

GENERAL NOTES

ANDREW S. HARD¹

January 12, 2015

CONTENTS

1	Introduction	3
2	Basic Statistics	3
2.1	Deduction and Induction	3
2.2	Conditional Probability and Bayes' Theorem	3
2.3	Correlation and Covariance	5
3	Information Theory	5
3.1	Entropy	6
4	Data Structures	7
4.1	List Structure	7
4.2	Tree Structure	7
4.3	Stack and Queue Structures	7
4.4	Graph Structure	7
4.5	Hash Maps	7
5	Numerical Methods	7
5.1	Errors	7
5.2	Root-Finding	8
5.3	Minimization	9
5.4	Regression	9
5.5	Matrix Algebra	9
6	Machine Learning	9
6.1	Introduction	9
6.2	Feature Selection	9
6.3	Regression (Again!)	9
6.4	Clustering	9
6.5	Support Vector Machines	9
6.6	Decision Trees	9
6.7	Lasso Method	9
7	Deep Learning	9
7.1	Overview	9
7.2	Perceptrons	9
7.3	MLP: Multi-Layer Perceptrons	9
7.4	RBM: Restricted Boltzmann Machines	9
7.5	CNN: Convolutional Neural Networks	11
7.6	RNN: Recurrent Neural Networks	11
7.7	LSTM: Long Short Term Memory	11
7.8	Auto-Encoders and Deep Belief Networks	11
7.9	Neural Turing Machines	11
8	Lasso Method	11

LIST OF FIGURES

Figure 1	A Boltzmann Machine diagram	10
----------	---------------------------------------	----

LIST OF TABLES

¹ *Department of Physics, University of Wisconsin, Madison, United States of America*

1 INTRODUCTION

As a fifth year graduate student, I have come to the realization that significant portions of my skillset and knowledge base are highly specialized. I intend to pursue a career outside of my field of study (experimental high-energy particle physics). In order to enhance my future job prospects, I decided that it would be useful to review basic concepts in computer science, statistics, mathematics, and machine learning. Most of the derivations and explanations will be borrowed from other sources. A list of references is currently being compiled.

2 BASIC STATISTICS

2.1 Deduction and Induction

Deduction: if $A \rightarrow B$ and A is true, then B is true.

Induction: if $A \rightarrow B$ and B is true, then A is more plausible.

2.2 Conditional Probability and Bayes' Theorem

The **conditional probability** of an event A assuming that B has occurred, denoted $P(A|B)$, is given by:

$$P(A|B) = \frac{P(A \cap B)}{P(B)}. \quad (1)$$

Via multiplication:

$$P(A \cap B) = P(A|B)P(B). \quad (2)$$

This makes sense. The probability of A and B is the probability of A *given* B times the probability of B . Or, just as reasonably, the probability of A and B is the probability of B *given* A times the probability of A . As an aside, the equation above can be generalized:

$$P(A \cap B \cap C) = P(A|B \cap C)P(B \cap C) = P(A|B \cap C)P(B|C)P(C). \quad (3)$$

Back to the derivation, since $P(A \cap B) = P(B \cap A)$, we can use Equation 2:

$$P(A|B)P(B) = P(B|A)P(A). \quad (4)$$

And re-arranging gives **Bayes' Theorem**:

$$P(A|B) = \frac{P(A)P(B|A)}{P(B)} = \frac{P(A)}{P(B)}P(B|A). \quad (5)$$

In this equation, $P(A)$ is the **prior probability**, the initial degree of belief in A , or the probability of A before B is observed. It is essentially the probability that hypothesis A is true before evidence is collected. The **posterior probability** $P(A|B)$ is the degree of belief after taking B into account.

In a nutshell, Bayes' theorem is a method for calculating the validity of beliefs based on available evidence. "Initial belief plus new evidence equals new and improved belief."

Bayes' theorem describes the probability of an event, based on conditions that might be relevant to the event. With the Bayesian probability interpretation the theorem expresses how a subjective degree of belief should rationally change to account for evidence (**Bayesian inference**). Note: in the **Frequentist Interpretation**, probability measures a proportion of outcomes, not probability of belief.

Bayes' theorem can also be interpreted in the following manner: A is some hypothesis, and B is the evidence for that hypothesis.

The equation makes intuitive sense, particularly in the last form in which it is presented. If $P(B) \gg P(A)$, the probability $P(A|B)$ will be small because B *usually* occurs without A occurring. Even if $P(B|A)$ were 100%, the posterior probability would still be small. On the other hand, if $P(A) \gg P(B)$, $P(A|B)$ is high, even if the conditional probability $P(B|A)$ is smaller.

Nate Silver has a discussion of Bayesian statistics in "The Signal And The Noise", p245. Let x be the initial estimate of the hypothesis likelihood (**prior probability**). In the book, x is the initial estimate that your spouse is cheating. y is a conditional probability - the probability of an observation being made given that the hypothesis is true. So y in his case was the probability of underwear appearing conditional on the spouse cheating. z is the probability of an observation being made given that the hypothesis is false. So z was the probability of underwear appearing if the spouse was not cheating. The **posterior probability** is the revised estimate of the hypothesis likelihood. In Silver's case, the likelihood that the spouse is cheating on you, given that underwear was found. The posterior probability is given by the expression:

$$x' = \frac{xy}{xy + z(1 - x)}. \quad (6)$$

In Silver's formulation, the denominator is equivalent to the total probability of an observation being made (total probability of underwear being found). So it is the probability of the hypothesis being true times the associated conditional probability of the observation of underwear plus the probability of the hypothesis being false times the associated conditional probability of the observation.

Note that this method can be applied recursively by substituting x' for x in the expression. We can also come up with a recurrent form of Bayes theorem as formulated in Equation 5 by letting $P(A_{n+1}) = P(A_n|B)$.

Yet another formulation, which is explained by the logic in the previous paragraph, is given by:

$$P(A|X) = \frac{P(X|A)P(A)}{P(X|A)P(A) + P(X|\text{not}A)P(\text{not}A)} \quad (7)$$

$P(A|X)$ is the probability of hypothesis A, given a positive observation X. $P(X|A)$ is the probability of the observation of X under hypothesis A. $P(A)$ is the probability of the hypothesis.

The derivation below is from mathworld's article on Bayes' Theorem. In this example, let A and S be sets. Furthermore, let

$$S \equiv \bigcup_{i=1}^N A_i, \quad (8)$$

so that A_i is an event in S , and $A_i \cap A_j = \emptyset \forall i \neq j$. Then:

$$A = A \cap S = A \cap \left(\bigcup_{i=1}^N A_i \right) = \bigcup_{i=1}^N (A \cap A_i) \quad (9)$$

Then we can compare the probabilities using the Law of Total Probability:

$$P(A) = P\left(\bigcup_{i=1}^N (A \cap A_i)\right) = \sum_{i=1}^N P(A \cap A_i) = \sum_{i=1}^N P(A_i)P(A|A_i) \quad (10)$$

Using substitution into Bayes' Theorem (Equation 5), we have a new form:

$$P(A_j|A) = \frac{P(A_j)P(A|A_j)}{\sum_{i=1}^N P(A_i)P(A|A_i)}. \quad (11)$$

What is the meaning of this? The denominator is simply $P(A)$, expressed as the sum of its constituents. So this reduces exactly to Bayes theorem as in Equation 5, with $A \rightarrow A_j$ and $B \rightarrow A$.

"In many cases, estimating the prior is just guesswork, allowing subjective factors to creep into your calculations. You might be guessing the probability of something that, unlike cancer, does not even exist, such as strings, multiverses, inflation or God. You might then cite dubious evidence to support your dubious belief. In this way, Bayes' theorem can promote pseudoscience and superstition as well as reason." From <http://blogs.scientificamerican.com/cross-check/bayes-s-theorem-what-s-the-big-deal/>

2.3 Correlation and Covariance

Correlation and covariance are similar. Both concepts describe the degree to which two random variables or sets of random variables tend to deviate from their expected values in similar ways.

Let X and Y be two random variables, with means μ_X and μ_Y and standard deviations σ_X and σ_Y , respectively. The **covariance** is defined (using $\langle \rangle$ to denote expectation value):

$$\sigma_{XY} = \langle (X - \langle X \rangle)(Y - \langle Y \rangle) \rangle. \quad (12)$$

Similarly, the **correlation** of the two variables is defined:

$$\rho_{XY} = \frac{\langle (X - \langle X \rangle)(Y - \langle Y \rangle) \rangle}{\sigma_X \sigma_Y}. \quad (13)$$

So correlation is dimensionless, while covariance is the multiple of the units of the two variables. **Variance** is simply the covariance of a variable with itself (σ_{XX} is usually denoted σ_X^2 , the square of the standard deviation). The correlation of a variable with itself is always 1, except in the degenerate case where the two variances are zero and the correlation does not exist.

3 INFORMATION THEORY

Information theory deals with quantifying information. It was developed by Claude Shannon with the goal of establishing limits on signal processing operations such as data compression, reliable storage, and communication (cite wiki).

The basic unit of information is the **bit**, a portmanteau of binary digit. It can have values of 0 or 1, and could represent boolean states (false/true) or activation states (off/on) among other things. **Bit-length** refers to the length of a binary number. Table 3 lists common data types from the IEEE Standard For Floating Point Arithmetic (IEEE 754).

data type	bit length	bit assignment
signed int	32	Bit string with range $[-2^{31}, 2^{31} - 1]$
unsigned int	32	Bit string with range $[0, 2^{32} - 1]$
float	32	1 sign, 8 exponent, 24 significand precision
double	64	1 sign, 11 exponent, 53 significand precision

3.1 Entropy

Shannon entropy H - the average number of bits per symbol needed to store or communicate a message. "Bits assumes \log_2 . "It is analogous to intensive physical entropies like molar and specific entropy S_0 but not extensive physical entropy S ."

There's also entropy H times the number of symbols in a message. This measure of entropy is in bits. It is analogous to absolute (extensive) physical entropy S .

Landauer's Principle: A change in $H \cdot N$ is equal to a change in physical entropy when the information system is perfectly efficient.

What does entropy mean? It is a measure of the uncertainty involved in predicting the value of a random variable. Example: a coin flip with 50 % probabilities of heads or tails provides less information than a die roll with a $\frac{1}{6}$ probability per side.

The Shannon entropy H in units of bits per symbol is given by

$$H = - \sum_i f_i \log_2(f_i) = - \sum_i \log_2\left(\frac{c_i}{N}\right). \quad (14)$$

In this equation, i is an index over distinct symbols, f_i is the frequency of symbols occurring in the message, and c_i is the number of times the i^{th} symbol occurs in the message of length N . HN is the entropy of a message. This has units of bits (not bits per symbol).

Example: a Bernoulli (binary) trial where $X = 1$ occurs with some probability. The entropy $H(X)$ of the trial is maximized when the two possible outcomes are equally probable ($P(X = 1) = 0.5$).

Boltzmann H	Shannon H
physical systems of particles	information systems of symbols
entropy per particle	entropy per symbol
probability of a microstate	probability of a symbol
Extensive physical entropy S	Entropy of file or message HN

4 DATA STRUCTURES

4.1 List Structure

4.1.1 *Basic List*

4.1.2 *Linked Lists*

4.1.3 *Circular Lists*

4.2 Tree Structure

4.2.1 *Binary trees*

4.2.2 *Tree Balancing*

4.2.3 *Red-Black trees*

4.3 Stack and Queue Structures

4.3.1 *Stacks*

4.3.2 *Queue*

4.3.3 *Priority Queue*

4.3.4 *Max/Min Heap*

4.4 Graph Structure

Define the **depth-first search**, **breadth-first search**, etc.

Dijkstra's algorithm is an algorithm for finding the shortest paths between nodes in a graph [?]. The original form of Dijkstra's algorithm has time complexity $O(|V|^2)$, where $|V|$ is the number of nodes. However, there is a priority-queue implementation of the algorithm that has worst-case performance of $O(|E| + |V| \log |V|)$, where $|E|$ is the number of edges.

Dijkstra's algorithm is known as a uniform-cost search, and is a type of best-first search.

4.5 Hash Maps

5 NUMERICAL METHODS

This section describes a number of methods for solving common problems in scientific programming, such as root-finding, minimization, regression, and matrix algebra. The discussions of each algorithm contain information on the algorithm complexity, the implementation, and the relative merits of the approach.

5.1 Errors

Many of the methods discussed provide numerical solutions rather than analytic solutions. This means that results will be approximations, and as such will have an associated error. **Approximation error** or **absolute error** is the discrepancy between an exact value and an approximation to that exact value. Such errors can occur because approximations are used instead of real data, or because the measurement of data is imprecise.

Absolute error is given by:

$$\epsilon = |v - v_{\text{approx.}}| \quad (15)$$

Relative error is defined:

$$\eta = \frac{\epsilon}{|v|} = |1 - \frac{v_{\text{approx.}}}{v}| \quad (16)$$

There are two primary error types for numerical methods: **roundoff errors** due to computer approximations, and **truncation errors** due to mathematical approximations. It is useful to avoid subtraction in numerical methods because subtraction leads to an increase in the relative significance of roundoff errors. Subtracting two very similar, very large numbers will lead to a very small difference. This is called **catastrophic cancellation**, and is demonstrated in the subtraction below.

$$0.12345678912345678 - 0.12345678900000000 = 0.00000000012345678. \quad (17)$$

But on a 10-digit floating-point machine, the calculation gives:

$$0.1234567891 - 0.1234567890 = 0.0000000001. \quad (18)$$

Two concepts related to error are accuracy and precision. **Accuracy** is how closely a computed or measured value agrees with the truth value. **Precision** is how closely individually computed or measured values agree with each other. For example, a darts player with tight clusters is a very precise thrower. But if the clusters are not centered around the target, they are inaccurate. In numerical methods, the term error can be used to represent both the inaccuracy and imprecision of calculations.

5.2 Root-Finding

Root-finding is a search for all values of \vec{x} such that $f(\vec{x}) = 0$. **Optimization** problems can be expressed as root-finding problems. To minimize or maximize a function $g(\vec{x})$, perform a root-finding search on $g'(\vec{x})$ (search for the values of \vec{x} such that $g'(\vec{x}) = 0$).

The **bisection method** is a simple root-finding method that repeatedly bisects an interval and selects the subinterval containing the root for further bisection. Advantages: the algorithm is simple to code recursively, and it is robust (always converges). Disadvantages: it is a relatively slow algorithm.

Given a function f that is continuous on the interval $[a, b]$ with $f(a)$ and $f(b)$ having opposite signs, the bisection method is guaranteed to converge. The absolute error is reduced by half during each application of the method. As a result, the bisection method *converges linearly*. It is significantly slower than other available algorithms. An upper bound on the absolute error (difference between the true and numerical result) can be determined by assuming the worst case scenario, that the root is at the midpoint of the interval $c = (a + b)/2$:

$$|c_n - c| \leq \frac{|b - a|}{2^n} \quad (19)$$

where n is the number of iterations. The number of iterations needed to meet a particular error tolerance ϵ is given by

$$\eta = \log_2\left(\frac{\epsilon_0}{\epsilon}\right) = \frac{\log \epsilon_0 - \log \epsilon}{\log 2} \quad (20)$$

where ϵ_0 is the size of the initial interval. The convergence of this algorithm is described as *linear* because $\epsilon_{n+1} = \text{constant} \times \epsilon_n$.

Next on the list: Newton-Raphson!

5.3 Minimization

5.4 Regression

Discuss linear regression

Nonlinear regression

5.5 Matrix Algebra

6 MACHINE LEARNING

6.1 Introduction

This section should contain a general description of data science principles, as well as use cases for different algorithms. Many tools for many problems.

Basic problems in machine learning are classification (labeling) and regression (function estimation).

How to do importance ranking of features? Feature importance, global loss function.

6.2 Feature Selection

Garbage in gives garbage out. Interactions between features. Interactions are not the same as correlations. Filters versus wrappers.

6.3 Regression (Again!)

6.4 Clustering

6.5 Support Vector Machines

6.6 Decision Trees

6.7 Lasso Method

7 DEEP LEARNING

7.1 Overview

Explain the use cases of each type of algorithm.

7.2 Perceptrons

7.3 MLP: Multi-Layer Perceptrons

7.4 RBM: Restricted Boltzmann Machines

<https://www.cs.toronto.edu/~hinton/absps/guideTR.pdf>

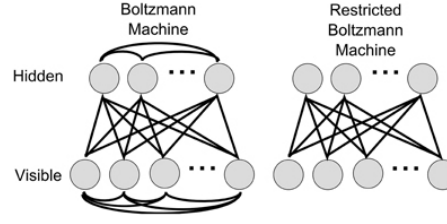


Figure 1: An example of a generic Boltzmann Machine versus a Restricted Boltzmann Machine, from http://www.frontiersin.org/files/Articles/58710/fnins-07-00178-r2/image_m/fnins-07-00178-g001.jpg.

"**Restricted Boltzmann machines (RBMs)** have been used as generative models of many different types of data including labeled or unlabeled images [?], windows of mel-cepstral coefficients that represent speech (Mohamed et al., 2009), bags of words that represent documents (Salakhutdinov and Hinton, 2009), and user ratings of movies (Salakhutdinov et al., 2007). In their conditional form they can be used to model high-dimensional temporal sequences such as video or motion capture data (Taylor et al., 2006) or speech (Mohamed and Hinton, 2010). Their most important use is as learning modules that are composed to form deep belief nets [?]."

A general Boltzmann machine is a network with symmetric weights fully connecting all input and hidden nodes. The details of the weights and nodes are described later in this section. A restricted Boltzmann machine separates nodes into layers and has no intra-layer connections. General Boltzmann machines, while theoretically providing a more diverse architecture, are difficult to train in practice. The difference between a general Boltzmann Machine and the Restricted Boltzmann Machine is illustrated in Figure 1.

The training procedure for RBMs is more complicated than back-propagation. They typically use a procedure called contrastive divergence [?]. There are several meta-parameters including learning rate, momentum, weight-cost, sparsity target, initial weights, number of hidden units, and mini-batch sizes.

For simplicity, the outline below will consider a two-layer network consisting of the visible unit layer representing the input data and a layer of hidden units called **feature detectors**. Deep neural networks can be trained by successively stacking RBMs, with the hidden layer of the first RBM used as an input for the second. Greedy pairwise training of layers has been demonstrated to be effective in learning deep architectures (reference greedy pair training). Input data are stochastic binary pixels (or continuous variables scaled to $[0, 1]$). Feature detectors are also stochastic binary units. The visible units are connected to the feature detectors using symmetrically weighted connections. A configuration of visible and hidden units has an energy (Hopfield, 1982) given by:

$$E(\mathbf{v}, \mathbf{h}) = - \sum_i a_i v_i - \sum_j b_j h_j - \sum_i \sum_j v_i h_j w_{ij} \quad (21)$$

for $i \in \text{visible}$ and $j \in \text{hidden}$. In this equation, v_i and h_j are the binary states of the visible unit i and hidden unit j , respectively. a_i and b_j are their biases and w_{ij} is the weight between them.

The energy can be used to define a probability of this pairing of visible and hidden units (hence the name "Boltzmann Machine"):

$$p(\mathbf{v}, \mathbf{h}) = \frac{1}{Z} e^{-E(\mathbf{v}, \mathbf{h})} \quad (22)$$

Z is the partition function:

$$Z = \sum_{\mathbf{v}, \mathbf{h}} e^{-E(\mathbf{v}, \mathbf{h})} \quad (23)$$

The network assigns a probability to a visible vector \mathbf{v} , which can be calculated by summing $p(\mathbf{v}, \mathbf{h})$ over all possible hidden vectors:

$$p(\mathbf{v}) = \frac{1}{Z} \sum_{\mathbf{h}} e^{-E(\mathbf{v}, \mathbf{h})} \quad (24)$$

A higher energy input is a less-probable input, according to the energy function. The network is trained to increase the probability (lowering the energy) assigned by the network to signal inputs while lowering the probability (increasing the energy) assigned to noise. To get the change in weights:

$$\frac{\delta \log p(\mathbf{v})}{\delta w_{ij}} = \langle v_i h_j \rangle_{\text{data}} - \langle v_i h_j \rangle_{\text{model}} \quad (25)$$

The angle brackets denote expectation values under the specified distribution (data or model). A learning rule for weights using stochastic gradient ascent on the log probability, and a learning rate ϵ , is then simply:

$$\delta w_{ij} = \epsilon (\langle v_i h_j \rangle_{\text{data}} - \langle v_i h_j \rangle_{\text{model}}). \quad (26)$$

Now all that remains is to get an unbiased sample of $\langle v_i h_j \rangle_{\text{data}}$ and $\langle v_i h_j \rangle_{\text{model}}$. The first one is straightforward, since there are no connections between hidden units of the same layer in an RBM. Given a randomly selected training image, \mathbf{v} , the binary state, h_j , of each hidden unit, j , is set to 1 with probability

$$p(h_j = 1 | \mathbf{v}) = \sigma(b_j + \sum_i v_i w_{ij}) \quad (27)$$

7.5 CNN: Convolutional Neural Networks

7.6 RNN: Recurrent Neural Networks

7.7 LSTM: Long Short Term Memory

7.8 Auto-Encoders and Deep Belief Networks

7.9 Neural Turing Machines

8 LASSO METHOD