

# EECS 126 Notes

Matthew Lacayo

April 11, 2021

These notes will follow along with the teachings of the book Introduction to Probability, Statistics, and Random Processes by Hossein Pishro-Nik. Additional notes will also be included that follow along with the lecture teaching of professor Tom Courtade from Spring 2021.

## Contents

<b>1</b>	<b>Multiple Random Variables</b>	<b>2</b>
1.1	Joint Distributions and Independence . . . . .	2
1.2	Moment Generating Functions . . . . .	2
1.3	Characteristic Functions . . . . .	3
1.4	Exercises Part 1 . . . . .	3
1.5	Markov and Chebyshev Inequalities . . . . .	4
<b>2</b>	<b>Limit Theorems and Convergence of Random Variables</b>	<b>6</b>
2.1	Limit Theorems . . . . .	6
2.2	Convergence of Random Variables . . . . .	7
<b>3</b>	<b>Random Processes</b>	<b>9</b>
3.1	Poisson Processes . . . . .	9
<b>4</b>	<b>Hypothesis Testing</b>	<b>11</b>
4.1	MAP and MLE . . . . .	11

# 1 Multiple Random Variables

## 1.1 Joint Distributions and Independence

### Definition 1.1: Joint PMF

Let  $X_1, X_2, \dots, X_n$  be  $n$  discrete random variables. The joint pmf of  $v$  is defined as:

$$P_{X_1, X_2, \dots, X_n}(X_1, X_2, \dots, X_n) = P(X_1 = x_1, X_2 = x_2, \dots, X_n = x_n)$$

For continuous random variables, the joint PDF can be defined similarly using integrals. The marginal PDF can be reconstructed by integrating out all other random variables from the joint PDF.

### Definition 1.2: Joint CDF

The joint CDF of  $n$  random variables  $X_1, X_2, \dots, X_n$  is defined as:

$$F_{X_1, X_2, \dots, X_n}(X_1, X_2, \dots, X_n) = F(X_1 \leq x_1, X_2 \leq x_2, \dots, X_n \leq x_n)$$

Earlier results regarding independence of 2 random variables still hold for  $n$  random variables. Namely, if a group of  $n$  random variables are independent, then their joint PDF is the same as the product of the marginal PDFs, the joint CDF is the product of the marginal CDFs, the joint PMF is the product of the marginal PMFs, and the expected value of the product of the random variables is the same as the product of the expected values. Moreover, the results regarding sums of independent random variables still hold, and so we can still make use of the equations we developed for expected value and variance.

## 1.2 Moment Generating Functions

### Definition 1.3: Moments

We define the  $n$ th moment of a random variable  $X$  to be  $E[X^n]$ .

### Definition 1.4: Moment Generating Function (MGF)

The moment generating function of a random variable  $X$  is defined to be the function  $M_X(s)$  where:

$$M_X(s) := E[e^{sX}]$$

The MGF of  $X$  exists if there exists a positive constant  $a$  such that  $M_X(s)$  is finite for all  $s \in [-a, a]$ .

The MGF is useful because it allows us to extract all the moments of  $X$  simply by taking derivatives (follows by considering the Taylor series expansion of  $e^{sX}$ ). Namely,

$$E[X^k] = \frac{d^k}{ds^k} M_X(s)|_{s=0}$$

The MGF also uniquely defines the distribution of  $X$ . I.e. if two random variables  $X, Y$  have the same MGF, then they also have the same distribution.

**Theorem 1.1**

If  $X_1, X_2, \dots, X_n$  are  $n$  independent random variables, then:

$$M_{X_1, X_2, \dots, X_n}(s) = M_{X_1}(s)M_{X_2}(s) \cdots M_{X_n}(s)$$

**Example 1.1**

Let  $X \sim \text{Binomial}(n, p)$ . Find the MGF of  $X$ .

Since  $X$  is binomial, we can think of it as the sum of  $n$  independent Bernoulli random variables:

$$X = X_1 + X_2 + \cdots + X_n$$

Since the  $X_i$  are i.i.d., we know that  $M_X(s) = M_{X_1}(s)M_{X_2}(s) \cdots M_{X_n}(s)$ .

$$M_{X_i}(s) = E[e^{sX_i}] = p(e^s) + (1-p)(e^0) = (1-p) + pe^s$$

Thus,

$$M_X(s) = (1-p + pe^s)^n$$

**Example 1.2**

Use MGFs to prove that if  $X \sim \text{Binomial}(m, p)$  and  $Y \sim \text{Binomial}(n, p)$  are independent, then  $X + Y \sim \text{Binomial}(m + n, p)$ .

Solution: The MGF of  $X + Y$  is:

$$M_X(s)M_Y(s) = (1-p + pe^s)^m \cdot (1-p + pe^s)^n = (1-p + pe^s)^{m+n}$$

Which is the MGF of  $Z \sim \text{Binomial}(m + n, p)$ .

**1.3 Characteristic Functions**

To be filled in.

**1.4 Exercises Part 1****Example 1.3**

Let  $X, Y, Z$  be three independent random variables with  $X \sim N(\mu, \sigma^2)$ , and  $Y, Z \sim \text{Uniform}(0, 2)$ . We also know that:

- $E[X^2Y + XYZ] = 13$
- $E[XY^2 + ZX^2] = 14$

Find  $\mu$  and  $\sigma$ .

Solution: From equation 1, we get that  $E[X^2] = 13 - \mu$ . From equation 2, we get that  $E[Y^2] = \frac{14 - E[X^2]}{\mu} = \frac{1 + \mu}{\mu}$ . We also know that  $\text{Var}(Y) = \text{Var}(Z) = \frac{1}{3}$ . Thus,  $E[Y^2] = 1 + \frac{1}{3} = \frac{4}{3}$ . This gives  $\mu = 3$ . Finally,  $\text{Var}(X) = \sigma^2 = 13 - \mu - \mu^2 = 1$ .

### Example 1.4

Let  $X_1, X_2, X_3$  be 3 i.i.d. *Bernoulli*( $p$ ) random variables and:

- $Y_1 = \max(X_1, X_2)$
- $Y_2 = \max(X_1, X_3)$
- $Y_3 = \max(X_2, X_3)$
- $Y = Y_1 + Y_2 + Y_3$

Find  $E[Y]$ .

Solution:  $E[Y] = 3E[Y_1]$ .  $E[Y_1] = (1 - (1 - p)^2) = p(2 - p)$ . Thus,  $E[Y] = 3p(2 - p)$ .

## 1.5 Markov and Chebyshev Inequalities

### Theorem 1.2: Markov's Inequality

If  $X$  is a nonnegative random variable, then for any  $a > 0$ :

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

### Theorem 1.3: Chebyshev's Inequality

If  $X$  is any random variable, then for any  $b > 0$  we have:

$$P(|X - E[X]| \geq b) \leq \frac{\text{Var}(X)}{b^2}$$

### Theorem 1.4

If  $X$  is a random variable, then for any  $a \in \mathbb{R}$  we have:

- $P(X \geq a) \leq e^{-sa} M_X(s)$ , for all  $s > 0$
- $P(X \leq a) \leq e^{-sa} M_X(s)$ , for all  $s < 0$

*Proof.* We start by seeing that:

- $P(X \geq a) = P(e^{sX} \geq e^{sa})$ , for all  $s > 0$
- $P(X \leq a) = P(e^{sX} \geq e^{sa})$ , for all  $s < 0$

For  $s > 0$ , we can use Markov's inequality to write:

$$P(e^{sX} \geq e^{sa}) \leq \frac{E[e^{sX}]}{e^{sa}} = e^{-sa} M_X(s)$$

A similar result follows for  $s < 0$ .

□

This bound is useful because it depends on a new parameter:  $s$ . That is to say, we can optimize over  $s$  to find the tightest bound.

### Theorem 1.5: Cauchy-Schwarz Inequality

For any two random variables  $X$  and  $Y$ , we have:

$$|E[XY]| \leq \sqrt{E[X^2]E[Y^2]}$$

Where equality holds if and only if  $X = \alpha Y$  for some constant  $\alpha \in \mathbb{R}$ .

## 2 Limit Theorems and Convergence of Random Variables

### 2.1 Limit Theorems

#### Definition 2.1: Sample Mean

For i.i.d. random variables  $X_1, X_2, \dots, X_n$ , the **sample mean**, denoted by  $\bar{X}$  or  $M_n$ , is defined as:

$$\bar{X} := \frac{X_1 + X_2 + \dots + X_n}{n}$$

Note that since each  $X_i$  is a random variable, the sample mean is also a random variable. It can easily be shown that:

- $E[\bar{X}] = E[X_i]$
- $Var(\bar{X}) = \frac{Var(X)}{n}$

#### Theorem 2.1: Weak Law of Large Numbers (WLLN)

Let  $X_1, X_2, \dots, X_n$  be i.i.d. random variables with a finite expected value  $E[X_i] = \mu < \infty$ . Then, for any  $\epsilon > 0$  we have:

$$\lim_{n \rightarrow \infty} P(|\bar{X} - \mu| \geq \epsilon) = 0$$

We can easily prove this if we make the assumption that the variance is finite, as we can then use chebyshevs inequality and show that the term on the right of the inequality goes to zero as n goes to infinity.

#### Definition 2.2: Normalized Variable

Suppose  $X_1, X_2, \dots, X_n$  are i.i.d. random variables with  $E[X_i] = \mu < \infty$  and  $Var(X_i) = \sigma^2 < \infty$ . We define the normalized random variable  $Z_n$  as:

$$Z_n = \frac{\bar{X} - E[\bar{X}]}{\frac{\sigma}{\sqrt{n}}} = \frac{X_1 + X_2 + \dots + X_n - n\mu}{\sqrt{n}\sigma}$$

Also, note that  $E[Z_n] = 0$  and  $Var(Z_n) = 1$ .

#### Theorem 2.2: The Central Limit Theorem (CLT)

Let  $X_1, X_2, \dots, X_n$  be i.i.d. random variables with  $E[X_i] = \mu < \infty$  and  $Var(X_i) = \sigma^2 < \infty$ . Then,  $Z_n$  converges in distribution to the standard normal random variable as n goes to infinity. That is, for all  $x \in \mathbb{R}$ :

$$\lim_{n \rightarrow \infty} P(Z_n \leq x) = \Phi(x)$$

Where  $\Phi(x)$  is the standard normal CDF.

### Example 2.1

In a communication system, each data packet consists of 1000 bits. Due to the noise, each bit may be received in error with probability 0.1. It is assumed bit errors occur independently. Find the probability that there are more than 120 errors in a certain data packet.

Solution: Let  $Y = X_1 + X_2 + \cdots + X_{1000}$  where each  $X_i \sim \text{Bernoulli}(0.1)$ . Thus,  $E[X_i] = 0.1$  and  $\text{Var}(X_i) = 0.09$ . We want to find  $P(Y \leq 120)$ .

$$P(Y \leq 120) = P\left(\frac{Y - n\mu}{\sigma\sqrt{n}} \leq \frac{120 - \mu}{\sigma\sqrt{n}}\right) = P\left(Z_n \leq \frac{120 - 100}{\sqrt{0.09}\sqrt{1000}}\right) \approx \Phi\left(\frac{20}{\sqrt{90}}\right)$$

Finally,  $1 - \Phi\left(\frac{20}{\sqrt{90}}\right) = 0.0175$

## 2.2 Convergence of Random Variables

### Definition 2.3: Convergence in Distribution

A sequence of random variables  $X_1, X_2, X_3, \dots$  converges **in distribution** to a random variable  $X$ , denoted as  $X_n \xrightarrow{d} X$ , if:

$$\lim_{n \rightarrow \infty} F_{X_n}(x) = F_X(x)$$

for all  $x$  at which  $F_X(x)$  is continuous.

### Example 2.2

Let  $X_2, X_3, X_4, \dots$  be a sequence of random variables such that:

$$F_{X_n}(x) = \begin{cases} 1 - (1 - \frac{1}{n})^{nx}, & \text{for } x > 0 \\ 0, & \text{otherwise} \end{cases}$$

Show that  $X_n \xrightarrow{d} \text{Expo}(1)$ .

Solution:

$$\lim_{n \rightarrow \infty} 1 - (1 - \frac{1}{n})^{nx} = 1 - e^{-x}$$

This is the CDF of  $\text{Expo}(1)$  as desired.

### Definition 2.4: Convergence in Probability

A sequence of random variables  $X_1, X_2, X_3, \dots$  converges **in probability** to a random variable  $X$ , denoted as  $X_n \xrightarrow{p} X$ , if for all  $\epsilon > 0$ :

$$\lim_{n \rightarrow \infty} P(|X_n - X| \geq \epsilon) = 0$$

### Example 2.3

Let  $X_n \sim \text{Expo}(n)$ . Show that  $X_n \xrightarrow{p} 0$  i.e. to the random variable that is always 0.

Solution:

$$\lim_{n \rightarrow \infty} P(|X_n| > \epsilon) = |e^{-n\epsilon}| = 0$$

as desired.

### Example 2.4

Reread the WLLN. This is an example of convergence in probability.

### Theorem 2.3: Almost Sure Convergence

A sequence of random variables  $X_1, X_2, X_3, \dots$  converges **almost surely** to a random variable  $X$ , denoted as  $X_n \xrightarrow{a.s.} X$ , if:

$$P(\{s \in S : \lim_{n \rightarrow \infty} X_n(s) = X(s)\}) = 1$$

This is saying that the set of samples for which  $X_n = X$  forms an event of probability 1 as  $n$  goes to infinity.

### Theorem 2.4: Strong Law of Large Numbers

Let  $X_1, X_2, \dots, X_n$  be i.i.d. random variables with  $E[X_i] = \mu < \infty$ . Let also:

$$M_n := \frac{X_1 + X_2 + \dots + X_n}{n}$$

Then  $M_n \xrightarrow{a.s.} \mu$ .



## 3 Random Processes

### 3.1 Poisson Processes

#### Definition 3.1: Counting Process

A random process  $\{N(t), t \in [0, \infty)\}$  is said to be a **counting process** if  $N(t)$  is the number of events occurred from time 0 up to and including time  $t$ . For a counting process, we assume:

- $N(0) = 0$
- $N(t) \in \{0, 1, 2, \dots\}$ , for all  $t \in [0, \infty)$
- For  $0 \leq s < t$ ,  $N(t) - N(s)$  shows the number of events that occur in the interval  $(s, t]$

We typically refer to the occurrence of an event as an 'arrival'.

#### Definition 3.2: Poisson Process

Let  $\lambda > 0$  be fixed. The counting process  $\{N(t), t \in [0, \infty)\}$  is called a **Poisson Process** with rate  $\lambda$  if all the following conditions hold:

- $N(0) = 0$
- $N(t)$  has independent increments (disjoint intervals are independent)
- The number of arrivals in any interval of length  $\tau > 0$  has  $Poisson(\lambda\tau)$  distribution.

Note that the number of arrivals in any given interval depends only on the length of the interval, not on the location of the interval.

#### Example 3.1

The number of customers arriving at a grocery store can be modeled by a Poisson process with intensity  $\lambda = 10$  customers per hour. Find the probability that there are 2 customers between 10:00 and 10:20.

Solution: We know that the length of the interval  $\tau = \frac{20}{60} = \frac{1}{3}$ . Thus, the distribution of arrivals in the 20 minute interval is  $X \sim Poisson(10 \cdot \frac{1}{3} = \frac{10}{3})$ . Finally,

$$P(X = 2) = \frac{e^{-\frac{10}{3}} (\frac{10}{3})^2}{2!} \approx 0.2$$

Now, let us consider arrival and interarrival times for a Poisson process.

### Example 3.2

Let  $N(t)$  be a Poisson process with rate  $\lambda$ . Let  $X_1$  be the time of the first arrival. Then,

$$\begin{aligned} P(X_1 > t) &= P(\text{no arrival in } (0, t]) \\ &= \frac{(\lambda t)^0}{0!} e^{-(\lambda t)} \\ &= e^{-\lambda t} \end{aligned}$$

Thus,

$$F_{X_1}(t) = 1 - e^{-\lambda t}$$

This is the CDF of Exponential( $\lambda$ ), and so  $X_1 \sim \text{Exponential}(\lambda)$ .

Let  $X_2$  be the time elapsed between the first and the second arrival. Since the probability of an arrival in disjoint intervals is independent, we also have that  $X_2 \sim \text{Exponential}(\lambda)$ . This yields the following definition.

### Definition 3.3: Interarrival Times for Poisson Processes

If  $N(t)$  is a Poisson process with rate  $\lambda$ , then the interarrival times  $X_1, X_2, \dots$  are independent and for  $i = 1, 2, 3, \dots$ :

$$X_i \sim \text{Exponential}(\lambda)$$

### Definition 3.4: Arrival Times for Poisson Processes

If  $N(t)$  is a Poisson process with rate  $\lambda$ , then the arrival times  $T_1, T_2, \dots$  have *Erlang*( $n, \lambda$ ) distribution. I.e.

$$T_i = X_1 + X_2 + \dots + X_i$$

Thus,

- $E[T_i] = \frac{n}{\lambda}$
- $Var(T_i) = \frac{n}{\lambda^2}$

Now, we will discuss merging and splitting Poisson processes. That is, we will consider taking two independent Poisson processes and adding them together, or taking one poisson process and splitting its arrivals into multiple 'streams' according to some defined distribution.

## 4 Hypothesis Testing

We now turn our attention to the methods to best utilize models, as well as how we can make models.

The goal of a "model" is to define the transition probabilities between some state of nature (or hypothesis), and some observation that we see potentially resulting from the underlying state of nature. For example, given that it is cloudy (state of nature), what is the probability that it rains (observation)? Since we cannot entirely capture the inner workings of nature to determine the exact ways that a given set of clouds could cause rain, we do not know the exact underlying probability of rain given that it is cloudy. However, we can seek ways to approximate the workings of nature by approximating the underlying probability. Moreover, we would like a model to define all relevant transition probabilities (such as probability of sun given cloudy, rain given no clouds, rain given some clouds, etc.) We refer to the probability distribution of the state of nature as the **prior**, and we usually denote the probability distribution with the symbol  $\pi$ . In our example, the prior would be the probability distribution of the different states of the sky. For example, we would have  $\pi(\text{cloudy}) = P(\text{cloudy})$ ,  $\pi(\text{sunny}) = P(\text{sunny})$ , etc.

Note that our discussion has been using the phrase "state of nature", however it will be helpful to begin thinking of this as a "hypothesis" instead. The justification for this choice will be more clear after reading the coming sections. However, for now, whenever you see the phrase hypothesis, try thinking of it as a "state of nature" instead if that makes more intuitive sense.

### 4.1 MAP and MLE

Say that we are given some observation, and we would like to know under what conditions would the likelihood of observing this outcome be maximized. In other words, after seeing some evidence, we now want to know which hypothesis we believe is the most likely candidate for producing the evidence. This could be helpful if we are not able to observe the inner working of a system and are only able to observe some output of that system, but we still want to be able to make some predictions about what is going on behind the scenes. Let us now formalize our goal. Assume that the possible hypotheses lie in a set  $X$ , and the possible observations lie in a set  $Y$ . Given some observation  $y \in Y$ , we want to find:

$$\arg \max_{x \in X} p_{x|y}(x|y)$$

How can we interpret these probabilities? We know that  $\pi(x)$  is the *prior* of  $x$ , and it represents the probability of the hypothesis being  $x$  (no evidence observed). We can also interpret  $p_{x|y}(x|y)$  as being our updated belief about the probability of the hypothesis being  $x$  now that we have observed an observation  $y$ . In other words, observing  $y$  gave us new information, and so we now want to update our beliefs on the hypotheses based on this information. Using this way of thinking is critical to understanding the utility of Bayes Law. Remember that the model gives us access to the transition probabilities between the hypotheses and the states of nature:  $p_{y|x}(y|x)$ , and we know that we can relate this to  $p_{x|y}(x|y)$  given Bayes Law, and so we can convert our goal into the following:

$$\arg \max_{x \in X} p_{x|y}(x|y) = \arg \max_{x \in X} \frac{p_{y|x}(y|x)p(x)}{p(y)} = \arg \max_{x \in X} \frac{p_{y|x}(y|x)\pi(x)}{p(y)}$$

Note that since our goal is to find  $x \in X$  that maximizes the expression, we can ignore the denominator. This is because the denominator depends only on  $y$ , and so maximizing the

entire expression amounts to maximizing what we can control, which in this case is just the numerator. Therefore, to find our guess for which hypothesis caused the evidence that we saw, we define:

$$\hat{X}_{MAP}(y) := \underbrace{\arg \max_{x \in X} p_{x|y}(x|y)}_{\text{"A Posteriori Probability"}} = \underbrace{\arg \max_{x \in X} p_{y|x}(y|x)\pi(x)}_{\text{Computable with model}}$$

This is called the **MAP Estimate**, and it is a function of the observation. MAP stands for **Maximum A Posteriori Probability**. A Posteriori Probability refers to the "after probability", meaning the updated probability distribution of the hypothesis given we observe some piece of evidence. The above equation shows how we can think about MAP. Our goal is to maximize the A Posteriori Probability, and we reduced this problem into one that we can solve by utilizing the model.

We call MAP a type of **estimator**. In particular, MAP is useful when we know the prior (the probability distribution of the hypotheses). However, what if we do not have access to this distribution? One possible solution is to assume that the prior is uniform among (us) all hypotheses: i.e.  $\pi(x) = 1/|X|$ . If we try to now compute the MAP estimate using this assumption, we will be able to treat  $\pi(x)$  as a constant (since it is the same for any  $x$ ), and so we will be able to pull it out of the argmax expression entirely. We will get a different type of estimate called the **MLE** or **Maximum Likelihood Estimate**:

$$\hat{X}_{MLE}(y) := \arg \max_{x \in X} p_{y|x}(y|x)$$

This has a very nice interpretation: we will estimate the hypothesis by choosing the hypothesis that has greatest probability of producing the observed outcome, regardless of the prior distribution of the hypotheses.