

# CONTENTS

<b>1</b>	<b>Math and Machine Learning Basics</b>	<b>2</b>
1.1	Linear Algebra (Quick Review) (Ch. 2)	3
1.1.1	Example: Principal Component Analysis	5
1.2	Probability & Information Theory (Quick Review) (Ch. 3)	7
1.3	Numerical Computation (Ch. 4)	9
1.4	Machine Learning Basics (Ch. 5)	12
1.4.1	Estimators, Bias and Variance (5.4)	12
1.4.2	Maximum Likelihood Estimation (5.5)	14
1.4.3	Bayesian Statistics (5.6)	15
1.4.4	Supervised Learning Algorithms (5.7)	17
<b>2</b>	<b>Deep Networks: Modern Practices</b>	<b>18</b>
2.1	Deep Feedforward Networks (Ch. 6)	19
2.1.1	Back-Propagation (6.5)	20
2.2	Regularization for Deep Learning (Ch. 7)	21
2.3	Optimization for Training Deep Models (Ch. 8)	23
2.4	Convolutional Neural Networks (Ch. 9)	27
2.5	Sequence Modeling (RNNs) (Ch. 10)	30
2.5.1	Review: The Basics of RNNs	30
2.5.2	RNNs as Directed Graphical Models	35
2.5.3	Challenge of Long-Term Deps. (10.7)	37
2.5.4	LSTMs and Other Gated RNNs (10.10)	38
2.6	Applications (Ch. 12)	39
2.6.1	Natural Language Processing (12.4)	39
2.6.2	Neural Language Models (12.4.2)	40
<b>3</b>	<b>Deep Learning Research</b>	<b>41</b>
3.1	Linear Factor Models (Ch. 13)	42
3.2	Autoencoders (Ch. 14)	45
3.3	Representation Learning (Ch. 15)	46
3.4	Structured Probabilistic Models for DL (Ch. 16)	47
3.4.1	Sampling from Graphical Models	49
3.4.2	Inference and Approximate Inference	49
3.5	Monte Carlo Methods (Ch. 17)	51
3.6	Confronting the Partition Function (Ch. 18)	53
3.7	Approximate Inference (Ch. 19)	54
3.8	Deep Generative Models (Ch. 20)	56

# MATH AND MACHINE LEARNING BASICS

## CONTENTS

1.1	Linear Algebra (Quick Review) (Ch. 2)	3
1.1.1	Example: Principal Component Analysis	5
1.2	Probability & Information Theory (Quick Review) (Ch. 3)	7
1.3	Numerical Computation (Ch. 4)	9
1.4	Machine Learning Basics (Ch. 5)	12
1.4.1	Estimators, Bias and Variance (5.4)	12
1.4.2	Maximum Likelihood Estimation (5.5)	14
1.4.3	Bayesian Statistics (5.6)	15
1.4.4	Supervised Learning Algorithms (5.7)	17

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

Table of Contents   Local

Written by Brandon McKinzie

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

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:

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

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

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

- 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 by<sup>1</sup>

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

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

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

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

- 1) **Change of basis:** The vector  $(Q^T \mathbf{x})$  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  $(Q^T \mathbf{x})_i$  by an amount  $\lambda_i$ , yielding the new vector  $(\Lambda(Q^T \mathbf{x}))$ .
- 3) **Change of basis:** Finally, we rotate this new vector back from the eigen-basis into its original basis, yielding the transformed result of  $Q\Lambda Q^T \mathbf{x}$ .

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

- **Positive definite:** all  $\lambda$  are positive; **positive semidefinite:** all  $\lambda$  are positive or zero.  
 → PSD:  $\forall \mathbf{x}, \mathbf{x}^T \mathbf{A} \mathbf{x} \geq 0$   
 → PD:  $\mathbf{x}^T \mathbf{A} \mathbf{x} = 0 \Rightarrow \mathbf{x} = \mathbf{0}$ .<sup>2</sup>
- 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  $D_{ii}$ .
- 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.

---

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

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

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

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

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

**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:

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

→ Large  $|\text{Cov}[f, g]|$  means the function values change a lot and both functions are far from their means at the same time.

→ **Correlation** normalizes the contribution of each variable in order to measure only how much the variables are related.

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

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

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

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

where  $m$  is the number of data points.

## Measure Theory.

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

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

## Functions of RVs.

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

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

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

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

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

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

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

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

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

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

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



## Numerical Computation (Ch. 4)

Table of Contents   Local

Written by Brandon McKinzie

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

The Hessian is the  
 Jacobian of the gradient.

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

3. Solve the re-designed unconstrained optimization problem:

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

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

---

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

## Machine Learning Basics (Ch. 5)

Table of Contents   Local

Written by Brandon McKinzie

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

## 1.4.1 ESTIMATORS, BIAS AND VARIANCE (5.4)

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

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

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

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

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

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

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

### Bias-Variance Tradeoff.

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

→ **Summary of main formulas.**

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

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

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

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

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

and desirable estimators are those with low MSE.

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

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

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

---

<sup>6</sup>Derivation:

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

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

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

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

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

---

### 1.4.2 MAXIMUM LIKELIHOOD ESTIMATION (5.5)

---

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

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

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

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

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

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

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

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

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

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

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

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

---

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

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

Distinction between frequentist and bayesian approach:

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

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

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

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

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

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

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

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

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

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

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

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

---

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

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

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

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

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

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

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

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

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

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

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

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



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

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

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

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

Equation 5.81 is the definition of logistic regression

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

**Support Vector Machines.** Driving by a linear function  $\mathbf{w}^T \mathbf{x} + \mathbf{b}$  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  $k(\mathbf{x}, \mathbf{x}^{(i)}) = \mathbf{x}^T \mathbf{x}^{(i)}$  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  $\boldsymbol{\alpha}$  with mostly zeros. The training examples with *nonzero*  $\alpha_i$  are known as **support vectors**.

# DEEP NETWORKS: MODERN PRACTICES

## CONTENTS

2.1	Deep Feedforward Networks (Ch. 6) . . . . .	19
2.1.1	Back-Propagation (6.5) . . . . .	20
2.2	Regularization for Deep Learning (Ch. 7) . . . . .	21
2.3	Optimization for Training Deep Models (Ch. 8) . . . . .	23
2.4	Convolutional Neural Networks (Ch. 9) . . . . .	27
2.5	Sequence Modeling (RNNs) (Ch. 10) . . . . .	30
2.5.1	Review: The Basics of RNNs . . . . .	30
2.5.2	RNNs as Directed Graphical Models . . . . .	35
2.5.3	Challenge of Long-Term Deps. (10.7) . . . . .	37
2.5.4	LSTMs and Other Gated RNNs (10.10) . . . . .	38
2.6	Applications (Ch. 12) . . . . .	39
2.6.1	Natural Language Processing (12.4) . . . . .	39
2.6.2	Neural Language Models (12.4.2) . . . . .	40

## Deep Feedforward Networks (Ch. 6)

Table of Contents    Local

Written by Brandon McKinzie

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

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

with  $\phi$  defining a hidden layer.

**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  $\alpha_i$  when  $z_i < 0$ :

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

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

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

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

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

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

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

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

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

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

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

---

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

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

## Regularization for Deep Learning (Ch. 7)

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

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

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

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

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

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

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

where  $H_{ij} = \frac{\partial^2 J}{\partial w_i \partial w_j} \big|_{w^*}$ . 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.** 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** 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!

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

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

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

Table of Contents   Local

Written by Brandon McKinzie

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

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

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

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

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

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

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

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

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

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

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

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

The key points to consider when choosing your batch size:

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

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

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

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

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

---

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



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

- **SGD.**

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

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

- **Momentum.**

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

- **Adam.** So-named to mean “adaptive moments.”<sup>13</sup> We now call  $\mathbf{r}$  the 2nd moment (variance) variable, and introduce  $\mathbf{s}$  as the 1st moment (mean) variable, where the moments are for the [true] gradient; the new variables act as estimates of the moments [since we estimate the gradient with a simple average over a minibatch]. Note that these moments are uncentered.

$$\mathbf{g} \leftarrow \frac{1}{m} \nabla_{\boldsymbol{\theta}} \sum_i^m L\left(f(\mathbf{x}^{(i)}; \boldsymbol{\theta}), \mathbf{y}^{(i)}\right) \quad (77)$$

$$\mathbf{s} \leftarrow \frac{1}{1 - \rho_1^{t-1}} [\rho_1 \mathbf{s} + (1 - \rho_1) \mathbf{g}] \quad (78)$$

$$\mathbf{r} \leftarrow \frac{1}{1 - \rho_2^{t-1}} [\rho_2 \mathbf{r} + (1 - \rho_2) \mathbf{g} \odot \mathbf{g}] \quad (79)$$

$$\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} - \frac{\epsilon}{\delta + \sqrt{\mathbf{r}}} \odot \mathbf{s} \quad (80)$$

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

---

<sup>13</sup>Review of moments: The  $n$ th moment of a real-valued continuous function  $f(x)$  of a real variable about a value  $c$  is

$$\mu_n = \int_{-\infty}^{\infty} (x - c)^n f(x) dx \quad (76)$$

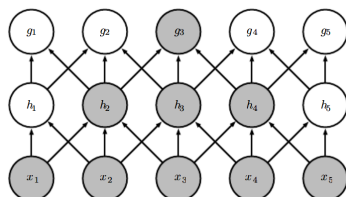
## Convolutional Neural Networks (Ch. 9)

We use a 2-D image  $I$  as our input (and therefore require a 2-D kernel  $K$ ). Note that most neural networks do not technically implement convolution<sup>14</sup>, but instead implement a related function called the *cross-correlation*, defined as

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

Convolution leverages the following three important ideas:

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



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

<sup>14</sup>Technically the convolution output is defined as

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

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

where 9.5 can be asserted due to commutativity of convolution.

**Pooling.** Helps make the representation approximately **invariant** to small translations of the input. The use of pooling can be viewed as adding an infinitely strong prior<sup>15</sup> that the function the layer learns must be invariant to small translations.

**Additional common tricks**<sup>16</sup>.

- **Local Response Normalization (LRN)**<sup>17</sup>. Purpose is to aid generalization ability. Let  $a_{x,y}^i$  denote the activity of a neuron computed by applying kernel  $i$  at position  $(x, y)$  and then applying the ReLU nonlinearity. The response-normalized activity  $b_{x,y}^i$  is given by the expression

$$b_{x,y}^i = \frac{a_{x,y}^i}{\left(k + \alpha \sum_j \left(a_{x,y}^j\right)^2\right)^\beta} \quad (81)$$

where  $j$  runs from  $[i - n/2]_+$  to  $\min(N - 1, i + n/2)$ , and  $N$  is the total number of kernels in the given layer<sup>18</sup>. Authors used  $k = 2$ ,  $n = 5$ ,  $\alpha = 10^{-4}$ ,  $\beta = 0.75$ .

- **Batch Normalization**<sup>19</sup>. BN Allows us to use much higher learning rates and be less careful about initialization. Algorithm defined in image below, where each element of the batch,  $x_i \equiv x_i^{(k)}$  (where we drop the  $k$  for notational simplicity), represents the  $k$ th activation output from the previous layer [for the  $i$ th sample in the batch] and about to be fed as input to the current layer.

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

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

<sup>15</sup>Where the distribution of this prior is over all possible functions learnable by the model.

<sup>16</sup>Collected on my own. In other words, not from the deep learning book, but rather a bunch of disconnected resources over time.

<sup>17</sup>From section 3.3 of Krizhevsky et al. (2012). AlexNet paper.

<sup>18</sup>In other words, the summation is over the adjacent kernel maps, with [total] window size  $n$  (manually chosen). The min/max just says  $n/2$  to the left (right) unless that would be past the leftmost (rightmost) kernel map.

<sup>19</sup>From "Batch Normalization: Accelerating Deep Network Training by Reducing Internal Covariate Shift" by Ioffe et al.

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

### Intuition of some math.

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

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

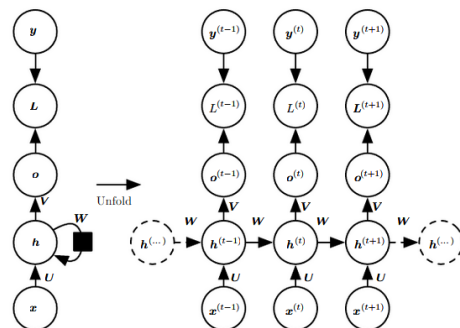
Table of Contents Local

Written by Brandon McKinzie

## 2.5.1 REVIEW: THE BASICS OF RNNs

## Notation/Architecture Used.

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

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

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

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

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

Note that this is an example of an RNN that maps input seqs to output seqs of the same length<sup>20</sup>. We can then compute, e.g., the log-likelihood loss  $L = \sum_t L^{(t)}$  over all time steps as:

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

Convince yourself this is identical to cross-entropy.

<sup>20</sup>Where “same length” is related to the number of timesteps (i.e.  $\tau$  input steps means  $\tau$  output steps), not anything about the actual shapes/sizes of each individual input/output.

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

## 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 **descendents**] to loss node(s).

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

$$\begin{aligned}
\left(\nabla_{\mathbf{o}^{(t)}} L\right)_i &= \frac{\partial L}{\partial \mathbf{o}_i^{(t)}} \\
&= \frac{\partial L}{\partial L^{(t)}} \cdot \frac{\partial L^{(t)}}{\partial \mathbf{o}_i^{(t)}} \\
&= (1) \cdot \frac{\partial L^{(t)}}{\partial \mathbf{o}_i^{(t)}} \\
&= \frac{\partial}{\partial \mathbf{o}_i^{(t)}} \left\{ -\log \left( \hat{\mathbf{y}}_{y^{(t)}}^{(t)} \right) \right\} \\
&= -\frac{\partial}{\partial \mathbf{o}_i^{(t)}} \left\{ \log \left( \frac{e^{\mathbf{o}_{y^{(t)}}^{(t)}}}{\sum_j e^{\mathbf{o}_j^{(t)}}} \right) \right\} \\
&= -\frac{\partial}{\partial \mathbf{o}_i^{(t)}} \left\{ \mathbf{o}_{y^{(t)}}^{(t)} - \log \left( \sum_j e^{\mathbf{o}_j^{(t)}} \right) \right\} \\
&= -\left\{ \mathbf{1}_{i,y^{(t)}} - \frac{\partial}{\partial \mathbf{o}_i^{(t)}} \log \left( \sum_j e^{\mathbf{o}_j^{(t)}} \right) \right\}
\end{aligned} \tag{82}$$

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)).

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

Before proceeding, **notice the following useful pattern:** If two nodes  $a$  and  $b$ , each containing  $n_a$  and  $n_b$  neurons, are fully connected by parameter matrix  $W_{n_b \times n_a}$  and directed like  $a \rightarrow b \rightarrow L$ , then<sup>21</sup>  $\nabla_a L = W^T \nabla_b L$ . Using this result, we can then iterate and take gradients back in time from  $t = \tau - 1$  to  $t = 1$  as follows:

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

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

$$\begin{aligned}
\frac{d}{dx} \tanh(x) &= 1 - \tanh^2(x) \\
(\text{diag}(\mathbf{a}))_{ii} &\triangleq a_i
\end{aligned}$$

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

---

<sup>21</sup>More generally,

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

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

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



previously done with vectors after a quick abstraction: For any tensor  $\mathbf{X}$  of arbitrary rank (e.g. if rank-4 then index like  $\mathbf{X}_{ijkl}$ ), use single variable (e.g.  $i$ ) to represent the complete tuple of indices<sup>22</sup>.

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

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

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

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

$$\nabla_{\mathbf{V}} L = \sum_t \nabla_{\mathbf{V}} L^{(t)} \quad (83a)$$

$$= \sum_t \nabla_{\mathbf{V}} L^{(t)}(\mathbf{o}_1^{(t)}, \dots, \mathbf{o}_{n_{out}}^{(t)}) \quad (83b)$$

$$= \sum_t \sum_i^{n_{out}} (\nabla_{\mathbf{o}^{(t)}} L)_i \nabla_{\mathbf{V}} \mathbf{o}_i^{(t)} \quad (83c)$$

$$= \sum_t \sum_i^{n_{out}} (\nabla_{\mathbf{o}^{(t)}} L)_i \nabla_{\mathbf{V}} \left\{ c_i + \sum_{j=1}^{n_{hid}} V_{ij} \mathbf{h}_j^{(t)} \right\} \quad (83d)$$

$$= \sum_t \sum_i^{n_{out}} (\nabla_{\mathbf{o}^{(t)}} L)_i \begin{bmatrix} 0 & 0 & \dots & 0 \\ \vdots & \vdots & \dots & \vdots \\ \mathbf{h}_1^{(t)} & \mathbf{h}_2^{(t)} & \dots & \mathbf{h}_{n_{hid}}^{(t)} \\ \vdots & \vdots & \dots & \vdots \\ 0 & 0 & \dots & 0 \end{bmatrix} \quad (83e)$$

$$= \sum_t (\nabla_{\mathbf{o}^{(t)}} L) (\mathbf{h}^{(t)})^T \quad (83f)$$

---

<sup>22</sup>More details on tensor derivatives: Consider the chain defined by  $\mathbf{Y} = g(\mathbf{X})$ , and  $z = f(\mathbf{Y})$ , where  $z$  is some vector. Then

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

where if 83e confuses you, see the footnote<sup>23</sup>.

- **W** ( $n_{hid} \times n_{hid}$ ). This one is a bit odd, since **W** is, in a sense, even more “shared” across time steps than **V**<sup>24</sup>. The authors here define/choose, when evaluating  $\nabla_{\mathbf{W}} h_i^{(t)}$  to only concern themselves with  $\mathbf{W} := \mathbf{W}^{(t)}$ , i.e. the direct connections to **h** at time  $t$ .

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

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

$$= \sum_t^\tau \sum_i^{n_{hid}} \left( \nabla_{\mathbf{h}^{(t)}} L \right)_i \left( \text{diag} \left( 1 - (\mathbf{h}^{(t)})^2 \right) \begin{bmatrix} 0 & 0 & \dots & 0 \\ \vdots & \vdots & \dots & \vdots \\ \mathbf{h}_1^{(t-1)} & \mathbf{h}_2^{(t-1)} & \dots & \mathbf{h}_{n_{hid}}^{(t-1)} \\ \vdots & \vdots & \dots & \vdots \\ 0 & 0 & \dots & 0 \end{bmatrix} \right) \quad (85b)$$

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

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

$$\nabla_{\mathbf{U}} L = \sum_t^\tau \nabla_{\mathbf{U}} L^{(t)} \quad (86a)$$

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

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

---

<sup>23</sup>The general lesson learned here is that, for some matrix  $\mathbf{W} \in \mathbb{R}^{a \times b}$  and vector  $\mathbf{x} \in \mathbb{R}^b$ ,

$$\sum_i \nabla_{\mathbf{W}} [(\mathbf{W}\mathbf{x})_i] = \begin{bmatrix} \mathbf{x}^T \\ \mathbf{x}^T \\ \vdots \\ \mathbf{x}^T \end{bmatrix} \quad (84)$$

where, of course, the output has the same dimensions as **W**.

<sup>24</sup>Specifically,  $\mathbf{h}^{(t)}$  is both

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

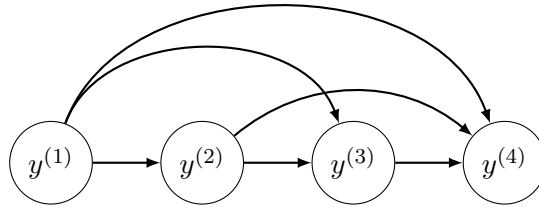
This is different than before, where we had  $\mathbf{o}^{(t)}$  not implicitly depending on earlier  $\mathbf{V}^{(t=i)}$ . In other words,  $\mathbf{h}^{(t)}$  is a descendant of all earlier (and current) **W**.

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.

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

$$P(\mathbb{Y}) \triangleq P(y^{(1)}, \dots, y^{(\tau)}) = \prod_{t=1}^{\tau} P(y^{(t)} \mid y^{(t-1)}, \dots, y^{(1)}) \quad (87)$$

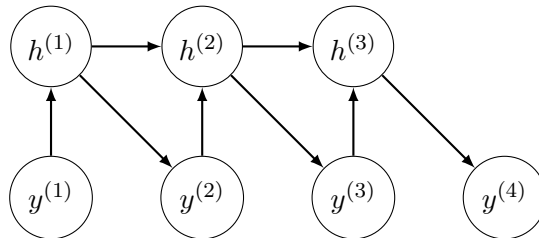
where I've drawn an example of the *complete graph* for  $\tau = 4$  below.



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

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

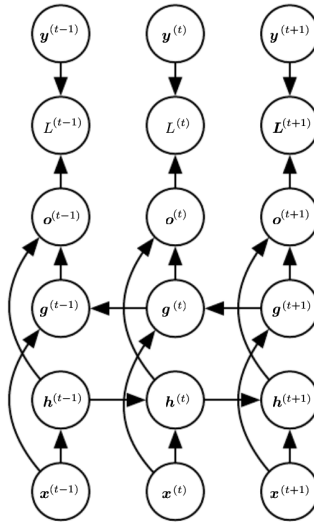
We could include latent variables  $\mathbf{h}$  at each timestep that capture the dependencies, reminiscent of a classic RNN:



Since in the RNN case all factors  $P(h^{(t)} | h^{(t-1)})$  are deterministic, we don't need any additional parameters to compute this probability<sup>25</sup>, other than the single  $m^2$  parameters needed to convert any  $h^{(t)}$  to the next  $h^{(t+1)}$  (which is shared across all transitions). Now, the number of parameters needed as a function of sequence length is constant, and as a function of  $k$  is just  $\mathcal{O}(k)$ .

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

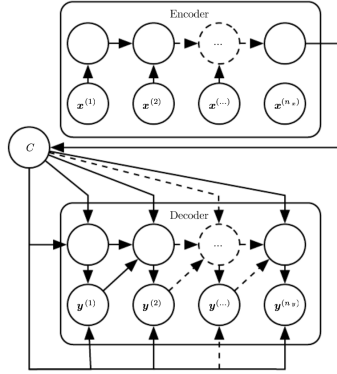
**Bidirectional RNNs.** In many applications, it is desirable to output a prediction of  $\mathbf{y}^{(t)}$  that may depend on *the whole sequence*. For example, in speech recognition, the interpretation of words/sentences can also depend on what is *about* to be said. Below is a typical bidirectional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

<sup>25</sup>Don't forget that, in a neural net, a variable  $y^{(t)}$  is represented by a *layer*, which itself is composed of  $k$  nodes, each associated with one of the  $k$  unique values that  $y^{(t)}$  could be.

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



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

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

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

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

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

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

**Any component of  $\mathbf{h}^{(0)}$  that isn't aligned with the largest eigenvector will eventually be discarded.**<sup>26</sup>

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

<sup>26</sup>Make sure to think about this from the right perspective. The largest value of  $t = \tau$  in the RNNs we've seen would correspond with either (1) the largest output sequence or (2) the largest input sequence (if fixed-vector output). After we extract the output from a given forward pass, we reset the clock and either back-propagate errors (if training) or get ready to feed another sequence.

be  $\mathcal{O}(v^n)^{27}$ . To obtain some desired variance  $v^*$  we may choose the individual weights with variance  $v = \sqrt[n]{v^*}$ .

---

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

---

While leaky units have connection weights that are either manually chosen constants or are trainable parameters, gated RNNs generalize this to connection weights that may change *at each time step*. Furthermore, gated RNNs can learn to both accumulate and *forget*, while leaky units are designed for just accumulation<sup>28</sup>

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

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

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

The final hidden state can then be computed via

$$h_i^{(t)} = \tanh(s_i^{(t)}) \odot q_i^{(t)} \quad (90)$$

---

<sup>27</sup>Quick sketch of (my) proof:

$$\text{Var} [w^{(i)}] = v = \mathbb{E} [(w^{(i)})^2] - \cancel{\mathbb{E} [w^{(i)}]^2} \quad (88)$$

$$\text{Var} \left[ \prod_t^n w^{(t)} \right] = \mathbb{E} \left[ \left( \prod_t^n w^{(t)} \right)^2 \right] = \prod_t^n \mathbb{E} [(w^{(t)})^2] = v^n \quad (89)$$

<sup>28</sup>Q: Isn't choosing to update with higher relative weight on the present the same as forgetting? A: Sort of. It's like "soft forgetting" and will inevitably erase more/less than desired (smeary). In this context, "forget" means to set the weight of a specific past cell to zero.

## Applications (Ch. 12)

Table of Contents   Local

*Written by Brandon McKinzie*

## 2.6.1 NATURAL LANGUAGE PROCESSING (12.4)

Begins on pg. 448

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

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

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

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

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

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

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

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

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

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

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

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 set<sup>29</sup>, even for modest  $n$ .

## 2.6.2 NEURAL LANGUAGE MODELS (12.4.2)

Designed to overcome curse of dimensionality by using a distributed representation of words. Recognize that any model trained on sentences of length  $n$  and then told to generalize to new sentences [also of length  $n$ ] must deal with a space<sup>30</sup> of possible sentences that is exponential in  $n$ . Such word representations (i.e. viewing words as existing in some high-dimensional space) are often called **word embeddings**. The idea is to map the words (or sentences) from the raw high-dimensional [vocab sized] space to a smaller feature space, where similar words are closer to one another. Using distributed representations may also be used with graphical models (think Bayes' nets) in the form multiple *latent variables*.

<sup>29</sup>For a given vocabulary, which usually has much more than  $n$  possible words, consider how many possible sequences of length  $n$ .

<sup>30</sup>Ok I tried re-wording that from the book's confusing wording but that was also a bit confusing. Let me break it down. Say you train on a thousand sentences each of length 5. For a given vocabulary of size VOCAB\_SIZE, the number of possible sequences of length 5 is  $(\text{VOCAB\_SIZE})^5$ , which can be quite a lot more than a thousand (not to mention the possibility of duplicate training examples). To the naive model, all points in this high-dimensional space are basically the same. A neural language model, however, tries to arrange the space of possibilities in a meaningful way, so that an unforeseen sample at test time can be said "similar" as some previously seen training example. It does this by *embedding* words/sentences in a lower-dimensional feature space.



# DEEP LEARNING RESEARCH

## CONTENTS

3.1	Linear Factor Models (Ch. 13)	42
3.2	Autoencoders (Ch. 14)	45
3.3	Representation Learning (Ch. 15)	46
3.4	Structured Probabilistic Models for DL (Ch. 16)	47
3.4.1	Sampling from Graphical Models	49
3.4.2	Inference and Approximate Inference	49
3.5	Monte Carlo Methods (Ch. 17)	51
3.6	Confronting the Partition Function (Ch. 18)	53
3.7	Approximate Inference (Ch. 19)	54
3.8	Deep Generative Models (Ch. 20)	56

## Linear Factor Models (Ch. 13)

Table of Contents   Local

Written by Brandon McKinzie

**Overview.** Much research is in building a *probabilistic model*<sup>31</sup> 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

$$p_{\text{model}}(x) = \sum_h \Pr(h) \Pr(x | h) = \mathbb{E}_h [p_{\text{model}}(x | h)] \quad (91)$$

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

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

Note that  $\mathbf{h}$  is a *vector* of arbitrary size, where we assume  $p(\mathbf{h})$  is a **factorial distribution**:  $p(\mathbf{h}) = \prod_i p(h_i)$ . This roughly means we assume the elements of  $\mathbf{h}$  are mutually independent<sup>32</sup>. The LFM describes the data-generation process as follows:

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

$$\mathbf{x} = \mathbf{W}\mathbf{h} + \mathbf{b} + \text{noise} \quad (92)$$

**Probabilistic PCA and Factor Analysis.**

- **Factor analysis:**

$$\mathbf{h} \sim \mathcal{N}(\mathbf{h}; \mathbf{0}, \mathbf{I}) \quad (93)$$

$$\text{noise} \sim \mathcal{N}(\mathbf{0}, \boldsymbol{\psi} \equiv \text{diag}(\boldsymbol{\sigma}^2)) \quad (94)$$

$$\mathbf{x} \sim \mathcal{N}(\mathbf{x}; \mathbf{b}, \mathbf{W}\mathbf{W}^T + \boldsymbol{\psi}) \quad (95)$$

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

<sup>31</sup>Whereas, before, we've been building *functions* of the input (deterministic).

<sup>32</sup>Note that, technically, this assumption isn't strictly the definition of mutual independence, which requires that every *subset* (i.e. not just the full set) of  $\{h_i \in \mathbf{h}\}$  follow this factorial property.

It is worth emphasizing the interpretation of  $\boldsymbol{\psi}$  as the matrix of **conditional variances**  $\sigma_i^2$ . *Huh?* Let's take a step back. The fact that we were able to separate the distributions in the above relations for  $\mathbf{h}$  and noise is from a built-in assumption that  $\Pr(x_i|\mathbf{h}, x_{j \neq i}) = \Pr(x_i|\mathbf{h})$ <sup>33</sup>.

### The Big Idea

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

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

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

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

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

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

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

$$\boldsymbol{\mu}_x = \mathbb{E}_{\mathbf{h}} [\mathbf{W}\mathbf{h} + \mathbf{b}] \quad (96)$$

$$= \int p(\mathbf{h})(\mathbf{W}\mathbf{h} + \mathbf{b}) d\mathbf{h} \quad (97)$$

$$= \mathbf{b} + \int \frac{1}{(2\pi)^{n/2}} e^{-\frac{1}{2} \sum_i h_i^2} \mathbf{W}\mathbf{h} d\mathbf{h} \quad (98)$$

$$= \mathbf{b} \quad (99)$$

$$\text{Cov}(\mathbf{x}) = \mathbb{E} [(\mathbf{x} - \mathbb{E}[\mathbf{x}])(\mathbf{x} - \mathbb{E}[\mathbf{x}])^T] \quad (100)$$

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

$$= \mathbb{E} [(\mathbf{W}\mathbf{h}\mathbf{h}^T \mathbf{W}^T)] + \boldsymbol{\psi} \quad (102)$$

$$= \mathbf{W}\mathbf{W}^T + \boldsymbol{\psi} \quad (103)$$

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

<sup>33</sup>Due to  $\langle \text{MATH} \rangle$ , this introduces a constraint that knowing the value of some element  $x_j$  doesn't alter the probability  $\Pr(x_i = W_i \cdot \mathbf{h} + b_i + \text{noise})$ . Given how we've defined the variable  $\mathbf{h}$ , this means that knowing noise<sub>*j*</sub> provides no clues about noise<sub>*i*</sub>. Mathematically, the noise must have a diagonal covariance matrix.

- **Thoughts.** Not really seeing why this is useful/noteworthy. Feels very contrived (many assumptions) and restrictive – it only applies if the dependencies between each  $x_i$  can be modeled with a random variable  $\mathbf{h}$  sampled from a unit variance Gaussian.
- **Probabilistic PCA:** Just factor analysis with  $\boldsymbol{\psi} = \sigma^2 \mathbf{I}$ . 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  $\mathbf{W}$ .

## Autoencoders (Ch. 14)

Table of Contents   Local

*Written by Brandon McKinzie*

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

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.

$$f(\cdot) \leftarrow I(\cdot) \quad (104)$$

$$\tilde{\mathbf{X}} = \mathbf{X} \quad (105)$$

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

$$f^{(k)} = \mathcal{L}(\tilde{\mathbf{X}}) \quad (106)$$

$$f(\cdot) \leftarrow f^{(k)}(f(\cdot)) \quad (107)$$

$$\tilde{\mathbf{X}} \leftarrow f^{(k)}(\tilde{\mathbf{X}}) \quad (108)$$

In English: just apply the regular learning/training process for each layer/stage **sequentially and individually**<sup>34</sup>.

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

---

<sup>34</sup>In other words, you proceed one layer at a time *in order*. You don’t touch layer  $i$  until the weights in layer  $i - 1$  have been learned.

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

Table of Contents   Local

Written by Brandon McKinzie

**Motivation.** In addition to classification, we can ask probabilistic models to perform other tasks such as density estimation ( $\mathbf{x} \rightarrow p(\mathbf{x})$ ), denoising, missing value imputation, or sampling. What these [other] tasks have in common is they require a *complete understanding of the input*. Let's start with the most naive approach of modeling  $p(\mathbf{x})$ , where  $\mathbf{x}$  contains  $n$  elements, each of which can take on  $k$  distinct values: we store a lookup table of all possible  $\mathbf{x}$  and the corresponding probability value  $p(\mathbf{x})$ . This requires  $k^n$  parameters<sup>35</sup>. Instead, we use graphs to describe model structure (direct/indirect interactions) to drastically reduce the number of parameters.

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

$$p(\mathbf{x}) = \prod_i p(x_i | \text{Pa}_{\mathcal{G}}(x_i)) \quad (109)$$

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

Clique potentials are constrained to be nonnegative.

$$\tilde{p}(\mathbf{x}) = \prod_{\mathcal{C} \in \mathcal{G}} \phi(\mathcal{C}) \quad (110)$$

<sup>35</sup>Consider the common NLP case where our vector  $\mathbf{x}$  contains  $n$  word tokens, each of which can take on any symbol in our vocabulary of size  $v$ . If we assign  $n = 100$  and  $v = 100,000$ , which are relatively common values for this case, this amounts to  $(1e5)^{1e2} = 10^{500}$  parameters.

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

$$p(\mathbf{x}) = \frac{1}{Z} \tilde{p}(\mathbf{x}) \quad (111)$$

$$Z = \int \tilde{p}(\mathbf{x}) d\mathbf{x} \quad (112)$$

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

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

$$\tilde{p}(\mathbf{x}) \triangleq \exp(-E(\mathbf{x})) \quad (113)$$

and  $E(\mathbf{x})$  is known as the **energy function**<sup>36</sup>. Many algorithms need to compute not  $p_{\text{model}}(\mathbf{x})$  but only  $\log \tilde{p}_{\text{model}}(\mathbf{x})$  (unnormalized log probabilities - logits!). For EBMs with latent variables  $\mathbf{h}$ , such algorithms are phrased in terms of the **free energy**:

$$\mathcal{F}(\mathbf{x} = x) = -\log \sum_{\mathbf{h}} \exp(-E(\mathbf{x} = x, \mathbf{h} = h)) \quad (114)$$

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

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

- **Separation.** For *undirected* models. If variables  $a$  and  $b$  are connected by a path involving only unobserved variables (an **active** path), then  $a$  and  $b$  are *not* separated. Otherwise, they are separated. Any paths containing at least one observed variable are called **inactive**.
- **D-Separation**<sup>37</sup>. For *directed* models. Although there are rules that help determine whether a path between  $a$  and  $b$  is d-separated, it is simplest to just determine whether  $a$  is independent from  $b$  given any observed variables along the path.

---

<sup>36</sup>Physics throwback: this mirrors the Boltzmann factor,  $\exp(-\varepsilon/\tau)$ , which is proportional to the probability of the system being in quantum energy state  $\varepsilon$ .

<sup>37</sup>The D stands for dependence.



---

### 3.4.1 SAMPLING FROM GRAPHICAL MODELS

---

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

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

---

### 3.4.2 INFERENCE AND APPROXIMATE INFERENCE

---

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

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

$$E(\mathbf{v}, \mathbf{h}) = -\mathbf{b}^T \mathbf{v} - \mathbf{c}^T \mathbf{h} - \mathbf{v}^T \mathbf{W} \mathbf{h} \quad (115)$$

where  $\mathbf{b}$ ,  $\mathbf{c}$ , and  $\mathbf{W}$  are unconstrained, real-valued, learnable parameters. One could interpret the values of the bias parameters as the affinities for the associated variable being its given value, and the value  $\mathbf{W}_{i,j}$  as the affinity of  $v_i$  being its value and  $h_j$  being its value at the same time<sup>38</sup>.

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

---

<sup>38</sup>More concretely, remember that  $\mathbf{v}$  is a one-hot vector representing some state that can assume  $\text{len}(\mathbf{v})$  unique values, and similarly for  $\mathbf{h}$ . Then  $\mathbf{W}_{i,j}$  gives the affinity for the state associated with  $v$  being its  $i$ th value and the state associated with  $h$  being its  $j$ th value.

that:

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

$$p(\mathbf{v} \mid \mathbf{h}) = \prod_i p(v_i \mid \mathbf{h}) \quad (117)$$

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

$$p(h_i = 1 \mid \mathbf{v}) = \sigma \left( c_i + \mathbf{v}^T \mathbf{W}_{:,i} \right) \quad (118)$$

allowing for efficient block Gibbs sampling.

## Monte Carlo Methods (Ch. 17)

Table of Contents   Local

Written by Brandon McKinzie

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

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

As we've seen before, the empirical average is an unbiased<sup>39</sup> estimator. Furthermore, the central limit theorem tells us that the distribution of  $\hat{s}_n$  converges to a normal distribution with mean  $s$  and variance  $\text{Var}[f(\mathbf{x})]/n$ .

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

$$p(\mathbf{x})f(\mathbf{x}) = q(\mathbf{x})\frac{p(\mathbf{x})f(\mathbf{x})}{q(\mathbf{x})} \quad (123)$$

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

$$\hat{s}_q = \frac{1}{n} \sum_{i=1, \mathbf{x}^{(i)} \sim q}^n \frac{p(\mathbf{x}^{(i)})f(\mathbf{x}^{(i)})}{q(\mathbf{x}^{(i)})} \quad (124)$$

---

<sup>39</sup>Recall that expectations on such an average are still taken over the underlying (assumed) probability distribution:

$$\mathbb{E}_p[\hat{s}_n] = \frac{1}{n} \sum_{i=1}^n \mathbb{E}_p[f(\mathbf{x}^{(i)})] \quad (120)$$

$$= \frac{1}{n} \sum_{i=1}^n s \quad (121)$$

$$= s \quad (122)$$

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

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

$$\text{Var} [\hat{s}_q] = \text{Var} \left[ \frac{p(\mathbf{x})f(\mathbf{x})}{q(\mathbf{x})} \right] / n \quad (125)$$

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

$$q^* = \frac{p(\mathbf{x}) | f(\mathbf{x}) |}{Z} \quad (126)$$

- **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:

$$\hat{s}_{BIS} = \frac{\sum_{i=1}^n \frac{\tilde{p}(\mathbf{x}^{(i)})}{\tilde{q}(\mathbf{x}^{(i)})} f(\mathbf{x}^{(i)})}{\sum_{i=1}^n \frac{\tilde{p}(\mathbf{x}^{(i)})}{\tilde{q}(\mathbf{x}^{(i)})}} \quad (127)$$

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

## Confronting the Partition Function (Ch. 18)

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

$$\log p_{\text{model}}(\mathbf{x}) = \log \tilde{p}_{\text{model}}(\mathbf{x}; \boldsymbol{\theta}) + c \quad (128)$$

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

$$p_{\text{joint}}(y=1) = \frac{1}{2} \quad (129)$$

$$p_{\text{joint}}(\mathbf{x} \mid y=1) = p_{\text{model}}(\mathbf{x}) \quad (130)$$

$$p_{\text{joint}}(\mathbf{x} \mid y=0) = p_{\text{noise}}(\mathbf{x}) \quad (131)$$

We can now use MLE on the optimization problem,

$$\boldsymbol{\theta}, c = \arg \max_{\boldsymbol{\theta}, c} \mathbb{E}_{\mathbf{x}, y \sim p_{\text{train}}} [\log p_{\text{joint}}(y \mid \mathbf{x})] \quad (132)$$

$$p_{\text{joint}}(y=1 \mid \mathbf{x}) = \frac{p_{\text{model}}(\mathbf{x})}{p_{\text{model}}(\mathbf{x}) + p_{\text{noise}}(\mathbf{x})} \quad (133)$$

$$= \frac{1}{1 + p_{\text{noise}}(\mathbf{x})/p_{\text{model}}(\mathbf{x})} \quad (134)$$

$$= \sigma(\log p_{\text{model}}(\mathbf{x}) - \log p_{\text{noise}}(\mathbf{x})) \quad (135)$$

## Approximate Inference (Ch. 19)

Table of Contents   Local

Written by Brandon McKinzie

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

Assume we have a probabilistic model consisting of observed variables  $\mathbf{v}$  and latent variables  $\mathbf{h}$ . We want to compute  $\log p(\mathbf{v}; \boldsymbol{\theta})$ , but it's too costly to marginalize out  $\mathbf{h}$ . Instead, we compute a lower bound  $\mathcal{L}(\mathbf{v}, \boldsymbol{\theta}, q)$  – often called the **evidence lower bound** (ELBO) or negative **variational free energy** – on  $\log p(\mathbf{v}; \boldsymbol{\theta})$ <sup>40</sup>:

$$\mathcal{L}(\mathbf{v}, \boldsymbol{\theta}, q) = \log p(\mathbf{v}; \boldsymbol{\theta}) - D_{KL}(q(\mathbf{h} | \mathbf{v}) || p(\mathbf{h} | \mathbf{v}; \boldsymbol{\theta})) \quad (136)$$

$$= -\mathbb{E}_{\mathbf{h} \sim q(\mathbf{h} | \mathbf{v})} [\log p(\mathbf{h}, \mathbf{v})] + H(q(\mathbf{h} | \mathbf{v})) \quad (137)$$

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

where the second form is the more canonical definition<sup>41</sup>. Note that  $\mathcal{L}(\mathbf{v}, \boldsymbol{\theta}, q)$  is a lower-bound on  $\log p(\mathbf{v}; \boldsymbol{\theta})$  by definition, since

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

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

<sup>40</sup>Recall that  $D_{KL}(P||Q) = \mathbb{E}_{x \sim P(x)} \left[ \log \frac{P(x)}{Q(x)} \right]$

<sup>41</sup>This can be derived easily from the first form. Hint:

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

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

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

$$q(\mathbf{h} \mid \mathbf{v}^{(i)}) = p(\mathbf{h} \mid \mathbf{v}^{(i)}; \boldsymbol{\theta}^{(0)}) \quad (138)$$

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

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

$$\arg \max_{\boldsymbol{\theta}} \sum_i \mathcal{L}(\mathbf{v}^{(i)}, \boldsymbol{\theta}, q(\mathbf{h} \mid \mathbf{v}^{(i)}; \boldsymbol{\theta}^{(0)})) \quad (139)$$

## Deep Generative Models (Ch. 20)

Table of Contents   Local

Written by Brandon McKinzie

**Boltzmann Machines** (20.1). An energy-based model over a  $d$ -dimensional binary random vector  $\mathbf{x} \in \{0, 1\}^d$ . The energy function is simply  $E(\mathbf{x}) = -\mathbf{x}^T \mathbf{U} \mathbf{x} - \mathbf{b}^T \mathbf{x}$ , i.e. parameters between all pairs of  $x_i, x_j$ , and bias parameters for each  $x_i$ <sup>42</sup>. In settings where our data consists of samples of fully observed  $\mathbf{x}$ , this is clearly limited to very simple cases, since e.g. the probability of some  $x_i$  being on is given by logistic regression from the values of the other units.

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

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

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

$$E(\mathbf{x}) = - \sum_{i=1}^d \sum_{j=i+1}^d x_i U_{i,j} x_j - \sum_{i=1}^d b_i x_i \quad (140)$$

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

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

$$\Delta E_i = \sum_{k=i+1}^d U_{i,k} x_k + \sum_{k=1}^{i-1} x_k U_{k,i} + b_i \quad (141)$$

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

$$\Delta E_i = \ln p_{i=on}(\mathbf{x}) - \ln p_{i=off}(\mathbf{x}) = -\ln \left( \frac{1 - p_{i=on}(\mathbf{x})}{p_{i=on}(\mathbf{x})} \right) \quad (142)$$

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

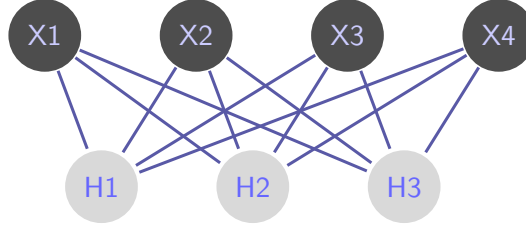
$$\therefore p_{i=on}(\mathbf{x}) = \frac{1}{1 + \exp(-\Delta E_i)} \quad (144)$$

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

<sup>42</sup> Authors are being lazy because it's assumed the reader is familiar (which is fair, I guess). i.e. they aren't mentioning that this formula implies that  $\mathbf{U}$  is either lower or upper triangular, and the diagonal is zero.



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



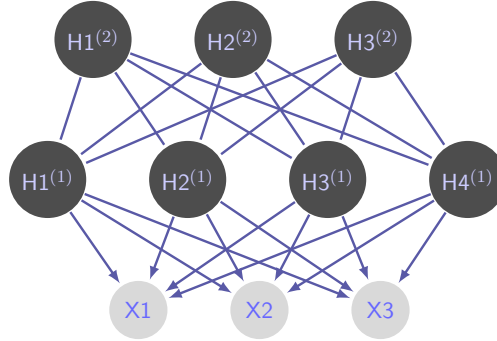
Although the joint distribution  $p(\mathbf{x}, \mathbf{h})$  has a potentially intractable partition function, the conditional distributions can be computed efficiently by exploiting independencies:

$$p(\mathbf{h} \mid \mathbf{x}) = \prod_{j=1}^{n_h} \sigma \left( [2\mathbf{h} - 1] \odot [\mathbf{c} + \mathbf{W}^T \mathbf{x}] \right)_j \quad (145)$$

$$p(\mathbf{x} \mid \mathbf{h}) = \prod_{i=1}^{n_x} \sigma \left( [2\mathbf{x} - 1] \odot [\mathbf{b} + \mathbf{W}\mathbf{h}] \right)_i \quad (146)$$

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

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



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

**Deep Boltzmann Machines** (20.4). Same as DBNs, but now all layers are undirected. Note that this is very close to the standard RBM, since we have a set of hidden and observed variables, except now we interpret certain subgroups of hidden units as being in a “layer”, thus allowing for connections between hidden units in adjacent layers. What’s interesting is that this still defines a bipartite graph, with odd-numbered layers on one side and even on the other<sup>43</sup>.

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

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

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

where  $\mathbf{L}$  is the *Cholesky decomposition*<sup>44</sup> of  $\boldsymbol{\Sigma}$ . In general, we think of the generator function  $g$  as providing a change of variables that transforms the distribution over  $\mathbf{z}$  into the desired distribution  $\mathbf{x}$ . Of course, there *is* an exact formula for doing this,

$$p_{\mathbf{x}}(\mathbf{x}) = \frac{p_{\mathbf{z}}(g^{-1}(\mathbf{x}))}{\left| \det \frac{\partial g}{\partial \mathbf{z}} \right|} \quad (147)$$

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

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

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

$$p(x_i | \mathbf{z}) = g(\mathbf{z})_i \quad p(\mathbf{x}) = \mathbb{E}_{\mathbf{z}} [p(\mathbf{x} | \mathbf{z})] \quad (148)$$

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

<sup>43</sup>Recall that this immediately implies that units in all odd layers are conditionally independent given the even layers (and vice-versa for even to odd).

<sup>44</sup>The [unique] Cholesky decomposition of a (real-symmetric) p.d. matrix  $\mathbf{A}$  is a decomposition of the form  $\mathbf{A} = \mathbf{L}\mathbf{L}^T$ , where  $\mathbf{L}$  is lower triangular.

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