

Machine Learning Study Guide

WILLIAM WATSON

Johns Hopkins University

billwatson@jhu.edu

1 Linear Algebra and Calculus

- 1.1 General Notation
- 1.2 Matrix Operations
- 1.3 Matrix Properties
- 1.4 Matrix Calculus

2 Convex Optimization

- 2.1 Convexity
- 2.2 Convex Optimization
 - 2.2.1 Gradient Descent
 - 2.2.2 Newton's Algorithm
- 2.3 Lagrange Duality and KKT Conditions

3 Probability and Statistics

- 3.1 Basics
- 3.2 Conditional Probability
- 3.3 Random Variables
- 3.4 Jointly Distributed Random Variables
- 3.5 Parameter Estimation
- 3.6 Probability Bounds and Inequalities

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:

$$I(x) = -\log P(x) \tag{1}$$

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

$$H(P) = \mathbb{E}_{x \sim P}[I(x)] = -\mathbb{E}_{x \sim P}[\log P(x)] \tag{2}$$

Which in the discrete setting is written as:

$$H(P) = - \sum_x P(x) \log P(x) \quad (3)$$

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:

$$\begin{aligned} D_{\text{KL}}(P\|Q) &= \mathbb{E}_{\mathbf{x} \sim P} \left[\log \frac{P(x)}{Q(x)} \right] \\ &= \mathbb{E}_{\mathbf{x} \sim P} [\log P(x) - \log Q(x)] \\ &= \sum_x P(x) \frac{\log P(x)}{\log Q(x)} \end{aligned} \quad (4)$$

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.

$$\begin{aligned} H(P, Q) &= H(P) + D_{\text{KL}}(P\|Q) \\ &= - \mathbb{E}_{\mathbf{x} \sim P} [\log Q(x)] \\ &= - \sum_x P(x) \log Q(x) \end{aligned} \quad (5)$$

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 Machine Learning Basics

5.1 Notation

5.2 Types of Learning

5.3 Metrics

5.3.1 Classification

5.3.2 Regression

5.4 Bias and Variance

6 Linear Regression

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

$$h_{\theta}(x) = \theta_0 + \theta_1 \cdot x_1 + \cdots + \theta_n \cdot x_n \quad (6)$$

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:

$$h(x) = \sum_{i=0}^n \theta_i \cdot x_i = \theta^T x \quad (7)$$

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:

$$J(\theta) = \frac{1}{2} \sum_{i=1}^m \left(h(x^{(i)}) - y^{(i)} \right)^2 \quad (8)$$

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:

$$\theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta) \quad (9)$$

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

$$\begin{aligned} \frac{\partial}{\partial \theta_j} J(\theta) &= \frac{\partial}{\partial \theta_j} \frac{1}{2} (h(x) - y)^2 \\ &= 2 \cdot \frac{1}{2} (h(x) - y) \cdot \frac{\partial}{\partial \theta_j} (h(x) - y) \\ &= (h(x) - y) \cdot \frac{\partial}{\partial \theta_j} \left(\sum_{i=0}^n \theta_i x_i - y \right) \\ &= (h(x) - y) x_j \end{aligned} \quad (10)$$

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

$$\theta_j := \theta_j + \alpha \left(y^{(i)} - h(x^{(i)}) \right) x_j^{(i)} \quad (11)$$

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).

$$\theta_j := \theta_j + \alpha \sum_{i=1}^m \left(y^{(i)} - h(x^{(i)}) \right) x_j^{(i)} \quad (12)$$

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:

$$J(\theta) = \frac{1}{2} (X\theta - \vec{y})^T (X\theta - \vec{y}) = \frac{1}{2} \sum_{i=1}^m \left(h(x^{(i)}) - y^{(i)} \right)^2 \quad (13)$$

We then take the derivative and find its roots.

$$\begin{aligned}
 \nabla_{\theta} J(\theta) &= \nabla_{\theta} \frac{1}{2} (X\theta - \vec{y})^T (X\theta - \vec{y}) \\
 &= \frac{1}{2} \nabla_{\theta} (\theta^T X^T X \theta - \theta^T X^T \vec{y} - \vec{y}^T X \theta + \vec{y}^T \vec{y}) \\
 &= \frac{1}{2} \nabla_{\theta} (\text{tr } \theta^T X^T X \theta - 2 \text{tr } \vec{y}^T X \theta) \\
 &= \frac{1}{2} (X^T X \theta + X^T X \theta - 2 X^T \vec{y}) \\
 &= X^T X \theta - X^T \vec{y}
 \end{aligned} \tag{14}$$

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

$$X^T X \theta = X^T \vec{y} \tag{15}$$

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

$$\theta = (X^T X)^{-1} X^T \vec{y} \tag{16}$$

6.3 Probabilistic Interpretation

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

$$y^{(i)} = \theta^T x^{(i)} + \epsilon^{(i)} \tag{17}$$

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:

$$p(y^{(i)}|x^{(i)}; \theta) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(y^{(i)} - \theta^T x^{(i)})^2}{2\sigma^2}\right) \tag{18}$$

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

$$L(\theta) = L(\theta; X, \vec{y}) = p(\vec{y}|X; \theta) \tag{19}$$

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

$$\begin{aligned}
 L(\theta) &= \prod_{i=1}^m p(y^{(i)}|x^{(i)}; \theta) \\
 &= \prod_{i=1}^m \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(y^{(i)} - \theta^T x^{(i)})^2}{2\sigma^2}\right)
 \end{aligned} \tag{20}$$

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:

$$\begin{aligned}
 \ell(\theta) &= -\log L(\theta) \\
 &= -\log \prod_{i=1}^m \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(y^{(i)} - \theta^T x^{(i)})^2}{2\sigma^2}\right) \\
 &= -\sum_{i=1}^m \log \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(y^{(i)} - \theta^T x^{(i)})^2}{2\sigma^2}\right) \\
 &= -m \log \frac{1}{\sqrt{2\pi}\sigma} + \frac{1}{\sigma^2} \cdot \frac{1}{2} \sum_{i=1}^m (y^{(i)} - \theta^T x^{(i)})^2
 \end{aligned} \tag{21}$$

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:

$$\frac{1}{2} \sum_{i=1}^m (y^{(i)} - \theta^T x^{(i)})^2 \tag{22}$$

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:

$$w^{(i)}(x) = \exp\left(-\frac{(x^{(i)} - x)^2}{2\tau^2}\right) \tag{23}$$

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:

$$h(x) = g(\theta^T x) = \frac{1}{1 + e^{-\theta^T x}} \tag{24}$$

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

$$g(z) = \frac{1}{1 + e^{-z}} \tag{25}$$

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:

$$\begin{aligned}
 g'(z) &= \frac{d}{dz} \frac{1}{1 + e^{-z}} \\
 &= \frac{1}{(1 + e^{-z})^2} (e^{-z}) \\
 &= \frac{1}{(1 + e^{-z})} \cdot \left(1 - \frac{1}{(1 + e^{-z})}\right) \\
 &= g(z)(1 - g(z))
 \end{aligned} \tag{26}$$

7.2 Cost Function

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

$$\begin{aligned}
 P(y = 1|x; \theta) &= h(x) \\
 P(y = 0|x; \theta) &= 1 - h(x)
 \end{aligned} \tag{27}$$

This can be written more compactly as:

$$p(y|x; \theta) = (h(x))^y (1 - h(x))^{1-y} \tag{28}$$

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

$$\begin{aligned}
 L(\theta) &= p_m^m | X; \theta \\
 &= \prod_{i=1}^m p(y^{(i)} | x^{(i)}; \theta) \\
 &= \prod_{i=1}^m \left(h(x^{(i)}) \right)^{y^{(i)}} \left(1 - h(x^{(i)}) \right)^{1-y^{(i)}}
 \end{aligned} \tag{29}$$

And taking the negative log likelihood to minimize:

$$\begin{aligned}
 \ell(\theta) &= -\log L(\theta) \\
 &= -\sum_{i=1}^m y^{(i)} \log h(x^{(i)}) + (1 - y^{(i)}) \log (1 - h(x^{(i)}))
 \end{aligned} \tag{30}$$

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:

$$\begin{aligned}
 \frac{\partial}{\partial \theta_j} \ell(\theta) &= -\left(y \frac{1}{g(\theta^T x)} - (1 - y) \frac{1}{1 - g(\theta^T x)} \right) \frac{\partial}{\partial \theta_j} g(\theta^T x) \\
 &= -\left(y \frac{1}{g(\theta^T x)} - (1 - y) \frac{1}{1 - g(\theta^T x)} \right) g(\theta^T x) (1 - g(\theta^T x)) \frac{\partial}{\partial \theta_j} \theta^T x \\
 &= -(y(1 - g(\theta^T x)) - (1 - y)g(\theta^T x)) x_j \\
 &= -(y - h(x)) x_j
 \end{aligned} \tag{31}$$

This therefore gives us the stochastic gradient ascent rule:

$$\theta_j := \theta_j + \alpha (y^{(i)} - h(x^{(i)})) x_j^{(i)} \tag{32}$$

8 Softmax Regression

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 Generative Learning: Gaussian Discriminant Analysis

13 Generative Learning: Naive Bayes

14 Tree-based Methods

15 K-Nearest Neighbors

16 K-Means Clustering

16.1 Hierarchical Clustering

16.2 Clustering Metrics

17 Expectation-Maximization

17.1 Mixture of Gaussians

17.2 Factor Analysis

18 Principal Component Analysis

19 Independent Component Analysis

20 Reinforcement Learning

20.1 Markov Decision Processes

20.2 Policy and Value Functions

20.3 Value Iteration Algorithm

20.4 Q-Learning

21 Hidden Markov Models

22 Deep Learning

22.1 Basics