

Machine Learning Study Guide

WILLIAM WATSON

Johns Hopkins University

billwatson@jhu.edu

Contents

1	Linear Algebra and Calculus	4	2.3	Lagrange Duality and KKT Conditions	13
1.1	General Notation	4	2.3.1	The Lagrangian	13
1.1.1	Identity Matrix	4	2.3.2	Primal and Dual Problems	13
1.1.2	Diagonal Matrix	4	2.3.3	Strong and Weak Duality	14
1.1.3	Orthogonal Matrix	4	2.3.4	Complementary Slackness	14
1.2	Matrix Operations	5	2.3.5	The KKT Conditions	14
1.2.1	Vector-Vector Products	5	3	Probability and Statistics	14
1.2.2	Vector-Matrix Products	5	3.1	Basics	14
1.2.3	Matrix-Matrix Products	5	3.1.1	Axioms of Probability	15
1.2.4	The Transpose	6	3.1.2	Permutation	15
1.2.5	The Trace	6	3.1.3	Combination	15
1.2.6	The Inverse	6	3.2	Conditional Probability	15
1.2.7	The Determinant	6	3.2.1	Bayes Rule	15
1.3	Matrix Properties	7	3.2.2	Independence	16
1.3.1	Norms	7	3.3	Random Variables	16
1.3.2	Linear Dependence and Rank . .	7	3.3.1	Cumulative Distribution Function (CDF)	16
1.3.3	Span, Range, and Nullspace . .	8	3.3.2	Probability Density Function (PDF)	16
1.3.4	Symmetric Matrices	8	3.3.3	Discrete PDF/CDF	16
1.3.5	Positive Semidefinite Matrices . .	8	3.3.4	Continuous PDF/CDF	17
1.3.6	Eigendecomposition	9	3.3.5	Expectation	17
1.3.7	Singular Value Decomposition . .	9	3.3.6	Variance and Standard Deviation . .	17
1.3.8	The Moore-Penrose Pseudoinverse .	9	3.4	Discrete Random Variables	18
1.4	Matrix Calculus	10	3.4.1	Bernoulli	18
1.4.1	The Gradient	10	3.4.2	Binomial	18
1.4.2	The Hessian	10	3.4.3	Geometric	18
1.4.3	Gradient Properties	11	3.4.4	Poisson	19
2	Convex Optimization	11	3.5	Continuous Random Variables	19
2.1	Convexity	11	3.5.1	Uniform	19
2.1.1	Convex Sets	11	3.5.2	Exponential	19
2.1.2	Convex Functions	11	3.5.3	Gaussian (Normal)	20
2.1.3	First-Order Conditions	11	3.6	Jointly Distributed Random Variables .	20
2.1.4	Second-Order Conditions	12	3.6.1	Marginal Density	20
2.1.5	Jensen's Inequality	12	3.6.2	Cumulative Distribution	20
2.1.6	Sublevel Sets	12	3.6.3	Conditional Density	20
2.2	Convex Optimization	12	3.6.4	Independence	20
2.2.1	Global Optimality	12	3.6.5	Expectation and Covariance . . .	20
2.2.2	Gradient Descent	13	3.6.6	Correlation	20
2.2.3	Newton's Algorithm	13	3.7	Parameter Estimation	20

3.7.1	Definitions	20	12 Margin Classification	28
3.7.2	Bias	20	13 Generative Learning: Gaussian Discriminant Analysis	28
3.7.3	Mean and Central Limit Theorem	20	13.1 Assumptions	28
3.7.4	Variance	20	13.2 Estimation	28
3.8	Probability Bounds and Inequalities . .	20	14 Generative Learning: Naive Bayes	28
3.8.1	Markov	20	14.1 Assumptions	28
3.8.2	Chebyshev	20	14.2 Estimation	28
3.8.3	Chernoff	20	15 Tree-based Methods	28
3.8.4	Hoeffding	20	16 K-Nearest Neighbors	28
4	Information Theory	20	16.1 Classification	28
5	Learning Theory	22	16.2 Regression	29
5.1	Bias and Variance	22	17 K-Means Clustering	29
5.2	Notation	22	17.1 Algorithm	29
5.2.1	Union Bound	22	17.2 Hierarchical Clustering	29
5.2.2	Hoeffding Inequality	22	17.3 Clustering Metrics	29
5.3	Training Error	22	18 Expectation-Maximization	30
5.4	Probably Approximately Correct (PAC)	22	18.1 Mixture of Gaussians	30
5.5	Hypothesis Classes	22	18.2 Factor Analysis	30
5.5.1	Shattering	22	19 Principal Component Analysis	30
5.5.2	Upper Bound Theorem	22	19.1 Eigenvalues, Eigenvectors, and the Spectral Theorem	30
5.5.3	VC Dimension	22	19.2 Algorithm	30
5.5.4	Vapnik Theorem	22	20 Independent Component Analysis	32
6	Linear Regression	22	21 Reinforcement Learning	32
6.1	LMS Algorithm	22	21.1 Markov Decision Processes	32
6.2	The Normal Equations	23	21.2 Policy and Value Functions	32
6.3	Probabilistic Interpretation	24	21.3 Value Iteration Algorithm	32
6.4	Locally Weighted Linear Regression . .	24	21.4 Q-Learning	32
7	Logistic Regression	25	22 Hidden Markov Models	32
7.1	The Logistic Function	25	23 Deep Learning: Basics	32
7.2	Cost Function	25	23.1 Basics	32
7.3	Gradient Descent	26	23.2 Activation Functions	32
8	Softmax Regression	26	23.3 Loss Functions	32
8.1	Softmax Function	26	23.4 Backpropagation	32
8.2	MLE and Cost Function	26	23.5 Regularization Methods	32
8.3	Gradient Descent	27	23.6 Optimization Algorithms	32
9	Generalized Linear Models	28	23.7 Convolutional Networks	32
9.1	Exponential Family	28	23.8 Recurrent Networks	32
9.2	Assumptions of GLMs	28	23.8.1 Elman RNN	32
9.3	Examples	28	23.8.2 Long Short-Term Memory	32
9.3.1	Ordinary Least Squares	28		
9.3.2	Logistic Regression	28		
9.3.3	Softmax Regression	28		
10	Perceptron	28		
11	Support Vector Machines	28		

23.8.3 Gated Recurrent Unit	32
23.8.4 Bidirectional RNNs	32
24 Deep Learning: Advanced	32
24.1 Autoencoders	32
24.1.1 Variational Autoencoders	32
24.2 General Adversarial Networks	32
24.3 Encoder-Decoder Models	32
24.3.1 Attention Models	32

1 Linear Algebra and Calculus

1.1 General Notation

A vector $x \in \mathbb{R}^n$ has n entries, and $x_i \in \mathbb{R}$ is the i -th entry:

$$x = \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{pmatrix} \in \mathbb{R}^n \quad (1)$$

We denote a matrix $A \in \mathbb{R}^{m \times n}$ with m rows and n columns, and $A_{ij} \in \mathbb{R}$ is the entry in the i -th row and j -th column:

$$A = \begin{pmatrix} A_{11} & \cdots & A_{1n} \\ \vdots & & \vdots \\ A_{m1} & \cdots & A_{mn} \end{pmatrix} \in \mathbb{R}^{m \times n} \quad (2)$$

Vectors can be viewed as a $n \times 1$ matrix.

1.1.1 Identity Matrix

The identity matrix $I \in \mathbb{R}^{n \times n}$ is a square matrix with ones along the diagonal and zero everywhere else:

$$I = \begin{pmatrix} 1 & 0 & \cdots & 0 \\ 0 & \ddots & \ddots & \vdots \\ \vdots & \ddots & \ddots & 0 \\ 0 & \cdots & 0 & 1 \end{pmatrix} \quad (3)$$

For all matrices $A \in \mathbb{R}^{n \times n}$ we have $A \times I = I \times A = A$.

1.1.2 Diagonal Matrix

A diagonal matrix $D \in \mathbb{R}^{n \times n}$ is a square matrix with nonzero values along the diagonal and zero everywhere else:

$$D = \begin{pmatrix} d_1 & 0 & \cdots & 0 \\ 0 & \ddots & \ddots & \vdots \\ \vdots & \ddots & \ddots & 0 \\ 0 & \cdots & 0 & d_n \end{pmatrix} \quad (4)$$

The diagonal matrix D is also written as $\text{diag}(d_1, \dots, d_n)$.

1.1.3 Orthogonal Matrix

Two vectors $x, y \in \mathbb{R}^n$ are orthogonal if $x^T y = 0$. A vector $x \in \mathbb{R}^n$ is normalized if $\|x\|_2 = 1$. A square matrix $U \in \mathbb{R}^{n \times n}$ is orthogonal if all its columns are orthogonal to each other and are normalized. Hence:

$$U^T U = I = U U^T \quad (5)$$

Hence the inverse of an orthogonal matrix is its transpose.

1.2 Matrix Operations

1.2.1 Vector-Vector Products

Given two vectors $x, y \in \mathbb{R}^n$, the inner product is:

$$x^T y = \sum_{i=1}^n x_i y_i \in \mathbb{R} \quad (6)$$

The outer product for a vector $x \in \mathbb{R}^m, y \in \mathbb{R}^n$ is:

$$xy^T = \begin{pmatrix} x_1 y_1 & \cdots & x_1 y_n \\ \vdots & \ddots & \vdots \\ x_m y_1 & \cdots & x_m y_n \end{pmatrix} \in \mathbb{R}^{m \times n} \quad (7)$$

1.2.2 Vector-Matrix Products

The product of a matrix $A \in \mathbb{R}^{m \times n}$ and vector $x \in \mathbb{R}^n$ is a vector $y = Ax \in \mathbb{R}^m$. If we write A by the rows, Ax is expressed as:

$$y = Ax = \begin{pmatrix} -a_1^T - \\ -a_2^T - \\ \vdots \\ -a_m^T - \end{pmatrix} x = \begin{pmatrix} a_1^T x \\ a_2^T x \\ \vdots \\ a_m^T x \end{pmatrix} \quad (8)$$

Here, the i -th entry of y is the inner product of the i -th row of A and x , $y_i = a_i^T x$. If we write A in column form:

$$y = Ax = \begin{pmatrix} | & | & \cdots & | \\ a_1 & a_2 & \cdots & a_n \\ | & | & \cdots & | \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{pmatrix} = a_1 x_1 + a_2 x_2 + \cdots + a_n x_n \quad (9)$$

Here, y is a linear combination of the columns of A , where the coefficients of the linear combination are given by the entries of x .

1.2.3 Matrix-Matrix Products

Given a matrix $A \in \mathbb{R}^{m \times n}$ and matrix $B \in \mathbb{R}^{n \times p}$, we can define $C = AB$ as follows:

$$C = AB = \begin{pmatrix} a_1^T b_1 & a_1^T b_2 & \cdots & a_1^T b_p \\ a_2^T b_1 & a_2^T b_2 & \cdots & a_2^T b_p \\ \vdots & \vdots & \ddots & \vdots \\ a_m^T b_1 & a_m^T b_2 & \cdots & a_m^T b_p \end{pmatrix} \quad (10)$$

Hence, each (i, j) -th entry of C is equal to the inner product of the i -th row of A and the j -th column of B . Compactly:

$$C_{ij} = a_i^T b_j = \sum_{k=1}^n a_{ik} b_{kj} \quad (11)$$

1.2.4 The Transpose

The transpose of a matrix $A \in \mathbb{R}^{m \times n}$ is $A^T \in \mathbb{R}^{n \times m}$ matrix whose entries are:

$$(A^T)_{ij} = A_{ji} \quad (12)$$

Properties of the transpose:

1. $(A^T)^T = A$
2. $(AB)^T = B^T A^T$
3. $(A + B)^T = A^T + B^T$

1.2.5 The Trace

The trace of a square matrix $A \in \mathbb{R}^{n \times n}$ is denoted $\text{tr}(A)$. It is the sum of diagonal elements in the matrix:

$$\text{tr } A = \sum_{i=1}^n A_{ii} \quad (13)$$

Properties of the trace:

1. For $A \in \mathbb{R}^{n \times n}$, $\text{tr } A = \text{tr } A^T$
2. For $A, B \in \mathbb{R}^{n \times n}$, $\text{tr}(A + B) = \text{tr } A + \text{tr } B$
3. For $A \in \mathbb{R}^{n \times n}$, $t \in \mathbb{R}$, $\text{tr}(tA) = t \cdot \text{tr } A$
4. For A, B such that AB is square, $\text{tr } AB = \text{tr } BA$
5. For A, B, C such that ABC is square, $\text{tr } ABC = \text{tr } BCA = \text{tr } CAB$, and so on

1.2.6 The Inverse

The inverse of a matrix A is noted A^{-1} and is the unique matrix such that:

$$A^{-1}A = I = AA^{-1} \quad (14)$$

Not all square matrices are invertible. In addition, assuming $A, B \in \mathbb{R}^{n \times n}$ are non-singular:

1. $(A^{-1})^{-1} = A$
2. $(AB)^{-1} = B^{-1}A^{-1}$
3. $(A^{-1})^T = (A^T)^{-1}$

1.2.7 The Determinant

The determinant of a square matrix $A \in \mathbb{R}^{n \times n}$, noted $|A|$ or $\det(A)$ is expressed recursively in terms of $A_{\setminus i, \setminus j}$, which is the matrix A without its i -th row and j -th column, as follows:

$$\det(A) = |A| = \sum_{j=1}^n (-1)^{i+j} A_{i,j} |A_{\setminus i, \setminus j}| \quad (15)$$

Remark that A is invertible if and only if $|A| \neq 0$. Also, $|AB| = |A||B|$ and $|A^T| = |A|$.

1.3 Matrix Properties

1.3.1 Norms

A norm of a vector x is any function $f : \mathbb{R}^n \rightarrow \mathbb{R}$ that satisfies 4 properties:

1. Non-negativity: For all $x \in \mathbb{R}^n$, $f(x) \geq 0$
2. Definiteness: $f(x) = 0$ if and only if $x = 0$
3. Homogeneity: For all $x \in \mathbb{R}^n$, $t \in \mathbb{R}$, $f(tx) = |t|f(x)$
4. Triangle Inequality: For all $x, y \in \mathbb{R}^n$, $f(x + y) \leq f(x) + f(y)$

However, most norms used come from the family of ℓ_p norms:

$$\ell_p = ||x||_p = \left(\sum_{i=1}^n x_i^p \right)^{\frac{1}{p}} \quad (16)$$

The p -norm is used in Holder's inequality. The Manhattan norm, used in LASSO regularization, is:

$$\ell_1 = ||x||_1 = \sum_{i=1}^n |x_i| \quad (17)$$

The euclidean norm, ℓ_2 is used in ridge regularization and distance measures:

$$\ell_2 = ||x||_2 = \sqrt{\sum_{i=1}^n x_i^2} \quad (18)$$

Finally, the infinity norm is used in uniform convergence:

$$\ell_\infty = ||x||_\infty = \max_i |x_i| \quad (19)$$

The Frobenius norm of a matrix is analogous to the ℓ_2 norm of a vector:

$$||A||_F = \sqrt{\sum_{ij} A_{ij}^2} \quad (20)$$

The dot product of two vectors can be expressed in terms of norms:

$$x^T y = ||x||_2 ||y||_2 \cos \theta \quad (21)$$

where θ is the angle between vectors x and y .

1.3.2 Linear Dependence and Rank

A set of vectors $\{x_1, x_2, \dots, x_n\} \subset \mathbb{R}^m$ is linearly independent if no vector can be represented as a linear combination of the remaining vectors. Conversely, if one vector in the set can be represented as a linear combination of the remaining vectors, then the vectors are said to be linearly dependent. Formally, for scalar values $\alpha_1, \dots, \alpha_{n-1} \in \mathbb{R}$:

$$x_n = \sum_{i=1}^{n-1} \alpha_i x_i \quad (22)$$

The rank of a given matrix A is noted $\text{rank}(A)$ and is the dimension of the vector space generated by its columns. This is equivalent to the maximum number of linearly independent columns of A . If $\text{rank}(A) = \min(m, n)$, then A is said to be full rank.

1.3.3 Span, Range, and Nullspace

The span of a set of vectors $\{x_1, x_2, \dots, x_n\}$ is the set of all vectors that can be expressed as a linear combination of $\{x_1, x_2, \dots, x_n\}$. Formally,

$$\text{span}(\{x_1, \dots, x_n\}) = \left\{ v : v = \sum_{i=1}^n \alpha_i x_i, \quad \alpha_i \in \mathbb{R} \right\} \quad (23)$$

The projection of a vector $y \in \mathbb{R}^n$ onto the span of $\{x_1, x_2, \dots, x_n\}$ is the vector $v \in \text{span}(\{x_1, \dots, x_n\})$ such that v is as close as possible to y as measured by the euclidean norm:

$$\text{Proj}(y; \{x_1, \dots, x_n\}) = \underset{v \in \text{span}(\{x_1, \dots, x_n\})}{\text{argmin}} \|y - v\|_2 \quad (24)$$

The range of a matrix $A \in \mathbb{R}^{m \times n}$ is the span of the columns of A :

$$\mathcal{R}(A) = \{v \in \mathbb{R}^m : v = Ax, x \in \mathbb{R}^n\} \quad (25)$$

The nullspace of a matrix $A \in \mathbb{R}^{m \times n}$ is the set of all vectors that equal 0 when multiplied by A :

$$\mathcal{N}(A) = \{x \in \mathbb{R}^n : Ax = 0\} \quad (26)$$

1.3.4 Symmetric Matrices

A square matrix $A \in \mathbb{R}^{n \times n}$ is symmetric if $A = A^T$. It is anti-symmetric if $A = -A^T$. For any matrix $A \in \mathbb{R}^{n \times n}$ the matrix $A + A^T$ is symmetric and the matrix $A - A^T$ is anti-symmetric. From this, any square matrix $A \in \mathbb{R}^{n \times n}$ can be represented as a sum of a symmetric matrix and an anti-symmetric matrix:

$$A = \underbrace{\frac{A + A^T}{2}}_{\text{Symmetric}} + \underbrace{\frac{A - A^T}{2}}_{\text{Antisymmetric}} \quad (27)$$

1.3.5 Positive Semidefinite Matrices

Given a square matrix $A \in \mathbb{R}^{n \times n}$ and a vector $x \in \mathbb{R}^n$, the scalar value $x^T Ax$ is called the quadratic form:

$$x^T Ax = \sum_{i=1}^n \sum_{j=1}^n A_{ij} x_i x_j \quad (28)$$

A symmetric matrix A is:

1. Positive definite (PD) if for all nonzero vectors, $x^T Ax > 0$
2. Positive semidefinite (PSD) if for all nonzero vectors, $x^T Ax \geq 0$
3. Negative definite (ND) if for all nonzero vectors, $x^T Ax < 0$
4. Negative semidefinite (NSD) if for all nonzero vectors, $x^T Ax \leq 0$
5. Indefinite if it is neither PSD nor NSD, i.e. if there exists a x_1, x_2 such that $x_1^T Ax_1 > 0$ and $x_2^T Ax_2 < 0$

Given any matrix $A \in \mathbb{R}^{m \times n}$, the gram matrix $G = A^T A$ is always PSD. If $m \geq n$ and A is full rank, then G is PD.

1.3.6 Eigendecomposition

Given a square matrix $A \in \mathbb{R}^{n \times n}$, we say that $\lambda \in \mathbb{C}$ is an eigenvalue of A and $v \in \mathbb{C}$ is the corresponding eigenvector if

$$Av = \lambda v, \quad v \neq 0 \quad (29)$$

The trace of A is equal to the sum of its eigenvalues:

$$\text{tr } A = \sum_{i=1}^n \lambda_i \quad (30)$$

The determinant of A is equal to the product of its eigenvalues:

$$|A| = \prod_{i=1}^n \lambda_i \quad (31)$$

Suppose matrix A has n linearly independent eigenvectors $\{v_1, \dots, v_n\}$ with corresponding eigenvalues $\{\lambda_1, \dots, \lambda_n\}$. Let V be a matrix with one eigenvalue per column: $V = [v_1, \dots, v_n]$. Let $\Lambda = \text{diag}(\lambda_1, \dots, \lambda_n)$. Then the eigendecomposition of A is given by:

$$A = V\Lambda V^{-1} \quad (32)$$

For when A is symmetric, then the eigenvalues of A are real, the eigenvectors of A are orthonormal, and hence V is an orthogonal matrix (henceforth renamed U). Hence $A = U\Lambda U^T$. A matrix whose eigenvalues are all positive is called positive definite. A matrix whose eigenvalues are all positive or zero valued is called positive semidefinite. Likewise, if all eigenvalues are negative, the matrix is negative definite, and if all eigenvalues are negative or zero valued, it is negative semidefinite. In addition, assuming sorted eigenvalues, for the following optimization problem:

$$\max_{x \in \mathbb{R}^n} x^T A x \quad \text{subject to } \|x\|_2^2 = 1 \quad (33)$$

The solution is v_1 the eigenvector corresponding to λ_1 . For the minimization problem:

$$\min_{x \in \mathbb{R}^n} x^T A x \quad \text{subject to } \|x\|_2^2 = 1 \quad (34)$$

the optimal solution for x is v_n , the eigenvector corresponding to eigenvalue λ_n .

1.3.7 Singular Value Decomposition

The singular value decomposition provides a way to factorize a $m \times n$ matrix A into singular vectors and singular values. It is defined as:

$$A = UDV^T \quad (35)$$

Suppose $A \in \mathbb{R}^{m \times n}$ matrix. Then $U \in \mathbb{R}^{m \times m}$, $D \in \mathbb{R}^{m \times n}$, $V \in \mathbb{R}^{n \times n}$ matrices. The matrices U and V are orthogonal, and matrix D is diagonal. The elements along the diagonal of D are known as the singular values of matrix A . The columns of U are known as the left-singular vectors while the columns of V are the right-singular vectors. The left-singular vectors of A are the eigenvectors of AA^T . The right-singular vectors of A are the eigenvectors of $A^T A$. We can use SVD to partially generalize matrix inversion to nonsquare matrices.

1.3.8 The Moore-Penrose Pseudoinverse

Matrix inversion is not defined for matrices that are not square. Note that for nonsingular A :

$$AA^{-1}A = A \quad (36)$$

However, if the inverse is not defined, we seek to find a matrix A^+ such that:

$$AA^+A = A \quad (37)$$

The moore-penrose pseudoinverse A^+ is defined as follows:

$$A^+ = \lim_{\alpha \rightarrow 0} (A^T A + \alpha I)^{-1} A^T \quad (38)$$

Practical algorithms use the singular value decomposition of A such that:

$$A^+ = VD^+U^T \quad (39)$$

where U, D, V are from the SVD of A , and the pseudoinverse of D^+ is obtained by taking the reciprocal of the nonzero diagonal elements of D . If A has more columns than rows, then using the pseudoinverse to solve a linear equation $Ax = y$ provides one of many solutions, but provides $x = A^+y$ with minimal euclidean norm $\|x\|_2$. When A has more rows than columns, the pseudoinverse gives us the x for which Ax is as close as possible to y , i.e. minimizing $\|Ax - y\|_2$.

1.4 Matrix Calculus

1.4.1 The Gradient

Let $f : \mathbb{R}^{m \times n} \rightarrow \mathbb{R}$ be a function and $A \in \mathbb{R}^{m \times n}$ be a matrix. The gradient of f with respect to A is a $m \times n$ matrix noted as $\nabla_A f(A)$ such that:

$$\nabla_A f(A) \in \mathbb{R}^{m \times n} = \begin{pmatrix} \frac{\partial f(A)}{\partial A_{11}} & \frac{\partial f(A)}{\partial A_{12}} & \cdots & \frac{\partial f(A)}{\partial A_{1n}} \\ \frac{\partial f(A)}{\partial A_{21}} & \frac{\partial f(A)}{\partial A_{22}} & \cdots & \frac{\partial f(A)}{\partial A_{2n}} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial f(A)}{\partial A_{m1}} & \frac{\partial f(A)}{\partial A_{m2}} & \cdots & \frac{\partial f(A)}{\partial A_{mn}} \end{pmatrix} \quad (40)$$

Or compactly for each ij entry:

$$\nabla_A f(A)_{ij} = \frac{\partial f(A)}{\partial A_{ij}} \quad (41)$$

However, the gradient of a vector $x \in \mathbb{R}^n$ is:

$$\nabla_x f(x) = \begin{pmatrix} \frac{\partial f(x)}{\partial x_1} \\ \frac{\partial f(x)}{\partial x_2} \\ \vdots \\ \frac{\partial f(x)}{\partial x_n} \end{pmatrix} \quad (42)$$

1.4.2 The Hessian

Let $f : \mathbb{R}^n \rightarrow \mathbb{R}$ be a function and $x \in \mathbb{R}^n$ be a vector. The hessian of f with respect to x is a $n \times n$ symmetric matrix noted as $H = \nabla_x^2 f(x)$ such that:

$$\nabla_x^2 f(x) \in \mathbb{R}^{n \times n} = \begin{pmatrix} \frac{\partial^2 f(x)}{\partial x_1^2} & \frac{\partial^2 f(x)}{\partial x_1 \partial x_2} & \cdots & \frac{\partial^2 f(x)}{\partial x_1 \partial x_n} \\ \frac{\partial^2 f(x)}{\partial x_2 \partial x_1} & \frac{\partial^2 f(x)}{\partial x_2^2} & \cdots & \frac{\partial^2 f(x)}{\partial x_2 \partial x_n} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial^2 f(x)}{\partial x_n \partial x_1} & \frac{\partial^2 f(x)}{\partial x_n \partial x_2} & \cdots & \frac{\partial^2 f(x)}{\partial x_n^2} \end{pmatrix} \quad (43)$$

Or compactly:

$$\nabla_x^2 f(x)_{ij} = \frac{\partial^2 f(x)}{\partial x_i \partial x_j} \quad (44)$$

Note that the hessian is only defined when $f(x)$ is real-valued.

1.4.3 Gradient Properties

For matrices A, B, C and vectors x, b :

1. $\nabla_x b^T x = b$
2. $\nabla_x x^T A x = 2Ax$ (if A symmetric)
3. $\nabla_x^2 x^T A x = 2A$ (if A symmetric)
4. $\nabla_A \text{tr}(AB) = B^T$
5. $\nabla_{A^T} f(A) = (\nabla_A f(A))^T$
6. $\nabla_A \text{tr}(ABA^T C) = CAB + C^T AB^T$
7. $\nabla_A |A| = |A|(A^{-1})^T$

2 Convex Optimization

2.1 Convexity

2.1.1 Convex Sets

A set C is convex if, for any $x, y \in C$ and $\alpha \in \mathbb{R}$ with $0 \leq \alpha \leq 1$, we have

$$\alpha x + (1 - \alpha)y \in C \quad (45)$$

This point is known as a convex combination of x and y .

2.1.2 Convex Functions

A function $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is convex if its domain $\mathcal{D}(f)$ is a convex set, and if, for all $x, y \in \mathcal{D}(f)$ and $\alpha \in \mathbb{R}$, $0 \leq \alpha \leq 1$, we have:

$$f(\alpha x + (1 - \alpha)y) \leq \alpha f(x) + (1 - \alpha)f(y) \quad (46)$$

A function is strictly convex if:

$$f(\alpha x + (1 - \alpha)y) < \alpha f(x) + (1 - \alpha)f(y) \quad (47)$$

2.1.3 First-Order Conditions

Suppose a function $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is differentiable, then f is convex if and only if $\mathcal{D}(f)$ is a convex set and for all $x, y \in \mathcal{D}(f)$ we have:

$$f(y) \geq f(x) + \nabla_x f(x)^T (y - x) \quad (48)$$

This is called the first-order approximation to the function f at point x . This is approximating f with its tangent line at point x . The first order condition for convexity says that f is convex if and only if the tangent line is a global underestimator of the function f . In other words, if we take our function and draw a tangent line at any point, then every point on this line will lie below the corresponding point on f .

2.1.4 Second-Order Conditions

Suppose a function $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is twice differentiable, i.e. the hessian $\nabla_x^2 f(x)$ is defined for all points x in the domain of f . Then f is convex if and only if $\mathcal{D}(f)$ is a convex set and its hessian is PSD:

$$\nabla_x^2 f(y) \succeq 0 \quad (49)$$

2.1.5 Jensen's Inequality

Suppose we start with the inequality in the basic definition of a convex function:

$$f(\alpha x + (1 - \alpha)y) \leq \alpha f(x) + (1 - \alpha)f(y) \quad (50)$$

We can extend this to convex combinations of more than one point:

$$f\left(\sum_{i=1}^k \alpha_i x_i\right) \leq \sum_{i=1}^k \alpha_i f(x_i) \quad (51)$$

For $\sum_k \alpha = 1$ and all $\alpha \geq 0$. This can be extended to integrals:

$$f\left(\int p(x) x dx\right) \leq \int p(x) f(x) dx \quad (52)$$

For $\int p(x) dx = 1$ and $p(x) \geq 0 \forall x$. This implies $p(x)$ is a probability density, and we can write our inequality in terms of expectations:

$$f(\mathbf{E}[x]) \leq \mathbf{E}[f(x)] \quad (53)$$

And is known as Jensen's inequality.

2.1.6 Sublevel Sets

Given a convex function $f : \mathbb{R}^n \rightarrow \mathbb{R}$ and a real number $\alpha \in \mathbb{R}$, the α -sublevel set is:

$$\{x \in \mathcal{D}(f) : f(x) \leq \alpha\} \quad (54)$$

In other words, the α -sublevel set is the set of all points x such that $f(x) \leq \alpha$.

2.2 Convex Optimization

A convex optimization problem seeks to optimize a convex function f with respect to a variable x and possible constraints.

$$\begin{aligned} & \text{minimize} && f(x) \\ & \text{subject to} && g_i(x) \leq 0, \quad i = 1, \dots, m \\ & && h_i(x) = 0, \quad i = 1, \dots, p \end{aligned} \quad (55)$$

where f is a convex function, g_i are convex functions, and h_i are affine functions. The optimal value p^* is equal to the minimum possible value of the objective function in the feasible region:

$$p^* = \min \{f(x) : g_i(x) \leq 0, i = 1, \dots, m, h_i(x) = 0, i = 1, \dots, p\} \quad (56)$$

The optimal point x^* is a point such that $f(x^*) = p^*$.

2.2.1 Global Optimality

A point x is locally optimal if it is feasible and if there exists some $R > 0$ such that all feasible points z with $\|x - z\|_2 \leq R$, satisfy $f(x) \leq f(z)$. This means that for x to be locally optimal, there exists no nearby points that have a lower objective value. A point x is globally optimal if it is feasible and for all feasible points z , $f(x) \leq f(z)$. For a convex optimization problem, all locally optimal points are globally optimal.

2.2.2 Gradient Descent

Using $\alpha \in \mathbb{R}$ as the learning rate, we can update a set of parameters θ with respect to minimizing a function f as follows:

$$\theta := \theta - \alpha \nabla f(\theta) \quad (57)$$

Stochastic gradient descent (SGD) is updating the parameter based on each training example, and batch gradient descent is on a batch of training examples.

2.2.3 Newton's Algorithm

Newton's algorithm is a numerical method using information from the second derivative to find θ such that $f'(\theta) = 0$.

$$\theta := \theta - \frac{f'(\theta)}{f''(\theta)} \quad (58)$$

For multidimensional parameters:

$$\theta := \theta - \alpha H^{-1} \nabla_{\theta} f(\theta) \quad (59)$$

Where H is the hessian matrix of second partial derivatives.

$$H_{ij} = \frac{\partial^2 f(\theta)}{\partial \theta_i \partial \theta_j} \quad (60)$$

2.3 Lagrange Duality and KKT Conditions

Generic differentiable convex optimization problems are of the form:

$$\begin{aligned} & \text{minimize} && f(x) \\ & \text{subject to} && g_i(x) \leq 0, \quad i = 1, \dots, m \\ & && h_i(x) = 0, \quad i = 1, \dots, p \end{aligned} \quad (61)$$

where $x \in \mathbb{R}^n$ is the optimization variable, $f : \mathbb{R}^n \rightarrow \mathbb{R}$ and $g_i : \mathbb{R}^n \rightarrow \mathbb{R}$ are differentiable convex functions, and $h_i : \mathbb{R}^n \rightarrow \mathbb{R}$ are affine functions.

2.3.1 The Lagrangian

Given a convex constrained minimization problem, the generalized lagrangian is a function $\mathcal{L} : \mathbb{R}^n \times \mathbb{R}^m \times \mathbb{R}^p \rightarrow \mathbb{R}$ defined as:

$$\mathcal{L}(x, \alpha, \beta) = f(x) + \sum_{i=1}^m \alpha_i g_i(x) + \sum_{i=1}^p \beta_i h_i(x) \quad (62)$$

We refer to x as the primal variables of the Lagrangian. The α and β are referred to as the dual variables or lagrange multipliers. The key idea behind Lagrangian duality is for any convex optimization problem, there always exist settings of the dual variables such that the unconstrained minimum of the Lagrangian with respect to the primal variables (keeping the dual variables fixed) coincides with the solution of the original constrained minimization problem.

2.3.2 Primal and Dual Problems

The primal problem is:

$$\min_x \underbrace{\left[\max_{\alpha, \beta: \alpha_i \geq 0, \forall i} \mathcal{L}(x, \alpha, \beta) \right]}_{\theta_{\mathcal{P}(x)}} = \min_x \theta_{\mathcal{P}(x)} \quad (63)$$

Here, $\theta_{\mathcal{D}} : \mathbb{R}^n \rightarrow \mathbb{R}$ is the primal objective, and the right hand side is the primal problem. In addition, $p^* = \theta_{\mathcal{D}}(x^*)$ is the optimal value of the primal objective.

The dual problem is:

$$\max_{\alpha, \beta: \alpha_i \geq 0, \forall i} \underbrace{\left[\min_x \mathcal{L}(x, \alpha, \beta) \right]}_{\theta_{\mathcal{D}}(x)} = \max_{\alpha, \beta: \alpha_i \geq 0, \forall i} \theta_{\mathcal{D}}(x) \quad (64)$$

Here, $\theta_{\mathcal{D}} : \mathbb{R}^m \times \mathbb{R}^p \rightarrow \mathbb{R}$ is the dual objective, and the right hand side is the dual problem. In addition, $s^* = \theta_{\mathcal{D}}(\alpha^*, \beta^*)$ is the optimal value of the dual objective.

2.3.3 Strong and Weak Duality

Weak duality for any pair of primal and dual problems, with d^* the optimal objective value of the dual problem and p^* as the optimal objective value of the primal problem we have:

$$d^* \leq p^* \quad (65)$$

Strong duality is when for any pair of primal and dual problems which satisfy certain technical conditions called constraint qualifications, then:

$$d^* = p^* \quad (66)$$

2.3.4 Complementary Slackness

If strong duality holds, then $\alpha_i^* g(x_i^*) = 0$ for each $i = 1, \dots, m$. This property is known as complementary slackness. This implies the following:

$$\alpha_i^* > 0 \implies g_i(x^*) = 0 \quad (67)$$

$$g_i(x^*) < 0 \implies \alpha_i^* = 0 \quad (68)$$

2.3.5 The KKT Conditions

Suppose that $x^* \in \mathbb{R}^n$, $\alpha^* \in \mathbb{R}^m$ and $\beta^* \in \mathbb{R}^p$ satisfy the following conditions:

1. Primal feasibility: $g_i(x^*) \leq 0$, for $i = 1, \dots, m$ and $h_i(x^*) = 0$ for $i = 1, \dots, p$
2. Dual feasibility: $\alpha_i^* \geq 0$ for $i = 1, \dots, m$
3. Complementary slackness: $\alpha_i^* g_i(x^*) = 0$ for $i = 1, \dots, m$
4. Lagrangian stationarity: $\nabla_x \mathcal{L}(x^*, \alpha^*, \beta^*) = \vec{0}$

Then x^* is primal optimal and (α^*, β^*) are dual optimal. Furthermore, if strong duality holds, then any primal optimal x^* and dual optimal (α^*, β^*) must satisfy the conditions 1 through 4. These conditions are known as the Karush-Kuhn-Tucker (KKT) conditions.

3 Probability and Statistics

3.1 Basics

The set of all possible outcomes of an experiment is known as the sample space and denoted by S . Any subset E of the sample space is known as an event. An event is a set consisting of possible outcomes of the experiment.

3.1.1 Axioms of Probability

For each event E , we denote $P(E)$ as the probability of event E occurring. $P(E)$ satisfies the following properties:

1. Every probability is between 0 and 1 included:

$$0 \leq P(E) \leq 1 \quad (69)$$

2. The probability that at least one event in the sample space will occur is 1:

$$P(S) = 1 \quad (70)$$

3. For any sequence of mutually exclusive events E_1, \dots, E_n we have:

$$P\left(\bigcup_{i=1}^n E_i\right) = \sum_{i=1}^n P(E_i) \quad (71)$$

3.1.2 Permutation

A permutation is an arrangement of r objects from a pool of n objects, in a given order. The number of such arrangements is given by $P(n, r)$:

$$P(n, r) = \frac{n!}{(n-r)!} \quad (72)$$

3.1.3 Combination

A combination is an arrangement of r objects from a pool of n objects, where order does not matter. The number of such arrangements is given by $C(n, r)$:

$$C(n, r) = \frac{P(n, r)}{r!} = \frac{n!}{r!(n-r)!} \quad (73)$$

Note that for $0 \leq r \leq n$ we have $P(n, r) \geq C(n, r)$.

3.2 Conditional Probability

Let B be an event with non-zero probability, the conditional probability of any event A given B is:

$$P(A|B) = \frac{P(A \cap B)}{P(B)} \quad (74)$$

3.2.1 Bayes Rule

For events A, B such that $P(B) > 0$, we have:

$$P(A|B) = \frac{P(B|A)P(A)}{P(B)} \quad (75)$$

From Bayes rule, we have:

$$P(A \cap B) = P(A|B)P(B) = P(A)P(B|A) \quad (76)$$

Let $\{A_i, i \in [1, n]\}$ be such that for all i , $A_i \neq \emptyset$. We say that $\{A_i\}$ is a partition if we have:

$$\forall i \neq j, A_i \cap A_j = \emptyset \quad \text{and} \quad \bigcup_{i=1}^n A_i = S \quad (77)$$

Remark that for any event B in the sample space, we have:

$$P(B) = \sum_{i=1}^n P(B|A_i)P(A_i) \quad (78)$$

Let $\{A_i, i \in [1, n]\}$ be a partition of the sample space, we can extend bayes rule as:

$$P(A_k|B) = \frac{P(B|A_k)P(A_k)}{\sum_{i=1}^n P(B|A_i)P(A_i)} \quad (79)$$

3.2.2 Independence

Two events A and B are independent if and only if we have:

$$P(A \cap B) = P(A)P(B) \quad (80)$$

3.3 Random Variables

A random variable X is a function that maps every element in a sample space to a real line.

3.3.1 Cumulative Distribution Function (CDF)

The cumulative distribution function F , which is monotonically non-decreasing and is such that $\lim_{x \rightarrow -\infty} F(X) = 0$ and $\lim_{x \rightarrow \infty} F(X) = 1$ is defined as:

$$F(x) = P(X \leq x) \quad (81)$$

In addition:

$$P(a < X \leq b) = F(b) - F(a) \quad (82)$$

3.3.2 Probability Density Function (PDF)

The probability density function is the derivative of the CDF. It has the following properties:

1. $f(x) \geq 0$
2. $\int_{-\infty}^{\infty} f(x) = 1$
3. $\int_{x \in A} f(x)dx = P(X \in A)$

3.3.3 Discrete PDF/CDF

If X is discrete, by denoting f as the PDF and F as the CDF, we have:

$$F(X) = \sum_{x_i \leq x} P(X = x_i) \quad (83)$$

$$f(x_j) = P(X = x_j) \quad (84)$$

And the following properties for the PDF:

$$0 \leq f(x_j) \leq 1 \quad (85)$$

$$\sum_j f(x_j) = 1 \quad (86)$$

3.3.4 Continuous PDF/CDF

If X is continuous, by denoting f as the PDF and F as the CDF, we have:

$$F(X) = \int_{-\infty}^x f(y)dy \quad (87)$$

$$f(x) = \frac{dF}{dx} \quad (88)$$

And the following properties for the PDF:

$$f(x) \geq 0 \quad (89)$$

$$\int_{-\infty}^{\infty} f(x)dx = 1 \quad (90)$$

3.3.5 Expectation

The expected value of a random variable, also known as the mean value or first moment, is denoted as $\mathbb{E}[X]$ or μ . It is the value obtained by averaging the results of a random variable. We use (D) for discrete, (C) for continuous.

$$(D) \quad \mathbb{E}[X] = \sum_{i=1}^n x_i f(x_i) \quad \text{and} \quad (C) \quad \mathbb{E}[X] = \int_{-\infty}^{+\infty} x f(x)dx \quad (91)$$

The expected value of a function of a random variable $g(X)$ is:

$$(D) \quad \mathbb{E}[g(X)] = \sum_{i=1}^n g(x_i) f(x_i) \quad \text{and} \quad (C) \quad \mathbb{E}[g(X)] = \int_{-\infty}^{+\infty} g(x) f(x)dx \quad (92)$$

The k -th moment, noted $\mathbb{E}[X^k]$ is the value of X^k that we expect to observe on average on infinitely many trials. The k -th moment is a case of the previous definition with $g : X \mapsto X^k$.

$$(D) \quad \mathbb{E}[X^k] = \sum_{i=1}^n x_i^k f(x_i) \quad \text{and} \quad (C) \quad \mathbb{E}[X^k] = \int_{-\infty}^{+\infty} x^k f(x)dx \quad (93)$$

A characteristic function $\psi(\omega)$ is derived from a probability density function $f(x)$ and is defined as:

$$(D) \quad \psi(\omega) = \sum_{i=1}^n f(x_i) e^{i\omega x_i} \quad \text{and} \quad (C) \quad \psi(\omega) = \int_{-\infty}^{+\infty} f(x) e^{i\omega x} dx \quad (94)$$

Remark that $e^{i\omega x} = \cos(\omega x) + i \sin(\omega x)$. The k -th moment can also be computed with the characteristic function as:

$$\mathbb{E}[X^k] = \frac{1}{i^k} \left[\frac{\partial^k \psi}{\partial \omega^k} \right]_{\omega=0} \quad (95)$$

3.3.6 Variance and Standard Deviation

The variance of a random variable, often noted $\text{Var}(X)$ or σ^2 , is a measure of the spread of its distribution function. It is determined as follows:

$$\text{Var}[X] = \mathbb{E}[(X - \mathbb{E}[X])^2] = \mathbb{E}[X^2] - \mathbb{E}[X]^2 \quad (96)$$

The standard deviation of a random variable, often noted σ , is a measure of the spread of its distribution function which is compatible with the units of the actual random variable. It is determined as follows:

$$\sigma = \sqrt{\text{Var}[X]} \quad (97)$$

Note that the variance for any constant a is $\text{Var}[a] = 0$, and $\text{Var}[af(X)] = a^2 \text{Var}[f(X)]$.

3.4 Discrete Random Variables

3.4.1 Bernoulli

For $X \sim \text{Bernoulli}(p)$, we define a binary event with a probability of p for a true event, and a false event with probability of $q = 1 - p$.

$$P(X = x) = \begin{cases} q = 1 - p & \text{if } x = 0 \\ p & \text{if } x = 1 \end{cases} \quad (98)$$

It can also be expressed as:

$$f(x; p) = p^k(1 - p)^{1-k} \quad \text{for } k \in \{0, 1\} \quad (99)$$

Other properties:

$$\mathbb{E}[X] = p \quad (100)$$

$$\text{Var}[X] = pq = p(1 - p) \quad (101)$$

$$\psi(\omega) = (1 - p) + pe^{i\omega} \quad (102)$$

3.4.2 Binomial

For $X \sim \text{Binomial}(n, p)$, the number of true events in n independent experiments, with true probability of p , false probability of $q = 1 - p$.

$$P(X = x) = \binom{n}{x} p^x (1 - p)^{n-x} \quad (103)$$

Other properties:

$$\mathbb{E}[X] = np \quad (104)$$

$$\text{Var}[X] = npq = np(1 - p) \quad (105)$$

$$\psi(\omega) = (pe^{i\omega} + (1 - p))^n \quad (106)$$

3.4.3 Geometric

For $X \sim \text{Geometric}(p)$, is the number of experiments with true probability of p until the first true event (number of trials to get one success).

$$P(X = x) = p(1 - p)^{x-1} \quad (107)$$

Other properties:

$$\mathbb{E}[X] = \frac{1}{p} \quad (108)$$

$$\text{Var}[X] = \frac{1 - p}{p^2} \quad (109)$$

$$\psi(\omega) = \frac{pe^{i\omega}}{1 - (1 - p)e^{i\omega}} \quad (110)$$

3.4.4 Poisson

For $X \sim \text{Poisson}(\lambda)$, for $\lambda > 0$, a probability distribution over the nonnegative integers used for modeling the frequency of rare events.

$$P(X = x) = \frac{\lambda^x}{x!} e^{-\lambda} \quad (111)$$

Other properties:

$$\mathbb{E}[X] = \lambda \quad (112)$$

$$\text{Var}[X] = \lambda \quad (113)$$

$$\psi(\omega) = e^{\lambda(e^{i\omega} - 1)} \quad (114)$$

3.5 Continuous Random Variables

3.5.1 Uniform

For $X \sim \text{Uniform}(a, b)$, we have equal probability density to every value between a and b .

$$f(x) = \begin{cases} \frac{1}{b-a} & \text{if } a \leq x \leq b \\ 0 & \text{otherwise} \end{cases} \quad (115)$$

Other properties:

$$\mathbb{E}[X] = \frac{a+b}{2} \quad (116)$$

$$\text{Var}[X] = \frac{(b-a)^2}{12} \quad (117)$$

$$\psi(\omega) = \frac{e^{i\omega b} - e^{i\omega a}}{(b-a)i\omega} \quad (118)$$

3.5.2 Exponential

For $X \sim \text{Exponential}(\lambda)$, $\lambda > 0$, is the decaying probability density over the nonnegative reals.

$$f(x) = \begin{cases} \lambda e^{-\lambda x} & \text{if } x \geq 0 \\ 0 & \text{otherwise} \end{cases} \quad (119)$$

Other properties:

$$\mathbb{E}[X] = \frac{1}{\lambda} \quad (120)$$

$$\text{Var}[X] = \frac{1}{\lambda^2} \quad (121)$$

$$\psi(\omega) = \frac{1}{1 - \frac{i\omega}{\lambda}} \quad (122)$$

3.5.3 Gaussian (Normal)

For $X \sim \text{Normal}(\mu, \sigma)$, denoted also $X \sim \mathcal{N}(\mu, \sigma)$.

$$f(x) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{1}{2}\left(\frac{x-\mu}{\sigma}\right)^2} \quad (123)$$

Other properties:

$$\mathbb{E}[X] = \mu \quad (124)$$

$$\text{Var}[X] = \sigma^2 \quad (125)$$

$$\psi(\omega) = e^{i\omega\mu - \frac{1}{2}\omega^2\sigma^2} \quad (126)$$

3.6 Jointly Distributed Random Variables

3.6.1 Marginal Density

3.6.2 Cumulative Distribution

3.6.3 Conditional Density

3.6.4 Independence

3.6.5 Expectation and Covariance

3.6.6 Correlation

3.7 Parameter Estimation

3.7.1 Definitions

3.7.2 Bias

3.7.3 Mean and Central Limit Theorem

3.7.4 Variance

3.8 Probability Bounds and Inequalities

3.8.1 Markov

3.8.2 Chebyshev

3.8.3 Chernoff

3.8.4 Hoeffding

4 Information Theory

Information Theory revolves around quantifying how much information is present in a signal. The basic intuition lies in the fact that learning an unlikely event has occurred is more informative than learning that a likely event has occurred. The basics are:

1. Likely events should have low information content, and in the extreme case, events that are guaranteed to happen should have no information content whatsoever.
2. Less likely events should have higher information content.

3. Independent events should have additive information.

We satisfy all three properties by defining self-information of an event x for a probability distribution P as:

$$I(x) = -\log P(x) \quad (127)$$

We can quantify the amount of uncertainty in a distribution using Shannon entropy:

$$H(P) = \mathbb{E}_{x \sim P}[I(x)] = -\mathbb{E}_{x \sim P}[\log P(x)] \quad (128)$$

Which in the discrete setting is written as:

$$H(P) = -\sum_x P(x) \log P(x) \quad (129)$$

In other words, the Shannon entropy of a distribution is the expected amount of information in an event drawn from that distribution. It gives a lower bound on the number of bits needed on average to encode symbols drawn from a distribution P . If we have two separate probability distributions $P(x)$ and $Q(x)$ over the same random variable x , we can measure how different these two distributions are using the Kullback-Leibler (KL) divergence:

$$\begin{aligned} D_{\text{KL}}(P||Q) &= \mathbb{E}_{x \sim P} \left[\log \frac{P(x)}{Q(x)} \right] \\ &= \mathbb{E}_{x \sim P} [\log P(x) - \log Q(x)] \\ &= \sum_x P(x) \frac{\log P(x)}{\log Q(x)} \end{aligned} \quad (130)$$

In the case of discrete variables, it is the extra amount of information needed to send a message containing symbols drawn from probability distribution P , when we use a code that was designed to minimize the length of messages drawn from probability distribution Q . The KL divergence is always non-negative, and is 0 if and only if P and Q are the same. We can relate the KL divergence to cross-entropy.

$$\begin{aligned} H(P, Q) &= H(P) + D_{\text{KL}}(P||Q) \\ &= -\mathbb{E}_{x \sim P} [\log Q(x)] \\ &= -\sum_x P(x) \log Q(x) \end{aligned} \quad (131)$$

Minimizing the cross-entropy with respect to Q is equivalent to minimizing the KL divergence, because Q does not participate in the omitted term (entropy is constant).

5 Learning Theory

5.1 Bias and Variance

5.2 Notation

5.2.1 Union Bound

5.2.2 Hoeffding Inequality

5.3 Training Error

5.4 Probably Approximately Correct (PAC)

5.5 Hypothesis Classes

5.5.1 Shattering

5.5.2 Upper Bound Theorem

5.5.3 VC Dimension

5.5.4 Vapnik Theorem

6 Linear Regression

Linear Regression seeks to approximate a real valued label y as a linear function of x :

$$h_{\theta}(x) = \theta_0 + \theta_1 \cdot x_1 + \cdots + \theta_n \cdot x_n \quad (132)$$

The θ_i 's are the parameters, or weights. If we include the intercept term via $x_0 = 1$, we can write our model more compactly as:

$$h(x) = \sum_{i=0}^n \theta_i \cdot x_i = \theta^T x \quad (133)$$

Here n is the number of input variables, or features. In Linear Regression, we seek to make $h(x)$ as close to y for a set of training examples. We define the cost function as:

$$J(\theta) = \frac{1}{2} \sum_{i=1}^m \left(h(x^{(i)}) - y^{(i)} \right)^2 \quad (134)$$

6.1 LMS Algorithm

We seek to find a set of θ such that we minimize $J(\theta)$ via a search algorithm that starts at some initial guess for our parameters and takes incremental steps to make $J(\theta)$ smaller until convergence. This is known as gradient descent:

$$\theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta) \quad (135)$$

Here, α is the learning rate. We can derive the partial derivative as:

$$\begin{aligned}
 \frac{\partial}{\partial \theta_j} J(\theta) &= \frac{\partial}{\partial \theta_j} \frac{1}{2} (h(x) - y)^2 \\
 &= 2 \cdot \frac{1}{2} (h(x) - y) \cdot \frac{\partial}{\partial \theta_j} (h(x) - y) \\
 &= (h(x) - y) \cdot \frac{\partial}{\partial \theta_j} \left(\sum_{i=0}^n \theta_i x_i - y \right) \\
 &= (h(x) - y) x_j
 \end{aligned} \tag{136}$$

Hence, for a single example (stochastic gradient descent):

$$\theta_j := \theta_j + \alpha \left(y^{(i)} - h(x^{(i)}) \right) x_j^{(i)} \tag{137}$$

This is called the LMS update rule. For a batched version, we can evaluate the gradient on a set of examples (batch gradient descent), or the full set (gradient descent).

$$\theta_j := \theta_j + \alpha \sum_{i=1}^m \left(y^{(i)} - h(x^{(i)}) \right) x_j^{(i)} \tag{138}$$

6.2 The Normal Equations

We can also directly minimize J without using an iterative algorithm. We define X as the matrix of all samples of size m by n . We let \vec{y} be a m dimensional vector of all target values. We can define our cost function J as:

$$J(\theta) = \frac{1}{2} (X\theta - \vec{y})^T (X\theta - \vec{y}) = \frac{1}{2} \sum_{i=1}^m \left(h(x^{(i)}) - y^{(i)} \right)^2 \tag{139}$$

We then take the derivative and find its roots.

$$\begin{aligned}
 \nabla_{\theta} J(\theta) &= \nabla_{\theta} \frac{1}{2} (X\theta - \vec{y})^T (X\theta - \vec{y}) \\
 &= \frac{1}{2} \nabla_{\theta} (\theta^T X^T X\theta - \theta^T X^T \vec{y} - \vec{y}^T X\theta + \vec{y}^T \vec{y}) \\
 &= \frac{1}{2} \nabla_{\theta} (\text{tr } \theta^T X^T X\theta - 2 \text{tr } \vec{y}^T X\theta) \\
 &= \frac{1}{2} (X^T X\theta + X^T X\theta - 2X^T \vec{y}) \\
 &= X^T X\theta - X^T \vec{y}
 \end{aligned} \tag{140}$$

To minimize J , we set its derivatives to zero, and obtain the normal equations:

$$X^T X\theta = X^T \vec{y} \tag{141}$$

Which solves θ for a value that minimizes $J(\theta)$ in closed form:

$$\theta = (X^T X)^{-1} X^T \vec{y} \tag{142}$$

6.3 Probabilistic Interpretation

Why does linear regression use the least-squares cost function? Assume that the target variables and inputs are related via:

$$y^{(i)} = \theta^T x^{(i)} + \epsilon^{(i)} \quad (143)$$

Here, $\epsilon^{(i)}$ is an error term for noise. We assume each $\epsilon^{(i)}$ is independently and identically distributed according to a Gaussian distribution with mean zero and some variance σ^2 . Hence, $\epsilon^{(i)} \sim \mathcal{N}(0, \sigma^2)$, so the density for any sample $x^{(i)}$ with label $y^{(i)}$ is $y^{(i)}|x^{(i)}; \theta \sim \mathcal{N}(\theta^T x^{(i)}, \sigma^2)$. This implies:

$$p(y^{(i)}|x^{(i)}; \theta) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(y^{(i)} - \theta^T x^{(i)})^2}{2\sigma^2}\right) \quad (144)$$

The probability of a dataset X is quantified by a likelihood function:

$$L(\theta) = L(\theta; X, \vec{y}) = p(\vec{y}|X; \theta) \quad (145)$$

Since we assume independence on each noise term (and samples), we can write the likelihood function as:

$$\begin{aligned} L(\theta) &= \prod_{i=1}^m p(y^{(i)}|x^{(i)}; \theta) \\ &= \prod_{i=1}^m \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(y^{(i)} - \theta^T x^{(i)})^2}{2\sigma^2}\right) \end{aligned} \quad (146)$$

To get the best choice of parameters θ , we perform maximum likelihood estimation such that $L(\theta)$ is maximized. Usually we take the negative log and minimize:

$$\begin{aligned} \ell(\theta) &= -\log L(\theta) \\ &= -\log \prod_{i=1}^m \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(y^{(i)} - \theta^T x^{(i)})^2}{2\sigma^2}\right) \\ &= -\sum_{i=1}^m \log \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(y^{(i)} - \theta^T x^{(i)})^2}{2\sigma^2}\right) \\ &= -m \log \frac{1}{\sqrt{2\pi}\sigma} + \frac{1}{\sigma^2} \cdot \frac{1}{2} \sum_{i=1}^m (y^{(i)} - \theta^T x^{(i)})^2 \end{aligned} \quad (147)$$

Hence, maximizing $L(\theta)$ is the same as minimizing the negative log likelihood $\ell(\theta)$, which for linear regression is the least squares cost function:

$$\frac{1}{2} \sum_{i=1}^m (y^{(i)} - \theta^T x^{(i)})^2 \quad (148)$$

Under the previous probabilistic assumptions on the data, least-squares regression corresponds to finding the maximum likelihood estimate of θ . This is thus one set of assumptions under which least-squares regression can be justified as performing maximum likelihood estimation. Note that θ is independent of σ^2 .

6.4 Locally Weighted Linear Regression

Locally Weighted Regression, also known as LWR, is a variant of linear regression that weights each training example in its cost function by $w^{(i)}(x)$, which is defined with parameter $\tau \in \mathbb{R}$ as:

$$w^{(i)}(x) = \exp\left(-\frac{(x^{(i)} - x)^2}{2\tau^2}\right) \quad (149)$$

Hence, in LWR, we do the following:

1. Fit θ to minimize $\sum_i w^{(i)} (y^{(i)} - \theta^T x^{(i)})^2$
2. Output $\theta^T x$

This is a non-parametric algorithm, where non-parametric refers to the fact that the amount of information we need to represent the hypothesis h grows linearly with the size of the training set.

7 Logistic Regression

We can extend this learning to classification problems, where we have binary labels y that are either 0 or 1.

7.1 The Logistic Function

For logistic regression, our new hypothesis for estimating the class of a sample x is:

$$h(x) = g(\theta^T x) = \frac{1}{1 + e^{-\theta^T x}} \quad (150)$$

where $g(z)$ is the logistic or sigmoid function:

$$g(z) = \frac{1}{1 + e^{-z}} \quad (151)$$

The sigmoid function is bounded between 0 and 1, and tends towards 1 as $z \rightarrow \infty$. It tends towards 0 when $z \rightarrow -\infty$. A useful property of the sigmoid function is the form of its derivative:

$$\begin{aligned} g'(z) &= \frac{d}{dz} \frac{1}{1 + e^{-z}} \\ &= \frac{1}{(1 + e^{-z})^2} (e^{-z}) \\ &= \frac{1}{(1 + e^{-z})} \cdot \left(1 - \frac{1}{(1 + e^{-z})}\right) \\ &= g(z)(1 - g(z)) \end{aligned} \quad (152)$$

7.2 Cost Function

To fit θ for a set of training examples, we assume that:

$$\begin{aligned} P(y = 1|x; \theta) &= h(x) \\ P(y = 0|x; \theta) &= 1 - h(x) \end{aligned} \quad (153)$$

This can be written more compactly as:

$$p(y|x; \theta) = (h(x))^y (1 - h(x))^{1-y} \quad (154)$$

Assume m training examples generated independently, we define the likelihood function of the parameters as:

$$\begin{aligned} L(\theta) &= p(\vec{y}|X; \theta) \\ &= \prod_{i=1}^m p(y^{(i)}|x^{(i)}; \theta) \\ &= \prod_{i=1}^m \left(h(x^{(i)}) \right)^{y^{(i)}} \left(1 - h(x^{(i)}) \right)^{1-y^{(i)}} \end{aligned} \quad (155)$$

And taking the negative log likelihood to minimize:

$$\begin{aligned}\ell(\theta) &= -\log L(\theta) \\ &= -\sum_{i=1}^m y^{(i)} \log h(x^{(i)}) + (1 - y^{(i)}) \log (1 - h(x^{(i)}))\end{aligned}\tag{156}$$

This is known as the binary cross-entropy loss function.

7.3 Gradient Descent

Lets start by working with just one training example (x,y), and take derivatives to derive the stochastic gradient ascent rule:

$$\begin{aligned}\frac{\partial}{\partial \theta_j} \ell(\theta) &= -\left(y \frac{1}{g(\theta^T x)} - (1 - y) \frac{1}{1 - g(\theta^T x)}\right) \frac{\partial}{\partial \theta_j} g(\theta^T x) \\ &= -\left(y \frac{1}{g(\theta^T x)} - (1 - y) \frac{1}{1 - g(\theta^T x)}\right) g(\theta^T x) (1 - g(\theta^T x)) \frac{\partial}{\partial \theta_j} \theta^T x \\ &= -(y(1 - g(\theta^T x)) - (1 - y)g(\theta^T x)) x_j \\ &= -(y - h(x)) x_j\end{aligned}\tag{157}$$

This therefore gives us the stochastic gradient ascent rule:

$$\theta_j := \theta_j + \alpha (y^{(i)} - h(x^{(i)})) x_j^{(i)}\tag{158}$$

We must use gradient descent for logistic regression since there is no closed form solution for this problem.

8 Softmax Regression

A softmax regression, also called a multiclass logistic regression, is used to generalize logistic regression when there are more than 2 outcome classes.

8.1 Softmax Function

The softmax function creates a probability distribution over a set of k classes for a training example x , with θ_k denoting the set of parameters to be optimized for the k -th class.

$$p(y = k|x; \theta) = \frac{\exp(\theta_k^T x)}{\sum_j \exp(\theta_j^T x)}\tag{159}$$

8.2 MLE and Cost Function

We can write the maximum likelihood function for softmax regression as:

$$L(\theta) = \prod_{i=1}^m \prod_k p(y = k|x; \theta)^{\mathbf{1}_{\{y_i=k\}}}\tag{160}$$

Where $\mathbf{1}\{y_i = k\}$ is the indicator function which is 1 if its argument is true, 0 otherwise. By taking the negative log likelihood:

$$\begin{aligned}
 \ell(\theta) &= -\log L(\theta) \\
 &= -\log \prod_{i=1}^m \prod_k p(y = k|x; \theta)^{\mathbf{1}\{y_i=k\}} \\
 &= -\sum_{i=1}^m \sum_k \left(\mathbf{1}\{y_i = k\} \cdot \left(\theta_k^T x_i - \log \left(\sum_j \exp(\theta_j^T x_i) \right) \right) \right) \\
 &= \sum_{i=1}^m -\theta_{y_i}^T x_i + \log \left(\sum_j \exp(\theta_j^T x_i) \right)
 \end{aligned} \tag{161}$$

This is known as the cross-entropy loss function.

8.3 Gradient Descent

To perform gradient descent, we must take the derivative of our cost function, but it is important to note that the derivative for the correct class is different than the other classes.

$$\begin{aligned}
 \nabla_{\theta_j} \ell(\theta) &= \nabla_{\theta_j} \left(\sum_{i=1}^m -\theta_{y_i}^T x_i + \log \left(\sum_k \exp(\theta_k^T x_i) \right) \right) \\
 &= \sum_{i=1}^m \nabla_{\theta_j} (-\theta_{y_i}^T x_i) + \nabla_{\theta_j} \left(\log \left(\sum_k \exp(\theta_k^T x_i) \right) \right) \\
 &= \sum_{i=1}^m \mathbf{1}\{y_i = j\} \cdot (-x_i) + \frac{\exp(\theta_j^T x_i)}{\sum_k \exp(\theta_k^T x_i)} \cdot x_i \\
 &= \sum_{i=1}^m \left(\frac{\exp(\theta_j^T x_i)}{\sum_k \exp(\theta_k^T x_i)} - \mathbf{1}\{y_i = j\} \right) \cdot x_i
 \end{aligned} \tag{162}$$

And our update equation for the j -th parameter weights is:

$$\theta_j := \theta_j - \alpha \left(\frac{\exp(\theta_j^T x_i)}{\sum_k \exp(\theta_k^T x_i)} - \mathbf{1}\{y_i = j\} \right) \cdot x_i \tag{163}$$

Note that since each class has a set of weights, our gradient is a matrix known as the jacobian \mathbf{J} , with k classes each with n feature weights.

$$\mathbf{J}_\theta = \begin{bmatrix} \frac{\partial \ell(\theta)}{\partial \theta_1} & \dots & \frac{\partial \ell(\theta)}{\partial \theta_k} \end{bmatrix} = \begin{bmatrix} \frac{\partial \ell(\theta)}{\partial \theta_{11}} & \dots & \frac{\partial \ell(\theta)}{\partial \theta_{k1}} \\ \vdots & \ddots & \vdots \\ \frac{\partial \ell(\theta)}{\partial \theta_{1n}} & \dots & \frac{\partial \ell(\theta)}{\partial \theta_{kn}} \end{bmatrix} \tag{164}$$

9 Generalized Linear Models

9.1 Exponential Family

9.2 Assumptions of GLMs

9.3 Examples

9.3.1 Ordinary Least Squares

9.3.2 Logistic Regression

9.3.3 Softmax Regression

10 Perceptron

11 Support Vector Machines

12 Margin Classification

13 Generative Learning: Gaussian Discriminant Analysis

13.1 Assumptions

13.2 Estimation

14 Generative Learning: Naive Bayes

14.1 Assumptions

14.2 Estimation

15 Tree-based Methods

16 K-Nearest Neighbors

The k -nearest neighbors algorithm, commonly known as k -NN, is a non-parametric approach where the response of a data point is determined by the nature of its k neighbors from the training set. It can be used in both classification and regression settings. The higher the parameter k , the higher the bias, and the lower the parameter k , the higher the variance.

16.1 Classification

In k -NN classification, the output is a class membership. An object is classified by a majority vote of its neighbors, with the object being assigned to the class most common among its k nearest neighbors (k is a positive integer, typically small). If $k = 1$, then the object is simply assigned to the class of that single nearest neighbor.

16.2 Regression

In k -NN regression, the output is the property value for the object. This value is the average of the values of its k nearest neighbors.

$$y = \frac{1}{k} \sum_{x_i \in \mathcal{N}_k(x)} y_i \quad (165)$$

where $\mathcal{N}_k(x)$ is the k nearest points around x .

17 K-Means Clustering

CLustering seeks to group similar points of data together in a cluster. We denote $c^{(i)}$ as the cluster for data point i and μ_j as the center for cluster j . We denote k as the number of clusters and n as the dimension of our data.

17.1 Algorithm

After randomly initializing the cluster centroids $\mu_1, \mu_2, \dots, \mu_k \in \mathbb{R}^n$, repeat until convergence:

1. For every data point i :

$$c^{(i)} = \arg \min_j \|x^{(i)} - \mu_j\|^2 \quad (166)$$

2. For each cluster j :

$$\mu_j = \frac{\sum_{i=1}^m 1_{\{c^{(i)}=j\}} x^{(i)}}{\sum_{i=1}^m 1_{\{c^{(i)}=j\}}} \quad (167)$$

The first step is known as cluster assignment, and the second updates the cluster center (i.e. the average of all points in the cluster). In order to see if it converges, use the distortion function:

$$J(c, \mu) = \sum_{i=1}^m \|x^{(i)} - \mu_{c^{(i)}}\|^2 \quad (168)$$

The distortion function J is non-convex, and coordinate descent of J is not guaranteed to converge to the global minimum (i.e. susceptible to local optima).

17.2 Hierarchical Clustering

Hierarchical clustering is a clustering algorithm with an agglomerative hierarchical approach that builds nested clusters in a successive manner. The types are:

1. Ward Linkage: minimize within cluster distance
2. Average Linkage: minimize average distance between cluster pairs
3. Complete Linkage: minimize maximum distance between cluster pairs

17.3 Clustering Metrics

In an unsupervised learning setting, it is often hard to assess the performance of a model since we don't have the ground truth labels as was the case in the supervised learning setting.

Silhouette coefficient By noting a and b the mean distance between a sample and all other points in the same class, and between a sample and all other points in the next nearest cluster, the silhouette coefficient s for a single sample is defined as follows:

$$s = \frac{b - a}{\max(a, b)} \quad (169)$$

Calinski-Harabaz Index By noting k the number of clusters, B_k and W_k the between and within-clustering dispersion matrices defined as:

$$B_k = \sum_{j=1}^k n_{c^{(j)}} (\mu_{c^{(j)}} - \mu)(\mu_{c^{(j)}} - \mu)^T \quad (170)$$

$$W_k = \sum_{i=1}^m (x^{(i)} - \mu_{c^{(i)}})(x^{(i)} - \mu_{c^{(i)}})^T \quad (171)$$

the Calinski-Harabaz index $s(k)$ indicated how well a clustering model defines its clusters, such that higher scores indicate more dense and well separated cluster assignments. It is defined as:

$$s(k) = \frac{\text{Tr}(B_k)}{\text{Tr}(W_k)} \times \frac{N - k}{k - 1} \quad (172)$$

18 Expectation-Maximization

18.1 Mixture of Gaussians

18.2 Factor Analysis

19 Principal Component Analysis

Principal Component Analysis is a dimension reduction technique that finds the variance maximizing the directions onto which to project the data.

19.1 Eigenvalues, Eigenvectors, and the Spectral Theorem

Recall that for a given matrix $A \in \mathbb{R}^{n \times n}$, λ is said to be an eigenvalue of A if there exists a vector $z \in \mathbb{R}^n \setminus \{0\}$, called an eigenvector, such that:

$$Az = \lambda z \quad (173)$$

The spectral theorem states that given matrix $A \in \mathbb{R}^{n \times n}$, if A is symmetric then A is diagonalizable by a real orthogonal matrix $U \in \mathbb{R}^{n \times n}$. Note $\Lambda = \text{diag}(\lambda_1, \dots, \lambda_n)$. Then $\exists \Lambda$ such that:

$$A = U\Lambda U^T \quad (174)$$

Note that the eigenvector associated with the largest eigenvalue is called the principal eigenvector of matrix A .

19.2 Algorithm

The PCA procedure projects the data onto k dimensions by maximizing the variance of the data as follows:

1. Normalize the data to have mean 0, standard deviation 1:

$$x^{(i)} \leftarrow \frac{x^{(i)} - \mu}{\sigma} \quad (175)$$

where μ and σ^2 are:

$$\mu = \frac{1}{m} \sum_{i=1}^m x^{(i)} \quad (176)$$

$$\sigma^2 = \frac{1}{m} \sum_{i=1}^m (x^{(i)} - \mu)^2 \quad (177)$$

2. Compute covariance matrix Σ , which is symmetric with real eigenvalues.

$$\Sigma = \frac{1}{m} \sum_{i=1}^m x^{(i)} x^{(i)T} \in \mathbb{R}^{n \times n} \quad (178)$$

3. Compute $u_1, \dots, u_k \in \mathbb{R}^n$ the k orthogonal principal eigenvectors of Σ , i.e. the orthogonal eigenvectors of the k largest eigenvalues.

4. Project the data on $\text{span}_{\mathbb{R}}(u_1, \dots, u_k)$ to create a vector $y^{(i)}$ from point $x^{(i)}$:

$$y^{(i)} = U^T x^{(i)} = \begin{pmatrix} u_1^T x^{(i)} \\ u_2^T x^{(i)} \\ \vdots \\ u_k^T x^{(i)} \end{pmatrix} \in \mathbb{R}^k \quad (179)$$

This procedure maximizes the variance among all k -dimensional spaces.

20 Independent Component Analysis

21 Reinforcement Learning

21.1 Markov Decision Processes

21.2 Policy and Value Functions

21.3 Value Iteration Algorithm

21.4 Q-Learning

22 Hidden Markov Models

23 Deep Learning: Basics

23.1 Basics

23.2 Activation Functions

23.3 Loss Functions

23.4 Backpropagation

23.5 Regularization Methods

23.6 Optimization Algorithms

23.7 Convolutional Networks

23.8 Recurrent Networks

23.8.1 Elman RNN

23.8.2 Long Short-Term Memory

23.8.3 Gated Recurrent Unit

23.8.4 Bidirectional RNNs

24 Deep Learning: Advanced

24.1 Autoencoders

24.1.1 Variational Autoencoders

24.2 General Adversarial Networks

24.3 Encoder-Decoder Models

24.3.1 Attention Models