alpha(w, c) = \max \left(\ln \frac{P(w, c)}{P(w)P_\alpha(c)}, 0 \right) \quad (167)$$

$$P_\alpha(c) = \frac{\text{count}(c)^\alpha}{\sum_{c'} \text{count}(c')^\alpha} \quad (168)$$

Although there are better methods than PPMI for weighted co-occurrence matrices, most notably **TF-IDF**, things like tf-idf are not generally used for measuring *word similarity*. For that, PPMI and significance-testing metrics like t-test and likelihood-ratio are more common. The **t-test** statistic, like PMI, measures how much more frequent the association is than chance.

$$t = \frac{\bar{x} - \mu}{\sqrt{s^2/N}} \quad (169)$$

$$\text{t-test}(a, b) = \frac{P(a, b) - P(a)P(b)}{\sqrt{P(a)P(b)}} \quad (170)$$

where \bar{x} is the observed mean, while μ is the expected mean [under our null-hypothesis of independence].

⁵²Computing PMI this way can be problematic for word pairs with small probability, especially if we have a small corpus. Recognize that PMI should never really be negative, but in practice this happens for such cases

Measuring Similarity. By far the most common similarity metric is the **cosine** of the angle between the vectors:

$$\text{cosine}(\mathbf{v}, \mathbf{w}) = \frac{\mathbf{v} \cdot \mathbf{w}}{|\mathbf{v}| |\mathbf{w}|} \quad (171)$$

Note that, since we've been defining vector elements as frequencies/PPMI values, they won't have negative elements, and thus our cosine similarities will be between 0 and 1 (not -1 and 1).

Alternatives to cosine:

- **Jaccard measure:** Described as "weighted number of overlapping features, normalized", but looks like a silly hack in my opinion:

$$\text{sim}_{Jac}(\mathbf{v}, \mathbf{w}) = \frac{\sum_{i=1}^N \min(\mathbf{v}_i, \mathbf{w}_i)}{\sum_{i=1}^N \max(\mathbf{v}_i, \mathbf{w}_i)} \quad (172)$$

- **Dice measure:** Another hack. This displeases me.

$$\frac{2 \times \sum_{i=1}^N \min(\mathbf{v}_i, \mathbf{w}_i)}{\sum_{i=1}^N (\mathbf{v}_i + \mathbf{w}_i)} \quad (173)$$

- **Jensen-Shannon Divergence:** An alternative to the KL-divergence⁵³, which represents the divergence of each distribution from the mean of the two:

$$\text{sim}_{JS}(\mathbf{v}||\mathbf{w}) = D\left(\mathbf{v} \left\| \frac{\mathbf{v} + \mathbf{w}}{2}\right.\right) + D\left(\mathbf{w} \left\| \frac{\mathbf{v} + \mathbf{w}}{2}\right.\right) \quad (175)$$

⁵³ Idea: if two vectors, \mathbf{v} and \mathbf{w} , each express a probability distribution (their values sum to one), then they are similar to the extent that these probability distributions are similar. The basis of comparing two probability distributions P and Q is the **Kullback-Leibler** divergence or relative entropy, defined as:

$$D(P||Q) = \sum_x P(x) \log \frac{P(x)}{Q(x)} \quad (174)$$

Semantics with Dense Vectors (Ch. 16): June 21, 2017

Table of Contents Local

Written by Brandon McKinzie

Overview. This chapter introduces three methods for generating short, dense vectors: (1) dimensionality reduction like SVD, (2) neural networks like skip-gram or CBOW, and (3) Brown clustering.

Dense Vectors via SVD. Method for finding more important dimensions of a dataset, “important” defined as dimensions wherein the data most varies. First applied (for language) for generating embeddings from term-document matrices in a model called **Latent Semantic Analysis** (LSA). LSA is just SVD on a $|V| \times c$ term-document matrix \mathbf{X} , factorized into $\mathbf{W}\mathbf{\Sigma}\mathbf{C}^T$. By using only the top $k < m$ dimensions of these three matrices, the product becomes a least-squares approx. to the original \mathbf{X} . It also gives us the reduced $|V| \times k$ matrix \mathbf{W}_k , where each row (word) is a k -dimensional vector (embedding). Voil, we have our dense vectors!

$$\mathbf{W} \in \mathbb{R}^{|V| \times m}$$

$$\mathbf{\Sigma} \in \mathbb{R}^{m \times m}$$

$$\mathbf{C}^T \in \mathbb{R}^{m \times c}$$

Note that LSA implementations typically use a particular weighting of each cell in the term-document matrix called the **local** and **global** weights.

$$[\text{local}] \quad \log f(i, j) + 1 \quad (176)$$

$$[\text{global}] \quad 1 + \frac{\sum_j p(i, j) \log p(i, j)}{\log D} \quad (177)$$

$f(i, j)$ is the raw frequency of word i in context j . D is number of docs.

For the case of a word-word matrix, it is common to use PPMI weighting.

Skip-Gram and CBOW. Neural models learn an embedding by starting with a random vector and then iteratively shifting a words embeddings to be more like the embeddings of neighboring words, and less like the embeddings of words that dont occur nearby⁵⁴. Word2vec, for example, learns embeddings by training to predict neighboring words⁵⁵.

- **Skip-Gram:** Learns two embeddings for each word w : the **word embedding** v (within matrix \mathbf{W}) and **context embedding** c (within matrix \mathbf{C}). Visually:

$$\mathbf{W} = \begin{pmatrix} \mathbf{v}_0^T \\ \mathbf{v}_1^T \\ \vdots \\ \mathbf{v}_{|V|}^T \end{pmatrix} \quad \mathbf{C} = (\mathbf{c}_0 \quad \mathbf{c}_1 \quad \cdots \quad \mathbf{c}_{|V|}) \quad (178)$$

For a context window of $L = 2$, and at a given word $\mathbf{v}^{(t)}$ inside the corpus⁵⁶, our goal is to predict the context [words] denoted as $[\mathbf{c}^{(t-2)}, \mathbf{c}^{(t-1)}, \mathbf{c}^{(t+1)}, \mathbf{c}^{(t+2)}]$.

- Example: Consider one of the context words, say $\mathbf{c}^{(t+1)} \triangleq \mathbf{c}_k$, where we also assume it's the k th word in our vocab. Also assume that our target word $\mathbf{v}^{(t)} \triangleq \mathbf{v}_j$ is the j th word in our vocab.
- Our task is to compute $\Pr[\mathbf{c}_k \mid \mathbf{v}_j]$. We do this with a softmax:

$$\Pr[\mathbf{c}_k \mid \mathbf{v}_j] = \frac{e^{\mathbf{c}_k^T \mathbf{v}_j}}{\sum_{i \in |V|} e^{\mathbf{c}_i^T \mathbf{v}_j}} \quad (179)$$

- **CBOW:** Continuous bag of words. Basically the mirror-image of skip-gram. Goal is to predict current word $\mathbf{v}^{(t)}$ from the context window of $2L$ words $[\mathbf{c}^{(t-2)}, \mathbf{c}^{(t-1)}, \mathbf{c}^{(t+1)}, \mathbf{c}^{(t+2)}]$.

As usual, the denominator of the softmax is computationally expensive, and usually we approximate it with **negative sampling**.

In the training phase, the algorithm walks through the corpus, at each target word choosing the surrounding context words as positive examples, and for each positive example also choosing k noise samples or negative samples: non-neighbor words. The goal will be to move the embeddings toward the neighbor words and away from the noise words.

⁵⁴Why? Why is this a sensible assumption? I see no reason a priori why it ought to be true.

⁵⁵Note that the prediction task is not the goal – it just happens to result in good word embeddings. Hacky.

⁵⁶Note that, technically, the position t of $\mathbf{v}^{(t)}$ is irrelevant for our computation; we are predicting those words based on which word $\mathbf{v}^{(t)}$ is in the vocabulary, not it's position in the corpus.

Example: Suppose we come along the following window (in “[]”) (L=2) in our corpus:

lemon, a [tablespoon of apricot preserves or] jam

Ultimately, we want dot products, $\mathbf{c}_i \cdot \text{vector}(\text{“apricot”})$, to be high for all four of the context words \mathbf{c}_i . We do negative sampling by sampling k random noise words according to their [unigram] frequency. So here, for e.g. $k = 2$, this would amount to 8 noise words, 2 for each context word. We want the dot products between “apricot” and these noise words to be low. For a given single context-word pair (w, c) , our training objective is to maximize:

$$\log \sigma(c \cdot w) + \sum_{i=1}^k \mathbb{E}_{w_i \sim p(w)} [\log \sigma(-w_i \cdot w)] \quad (180)$$

In practice, common to use $p^{3/4}(w)$ instead of $p(w)$

Again, the above is for a single context-target word-pair and, accordingly, the summation is only over $k = 2$ (for our example). Don’t try to split the expectation into a summation or anything – just view it as an expected value. To iteratively shift parameters, we use an optimizer like SGD.

The actual model architecture is a typical neural net, progressing as follows:

$$\text{“apricot”} \rightarrow \mathbf{w}^{\text{one-hot}} = [0 \ 0 \ \dots \ 1 \ \dots \ 0] \quad (181)$$

$$\rightarrow \mathbf{h} = \mathbf{W}^T \mathbf{w}^{\text{one-hot}} \quad (182)$$

$$\rightarrow \mathbf{o} = \mathbf{C}^T \mathbf{h} = [c_0^T \mathbf{h}, c_1^T \mathbf{h}, \dots, c_{|V|}^T \mathbf{h}]^T \quad (183)$$

$$\rightarrow \mathbf{y} = \text{softmax}(\mathbf{o}) = [\text{Pr}[\mathbf{c}_0 | \mathbf{h}], \text{Pr}[\mathbf{c}_2 | \mathbf{h}], \dots, \text{Pr}[\mathbf{c}_{|V|} | \mathbf{h}]]^T \quad (184)$$

$$(185)$$

Brown Clustering. An agglomerative clustering algorithm for deriving vector representations of words by clustering words based on their associations with the preceding or following words. Makes use of the **class-based language model** (CBLM), wherein each word w belongs to some class $c \in C$ via the probability $P(w | c)$. CBLMs define

$$P(w_i | w_{i-1}) = P(c_i | c_{i-1})P(w_i | c_i) \quad (186)$$

$$P(\text{corpus} | C) = \prod_{i=1}^n P(w_i | w_{i-1}) \quad (187)$$

A naive and extremely inefficient version of Brown clustering, a hierarchical clustering algorithm, is as follows:

1. Each word is initially assigned to its own cluster.
2. For each cluster pair $(c_i, c_{j \neq i})$, compute the value of eq 187 that would result from merging c_i and c_j into a single cluster. The pair whose merger results in the smallest decrease in eq 187 is merged.
3. Clustering proceeds until all words are in one big cluster.

This process builds a binary tree from the bottom-up, and the binary string corresp. to a word's traversal from leaf-to-root is its representation.

APPENDIX A - QUESTIONS AND STUFF I ALWAYS FORGET

Questions:

- **Q:** In general, how can one tell if a matrix \mathbf{A} has an eigenvalue decomposition?
[insert more conceptual matrix-related questions here . . .]
- **Q:** Let \mathbf{A} be real-symmetric. What can we say about \mathbf{A} ?
 - Proof that eigendecomposition $\mathbf{A} = \mathbf{Q}\mathbf{\Lambda}\mathbf{Q}^T$ exists: Wow this is apparently quite hard to prove according to many online sources. Guess I don't feel so bad now that it wasn't (and still isn't) obvious.
 - Eigendecomposition not unique. This is apparently because two or more eigenvectors may have same eigenvalue.

This is the principal axis theorem: if \mathbf{A} symmetric, then orthonorm basis of e-vects exists.

Stuff I Forget:

- Existence of eigenvalues/eigenvectors. Let $\mathbf{A} \in \mathbb{R}^{n \times n}$.
 - λ is an eigenvalue of \mathbf{A} iff it satisfies $\det(\lambda\mathbf{I} - \mathbf{A}) = 0$. Why? Because it is an equivalent statement as requiring that $(\lambda\mathbf{I} - \mathbf{A})\mathbf{x} = 0$ has a nonzero solution for \mathbf{x} .
 - The following statements are equivalent:
 - * \mathbf{A} is diagonalizable.
 - * \mathbf{A} has n linearly independent eigenvectors.
 - The **eigenspace** of \mathbf{A} corresponding to λ is the solution space of the homogeneous system $(\lambda\mathbf{I} - \mathbf{A})\mathbf{x} = 0$.
 - \mathbf{A} has at most n distinct eigenvalues.
- Diagonalizability notes from 5.2 of advanced linear alg. book (261). Recall that \mathbf{A} is defined to be diagonalizable if and only if there exists an ordered basis β for the space consisting of eigenvectors of \mathbf{A} .
 - If the standard way of finding eigenvalues leads to k distinct λ_i , then the corresponding set of k eigenvectors v_i are guaranteed to be linearly independent (but might not span the full space).
 - If \mathbf{A} has n linearly independent eigenvectors, then \mathbf{A} is diagonalizable.
 - The characteristic polynomial of any diagonalizable linear operator splits (can be factored into product of linear factors). The **algebraic multiplicity** of an eigenvalue λ is the largest positive integer k for which $(t - \lambda)^k$ is a factor of $f(t)$.
- Logistic regression is called a *linear model* because it classifies based on the linear $\theta^T \mathbf{x}$ before feeding that to the logistic sigmoid function.

Most info here comes from chapter 5 of your "Elementary Linear Algebra" textbook (around pg305)

Recall that a linear operator is a special case of a linear map where the input space is the same as the output space.