

IV. Statistical Estimation and Classification

Probability: An Extremely Concise Review

1. A scalar-valued random variable X is completely characterized by its distribution function

$$F_X(u) = \mathbb{P}(X \leq u).$$

This is also called the **cumulative distribution function** (cdf). $F_X(u)$ is monotonically increasing in u ; it goes to one as $u \rightarrow \infty$ and goes to zero as $u \rightarrow -\infty$.

2. If F_X is differentiable, then we can also characterize X using its **probability density function** (pdf)

$$f_X(x) = \left. \frac{dF_X(u)}{du} \right|_{u=x}$$

The density has the properties $f_X(x) \geq 0$ and

$$\int_{-\infty}^{\infty} f_X(x) dx = 1.$$

Events of interest are subsets¹ of the real line — given such an event/subset \mathcal{E} , we can compute the probability of \mathcal{E} occurring as

$$\mathbb{P}(\mathcal{E}) = \int_{x \in \mathcal{E}} f_X(x) dx.$$

¹Technically, it must be a subset of the real line that can be written as some combination of countable unions, countable intersections, and complements of intervals. You really have to know something about real analysis to construct a set that does not meet this criteria.

It is possible that a pdf exists even if F_X is not differentiable everywhere, for example:

$$F_X(u) = \begin{cases} 0, & u < 0, \\ u, & 0 \leq u \leq 1, \\ 1, & u \geq 1 \end{cases} \quad \text{has pdf} \quad f_X(x) = \begin{cases} 0, & x < 0, \\ 1, & 0 \leq x \leq 1, \\ 0, & x > 1. \end{cases}$$

3. The **expectation** of a function $g(X)$ of a random variable is

$$E[g(X)] = \int_{-\infty}^{\infty} g(x) f_X(x) dx.$$

This is the “average value” of $g(X)$ in that given a series of realizations $X = x_1, X = x_2, \dots$, of X ,

$$\frac{1}{M} \sum_{m=1}^M g(x_m) \rightarrow E[g(X)], \quad \text{as } M \rightarrow \infty.$$

This fact is known as the **(weak) law of large numbers**.

4. The **moment** of X of degree p is the expectation of the monomial $g(x) = x^p$. The zeroth moment is always 1:

$$E[X^0] = E[1] = \int_{-\infty}^{\infty} f_X(x) dx = 1,$$

and the first moment is the **mean**:

$$E[X] = \int_{-\infty}^{\infty} x f_X(x) dx.$$

The **variance** is the second moment minus the mean squared:

$$\text{var}(X) = E[X^2] - (E[X])^2 = E[(X - E[X])^2].$$

This is sometime referred to as the “variation around the mean”. Aside from the zeroth moment, there is nothing that says that the integrals above must converge; it is easy to construct examples of well-defined random variables where $E[X] = \infty$.

5. A pair of random variables (X, Y) are completely described by their **joint distribution function** (joint cdf)²

$$F_{X,Y}(u, v) = P(X \leq u, Y \leq v).$$

Again, if $F_{X,Y}$ is continuously differentiable, (X, Y) is also characterized by the density

$$f_{X,Y}(x, y) = \left. \frac{\partial F_{X,Y}(u, v)}{\partial u \partial v} \right|_{(u,v)=(x,y)}.$$

In this case, events of interest correspond to regions in the plane \mathbb{R}^2 , and the probability of an event occurring is the integral of the density over this region.

6. From the joint pdf $f_{X,Y}(x, y)$, we can recover the individual **marginal pdfs** for X and Y using

$$\begin{aligned} f_X(x) &= \int_{-\infty}^{\infty} f_{X,Y}(x, y) dy, \\ f_Y(y) &= \int_{-\infty}^{\infty} f_{X,Y}(x, y) dx. \end{aligned}$$

The pair of densities $f_X(x)$, $f_Y(y)$ tell us how X and Y behave individually, but not how they *interact*.

²For fixed $u, v \in \mathbb{R}$, the notation $P(X \leq u, Y \leq v)$ should be read as “the probability that X is $\leq u$ and Y is $\leq v$.”

7. If X and Y do interact in a meaningful way, then observing one of them affects the distribution of the other. If we observe $X = x$, then with this knowledge, the density for Y becomes

$$f_Y(y|X = x) = \frac{f_{X,Y}(x, y)}{f_X(x)}.$$

This is a density over y ; it is easy to check that it is positive everywhere and that it integrates to one. $f_Y(y|X = x)$ is called the **conditional density** for Y given $X = x$.

8. We call X and Y **independent** if observing X tells us nothing about Y (and vice versa). This means

$$f_Y(y|X = x) = f_Y(y), \quad \text{for all } x \in \mathbb{R},$$

and

$$f_X(x|Y = y) = f_X(x), \quad \text{for all } y \in \mathbb{R}.$$

(If one of the statements above is true, then the other follows automatically.) Equivalently, independence means that the joint pdf is *separable*:

$$f_{X,Y}(x, y) = f_X(x) f_Y(y).$$

9. We can always factor the joint pdf in two different ways:

$$f_X(x)f_Y(y|X = x) = f_{X,Y}(x, y) = f_Y(y)f_X(x|Y = y).$$

At this point, we should be comfortable enough with what is going on that we can use $f_Y(y|x)$ as short-hand notation for $f_Y(y|X = x)$. Then we can rewrite the above in its more common form as

$$f_X(x)f_Y(y|x) = f_{X,Y}(x, y) = f_Y(y)f_X(x|y).$$

This factorization also gives us a handy way to compute the marginals:

$$f_X(x) = \int_{-\infty}^{\infty} f_Y(y) f_X(x|y) dy.$$

It also yields **Bayes' equation**

$$f_X(x|y) = \frac{f_Y(y|x) f_X(x)}{f_Y(y)},$$

which is a fundamental relation for statistical inference.

10. All of the above extends in the obvious way to more than two random variables. A **random vector**

$$X = \begin{bmatrix} X_1 \\ X_2 \\ \vdots \\ X_D \end{bmatrix}$$

is completely characterized by the density $f_X(\mathbf{x}) = f_X(x_1, \dots, x_D)$ on \mathbb{R}^D . In general, we can factor the joint pdf as

$$f_X(\mathbf{x}) = f_{X_1}(x_1) f_{X_2}(x_2|x_1) f_{X_3}(x_3|x_2, x_1) \cdots f_{X_D}(x_D|x_1, \dots, x_{D-1}).$$

11. The p th moment of a random vector X that maps into \mathbb{R}^D is the collection of expectations of all monomials of order p . The mean of a random vector is a vector of length D :

$$\mathbb{E}[X] = \begin{bmatrix} \mathbb{E}[X_1] \\ \vdots \\ \mathbb{E}[X_D] \end{bmatrix},$$

the second moment is the $D \times D$ matrix of all correlations between entries:

$$\mathbb{E}[XX^T] = \begin{bmatrix} \mathbb{E}[X_1^2] & \mathbb{E}[X_1X_2] & \cdots & \mathbb{E}[X_1X_D] \\ \vdots & & \ddots & \vdots \\ \mathbb{E}[X_DX_1] & \cdots & & \mathbb{E}[X_D^2] \end{bmatrix},$$

the third moment is the $D \times D \times D$ tensor $\mathbb{E}[X \otimes X \otimes X]$, where

$$(\mathbb{E}[X \otimes X \otimes X])(i, j, k) = \mathbb{E}[X_iX_jX_k],$$

and so on. The **covariance** matrix contains all the pairs of second moments:

$$R_{i,j} = \mathbb{E}[(X_i - \mathbb{E}[X_i])(X_j - \mathbb{E}[X_j])].$$

If $\boldsymbol{\mu}_X = \mathbb{E}[X]$ is the mean vector, we can write the covariance matrix succinctly in terms of the second moment as

$$\mathbf{R} = \mathbb{E}[XX^T] - \boldsymbol{\mu}_X\boldsymbol{\mu}_X^T$$

12. Given independent observations $X = \mathbf{x}_1, X = \mathbf{x}_2, \dots, X = \mathbf{x}_M$ of a random vector X with unknown (or partially known) distribution, a completely reasonable way to estimate the mean vector is using

$$\hat{\boldsymbol{\mu}} = \frac{1}{M} \sum_{m=1}^M \mathbf{x}_m.$$

If the mean $\boldsymbol{\mu}_X = \mathbb{E}[X]$ is known but the covariance is not, we can estimate the covariance using

$$\hat{\mathbf{R}} = \left(\frac{1}{M} \sum_{m=1}^M \mathbf{x}_m\mathbf{x}_m^T \right) - \boldsymbol{\mu}_X\boldsymbol{\mu}_X^T.$$

If both the mean and covariance are unknown, we first estimate the mean vector as above, then take

$$\hat{\mathbf{R}} = \left(\frac{1}{M-1} \sum_{m=1}^M \mathbf{x}_m \mathbf{x}_m^T \right) - \hat{\boldsymbol{\mu}} \hat{\boldsymbol{\mu}}^T.$$

The difference in the scaling is to ensure that $E[\hat{\mathbf{R}}] = \mathbf{R}$ in both cases.

The Weak Law of Large Numbers

The WLLN is absolutely fundamental to machine learning (and really to all of probability and statistics). It basically formalizes the notion that given a series of independent samples of a random variable X , we can approximate $E[X]$ by averaging the samples. The WLLN states that if X_1, X_2, \dots are independent copies of a random variable X ,

$$\frac{1}{N} \sum_{n=1}^N X_n \rightarrow E[X] \quad \text{as } N \rightarrow \infty.$$

The only condition for this convergence is that X has finite variance.

We start by stating the main result precisely. Let X be a random variable with pdf $f_X(x)$, mean $E[X] = \mu$, and variance $\text{var}(X) = \sigma^2 < \infty$. We observe *samples* of X labeled X_1, X_2, \dots, X_N . The X_i are independent of one another, and they all have the same distribution as X . We will show that the sample mean formed from a sample of size N :

$$M_N = \frac{1}{N}(X_1 + X_2 + \dots + X_N),$$

obeys³

$$P(|M_N - \mu| > \epsilon) \leq \frac{\sigma^2}{N\epsilon^2},$$

where $\epsilon > 0$ is an arbitrarily small number. In the expression above, M_N is the only thing which is random; μ and σ^2 are fixed underlying properties of the distribution, N is the amount of data we see, and ϵ is something we can choose arbitrarily.

³This is a simple example of a *concentration bound*. It is not that tight; we will later counter inequalities of this type that are much more precise. But it is relatively simple and will serve our purpose here.

Notice that no matter how small ϵ is, the probability on the right hand side above goes to zero as $N \rightarrow \infty$. That is, for any fixed $\epsilon > 0$,

$$\lim_{N \rightarrow \infty} P(|M_N - \mu| > \epsilon) = 0.$$

This result follows from two simple but important tools known as the *Markov* and *Chebyshev* inequalities.

Markov inequality

Let X be a random variable that only takes positive values:

$$f_X(x) = 0, \quad \text{for } x < 0, \quad \text{or} \quad F_X(0) = 0.$$

Then

$$P(X \geq a) \leq \frac{E[X]}{a} \quad \text{for all } a > 0.$$

For example, the probability that X is more than 5 times its mean is $1/5$, 10 times the mean is $1/10$, etc. And this holds for **any distribution**.

The Markov inequality is easy to prove:

$$\begin{aligned} E[X] &= \int_0^{\infty} x f_X(x) dx \\ &\geq \int_a^{\infty} x f_X(x) dx \\ &\geq \int_a^{\infty} a f_X(x) dx \\ &= a \cdot P(X \geq a) \end{aligned}$$

and so $P(X \geq a) \leq \frac{E[X]}{a}$.

Again, this is a very general statement in that we have assumed nothing about X other than it is positive. The price for the generality is that the bound is typically very loose, and does not usually capture the behavior of $P(X \geq a)$. We can, however, cleverly apply the Markov inequality to get something slightly more useful.

Chebyshev inequality

The main use of the Markov inequality turns out to be its use in deriving other, more accurate deviation inequalities. Here we will use it to derive the **Chebyshev inequality**, from which the weak law of large numbers will follow immediately.

Chebyshev inequality: If X is a random variable with mean μ and variance σ^2 , then

$$P(|X - \mu| > c) \leq \frac{\sigma^2}{c^2} \quad \text{for all } c > 0.$$

The Chebyshev inequality follows immediately from the Markov inequality in the following way. No matter what range of values X takes, the quantity $|X - \mu|^2$ is always positive. Thus

$$P(|X - \mu|^2 > c^2) \leq \frac{E[|X - \mu|^2]}{c^2} = \frac{\sigma^2}{c^2}.$$

Since squaring $(\cdot)^2$ is monotonic (invertible) over positive numbers,

$$P(|X - \mu|^2 > c^2) = P(|X - \mu| > c) \leq \frac{\sigma^2}{c^2}.$$

We now have a bound which depends on the mean and the variance of X ; this leads to a more accurate approximation of the probability.

Simple proof of the weak law of large numbers

We now turn to the behavior of the the sample mean

$$M_N = \frac{X_1 + X_2 + \cdots + X_N}{N},$$

where again the X_i are iid random variables with $E[X_i] = \mu$ and $\text{var } X_i = \sigma^2$. We know that

$$E[M_N] = \frac{E[X_1] + E[X_2] + \cdots + E[X_N]}{N} = \frac{N\mu}{N} = \mu,$$

and since the X_i are independent,

$$\text{var}(M_N) = \frac{\text{var}(X_1) + \text{var}(X_2) + \cdots + \text{var}(X_N)}{N^2} = \frac{N\sigma^2}{N^2} = \frac{\sigma^2}{N}.$$

For any $\epsilon > 0$, a direct application of the Chebyshev inequality tells us that

$$P(|M_N - \mu| > \epsilon) \leq \frac{\sigma^2}{N\epsilon^2}.$$

The point is that this gets arbitrarily small as $N \rightarrow \infty$ no matter what ϵ was chosen to be. We have established, in some sense, that even though $\{M_N\}$ is a sequence of random numbers, it converges to something deterministic, namely μ .

WLLN: Let X_1, X_2, \dots be iid random variables as above. For **every** $\epsilon > 0$, we have

$$P(|M_N - \mu| > \epsilon) = P\left(\left|\frac{X_1 + \dots + X_N}{N} - \mu\right| > \epsilon\right) \longrightarrow 0,$$

as $N \rightarrow \infty$.

One of the philosophical consequences of the WLLN is that it tells us that probabilities can be estimated through **empirical frequencies**. Suppose I want to estimate the probability of an event A occurring related to some probabilistic experiment. We run a series of (independent) experiments, and set $X_i = 1$ if A occurred in experiment i , and $X_i = 0$ otherwise. Then given X_1, \dots, X_N , we estimate the probability of A in a completely reasonable way, by computing the percentage of times it occurred:

$$p_{\text{empirical}} = \frac{X_1 + \dots + X_N}{N}.$$

The WLLN tells us that

$$p_{\text{empirical}} \rightarrow P(A), \quad \text{as } N \rightarrow \infty.$$

This lends some mathematical weight to our interpretation of probabilities as *relative frequencies*.

All of the above of course applies to functions of random variables. That is, if X is a random variable, and $g(X)$ is a function of that random variable with

$$\text{var}(g(X)) = E[(g(X) - E[g(X)])^2] < \infty,$$

then given independent realizations X_1, \dots, X_N , we have

$$\frac{1}{N} \sum_{n=1}^N g(X_n) \rightarrow \mathbb{E}[g(X)]$$

as $N \rightarrow \infty$.

Minimum Mean-Square Error Estimation

Now we will take our first look at estimating variables that are themselves random, subject to a known probability law.

We start our discussion with a very basic problem. Suppose Y is a scalar random variable with a known pdf $f_Y(y)$. Here is a fun game: you guess what Y is going to be, then I draw a realization of Y corresponding to its probability law, then we see how close you were with your guess.

What is your best guess?

Well, that of course depends on what exactly we mean by “best”, i.e. what price I pay for being a certain amount off. But if we penalize the *mean-squared error*, we know exactly how to minimize it.

Let g be your guess. The error in your guess is of course random (since the realization of Y is random), and so is the squared-error $(Y - g)^2$. We want to choose g so that the mean of the squared error is as small as possible:

$$\underset{g}{\text{minimize}} \quad \mathbb{E}[(Y - g)^2].$$

Expanding the squared error makes it clear how to do this:

$$\mathbb{E}[(Y - g)^2] = \mathbb{E}[Y^2] - 2g \mathbb{E}[Y] + g^2.$$

No matter what the first moment $\mathbb{E}[Y]$ and second moment $\mathbb{E}[Y^2]$ are (as long as they are finite), the expression above is a convex quadratic function in g , and hence is minimized when its first derivative (w.r.t. g) is zero, i.e. when

$$-2 \mathbb{E}[Y] + 2g = 0 \quad \Rightarrow \quad \hat{g} = \mathbb{E}[Y].$$

The mean squared error for this choice \hat{g} is of course exactly the variance of Y ,

$$\mathbb{E}[(Y - \hat{g})^2] = \mathbb{E}[(Y - \mathbb{E}[Y])^2] = \text{var}(Y).$$

The story gets more interesting (and relevant) when we have multiple random variables, some of which we observe, some of which we do not. Suppose that two random variables (Y, Z) have joint pdf $f_{Y,Z}(y, z)$. Suppose that a realization of (Y, Z) is drawn, and I get to observe Z . What have I learned about Y ?

If Y and Z are independent, then the answer is of course nothing. But if they are not independent, then the marginal distribution of Y changes. In particular, before the random variables were drawn, the (marginal) pdf for Y was

$$f_Y(y) = \int f_{Y,Z}(y, z) dz.$$

After we observe $Z = z$, we have

$$f_Y(y|Z = z) = \frac{f_{Y,Z}(y, z)}{f_Z(z)} = \frac{f_{Y,Z}(y, z)}{\int f_{Y,Z}(y, z) dy}.$$

Y is still a random variable, but its distribution depends on the value z that was observed for Z .

Now, given that I have observed $Z = z$, what is the best guess for Y ? If by “best” we mean that which minimizes the mean squared error, it is the conditional mean. That is, the minimizer of

$$\underset{g}{\text{minimize}} \quad \mathbb{E}[(Y - g)^2|Z = z]$$

is

$$\hat{g} = E[Y|Z = z].$$

Notice that unlike before, \hat{g} is not pre-determined, it depends on the outcome $Z = z$. We might denote

$$\hat{g}(z) = E[Y|Z = z].$$

For a particular choice of z , the mean-squared error is the conditional variance

$$\begin{aligned} E[(Y - \hat{g}(z))^2|Z = z] &= E[(Y - E[Y|Z = z])^2|Z = z] \\ &= E[Y^2|Z = z] - (E[Y|Z = z])^2 \\ &= \text{var}(Y|Z = z). \end{aligned}$$

So in general, not only does \hat{g} depend on z , but its performance (its mean-square error) also depends on z .

We can also average over the draw of Z . First, note that since Z is a random variable, \hat{g} is a priori also random, we might say

$$\hat{g}(Z) = E[Y|Z].$$

Let me pause here because this is the point where many people start to get confused. The quantities

$$E[Y] = \int_{-\infty}^{\infty} y f_Y(y) dy, \quad \text{and} \quad E[Y|Z = z] = \int_{-\infty}^{\infty} y f_Y(y|Z = z) dy$$

are deterministic, but we re-emphasize that

$$E[Y|Z] = \int_{-\infty}^{\infty} y f_Y(y|Z) dy$$

is *random*. The above integrates out the randomness of Y , but not that of Z . (Note that in $E[Y|Z = z]$ the randomness in Z is removed through direct observation.)

The “average estimate” is now

$$\begin{aligned} E[g(Z)] &= E[E[Y|Z]] \\ &= \int E[Y|Z = z] f_Z(z) dz \\ &= E[Y]. \end{aligned}$$

The identity $E[E[Y|Z]] = E[Y]$ is known as the law of *iterated expectation* or *total expectation*. The inside E above integrates out the randomness in Y while the outside one integrates over Z , the result is a deterministic quantity.

Since $E[g(Z)] = E[Y]$, on average we are doing the same thing as if we didn’t observe Z at all. But since we are adapting g to the draw of Z , we get better average performance. The mean square error (which is random with Z) is

$$E[(Y - \hat{g}(Z))^2|Z] = E[Y^2|Z] - (E[Y|Z])^2 = \text{var}(Y|Z).$$

The average performance is then

$$E[E[(Y - \hat{g}(Z))^2|Z]] = E[\text{var}(Y|Z)] \leq \text{var}(Y).$$

The last inequality, which follows from the law of total variance discussed below, means that on average, an (optimal) estimator using knowledge of Z will outperform an (optimal) estimator without knowledge of Z . Of course, when Y and Z are uncorrelated (meaning $E[YZ] = E[Y]E[Z]$) we will have $E[\text{var}(Y|Z)] = \text{var}(Y)$ and knowing Z makes no difference.

The law of total variance

Recall that for any random variable Y ,

$$\text{var}(Y) = E[Y^2] - (E[Y])^2. \quad (1)$$

As we have seen above, $E[Y|Z]$ is a random variable (where now the randomness is being caused by Z). Hence it also has a mean

$$E[E[Y|Z]] = E[Y],$$

and a variance

$$\begin{aligned} \text{var}(E[Y|Z]) &= E[(E[Y|Z])^2] - (E[E[Y|Z]])^2 \\ &= E[(E[Y|Z])^2] - (E[Y])^2. \end{aligned} \quad (2)$$

The quantity (again as we have seen above) $\text{var}(Y|Z)$ is also a random variable; we can write its mean as

$$\begin{aligned} E[\text{var}(Y|Z)] &= E[E[(Y - E[Y|Z])^2 | Z]] \\ &= E[E[Y^2|Z] - E[(E[Y|Z])^2]] \\ &= E[Y^2] - E[(E[Y|Z])^2]. \end{aligned} \quad (3)$$

Adding together (2) and (3) and applying (1) gives us the cute expression

$$\text{var}(Y) = E[\text{var}(Y|Z)] + \text{var}(E[Y|Z]).$$

This is known as the **law of total variance**. It basically says that you can decompose the variance in a random variable by the expected variance in “using Z to predict Y ” and the variance of the prediction (condition expectation) itself. Note that all of the quantities above are non-negative, so we have the inequalities

$$E[\text{var}(Y|Z)] \leq \text{var}(Y) \quad \text{and} \quad \text{var}(E[Y|Z]) \leq \text{var}(Y).$$