

Recitation Notes for PSTAT 120B

Haosheng Zhou

Jan, 2023

Week 1

This week we will review some of the important properties taught in 120A as a preparation for homework 0.

The first concept to review is the **continuous and discrete random variables**. Generally, a random variable

$$X(\omega) : \Omega \rightarrow \mathbb{R} \quad (1)$$

is a mapping from the sample space to the real numbers, i.e. it assigns a value to each possible outcome of random experiment. Discrete random variables can take countably many values while continuous random variables can take uncountably many values. For example, if we want to consider the random variable X as the outcome after rolling one dice, then we have to first specify the **sample space**, i.e. the set of all possible outcomes rolling one dice, which should be $\Omega = \{1, 2, \dots, 6\}$. As a result, such random variable X is defined as

$$X(\omega) = \omega \quad (2)$$

an identity map. Since X can only take values in $\{1, 2, \dots, 6\}$, a finite set, it's a discrete random variable.

To describe a single random variable, we have the **cumulative distribution function (CDF)** defined for any random variable X as

$$F(x) = \mathbb{P}(X \leq x) \quad (3)$$

Such F is always right-continuous, increasing and $F(-\infty) = 0, F(+\infty) = 1$ (try to explain the meaning of those properties). In particular, for continuous random variable such F is continuous and for discrete random variable such F is a step function. For continuous random variables, assume that F is nice enough to be differentiable so $F' = f$ gives the **density** that characterizes the distribution of the continuous random variable (for random vectors, those concepts can be generalized).

To describe the relationship between two random variables, the most important property is **independence**. We call X, Y independent if

$$\forall x, y \in \mathbb{R}, \mathbb{P}(X \leq x, Y \leq y) = \mathbb{P}(X \leq x) \mathbb{P}(Y \leq y) \quad (4)$$

which can also be explained in the sense of conditional probability (try to write the equality in the conditional form). For discrete r.v. X, Y , they are independent if and only if $\forall x, y \in \mathbb{R}, \mathbb{P}(X = x, Y = y) = \mathbb{P}(X = x) \mathbb{P}(Y = y)$ and for continuous r.v. X, Y , they are independent if and only if $f_X(x)f_Y(y) = f_{X,Y}(x, y)$ a.e. (think about why the criterion for discrete r.v. does not hold for continuous r.v.).

The important concept to mention is the **expectation** of continuous or discrete random variables. For discrete random variable X , assume that its distribution is given by

$$p_k = \mathbb{P}(X = a_k) \quad (k = 0, 1, \dots) \quad (5)$$

so the expectation is formed as

$$\mathbb{E}X = \sum_{k=0}^{\infty} a_k \cdot \mathbb{P}(X = a_k) = \sum_{k=0}^{\infty} a_k \cdot p_k \quad (6)$$

i.e., the **sum** of the product of the possible value a_k taken by X and the probability of X taking value a_k .

For discrete random variable X , assume that its density is $f(x)$, so the expectation is formed as

$$\mathbb{E}X = \int_{\mathbb{R}} x f(x) dx \quad (7)$$

i.e., the **integral** of the product of the possible value x taken by X and $f(x)$, the likelihood of X taking value x . In the homework, we will be asked to prove the **linearity of expectation** by using those definitions.

Another important concept is the **variance**, defined as

$$\text{Var}(X) = \mathbb{E}(X - \mathbb{E}X)^2 \quad (8)$$

the connection between variance and expectation can be given by the useful formula that

$$\text{Var}(X) = \mathbb{E}X^2 - (\mathbb{E}X)^2 \quad (9)$$

for two random variables, we can define the **covariance** to describe their relationship

$$\text{cov}(X, Y) = \mathbb{E}[(X - \mathbb{E}X)(Y - \mathbb{E}Y)] \quad (10)$$

and a similar identity holds that

$$\text{cov}(X) = \mathbb{E}XY - (\mathbb{E}X)(\mathbb{E}Y) \quad (11)$$

note that $\text{cov}(X, X) = \text{Var}(X)$ and that $\text{cov}(X, Y)$ is **bilinear**, i.e. $\text{cov}(aX + bY, Z) = a \cdot \text{cov}(X, Z) + b \cdot \text{cov}(Y, Z)$, $\text{cov}(Z, aX + bY) = a \cdot \text{cov}(Z, X) + b \cdot \text{cov}(Z, Y)$ and **symmetric**, i.e. $\text{cov}(X, Y) = \text{cov}(Y, X)$. This is especially useful when computing the variance of a linear combination. For example, if we want to write $\text{Var}(2X + 3Y)$ in terms of $\text{Var}(X), \text{Var}(Y)$,

$$\text{Var}(2X + 3Y) = \text{cov}(2X + 3Y, 2X + 3Y) \quad (12)$$

$$= 2\text{cov}(X, 2X + 3Y) + 3\text{cov}(Y, 2X + 3Y) \quad (13)$$

$$= 2[2\text{cov}(X, X) + 3\text{cov}(X, Y)] + 3[2\text{cov}(Y, X) + 3\text{cov}(Y, Y)] \quad (14)$$

$$= 4\text{Var}(X) + 12\text{cov}(X, Y) + 9\text{cov}(Y, Y) \quad (15)$$

you are asked to prove a more general version of this property in the homework.

Finally, let's talk about **normal distribution**. We say $X \sim N(\mu, \sigma^2)$ if it has density

$$f(x) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x-\mu)^2}{2\sigma^2}} \quad (x \in \mathbb{R}) \quad (16)$$

for the two parameters μ, σ^2 of normal distribution, a direct interpretation is that $\mathbb{E}X = \mu, \text{Var}(X) = \sigma^2$. You can try to prove those properties on your own by applying the definitions of expectation and variance to calculate the integrals. A trick will be that when calculating the integral

$$\mathbb{E}X = \int_{\mathbb{R}} x \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x-\mu)^2}{2\sigma^2}} dx \quad (17)$$

use the change of variables $u = \frac{x-\mu}{\sigma}$ to make the life easier

$$\mathbb{E}X = \int_{\mathbb{R}} x \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x-\mu)^2}{2\sigma^2}} dx \quad (18)$$

$$= \frac{1}{\sqrt{2\pi}\sigma} \sigma \int_{\mathbb{R}} (\sigma u + \mu) e^{-\frac{u^2}{2}} du \quad (19)$$

$$= \frac{1}{\sqrt{2\pi}} \int_{\mathbb{R}} (\sigma u + \mu) e^{-\frac{u^2}{2}} du \quad (20)$$

$$= \frac{\mu}{\sqrt{2\pi}} \int_{\mathbb{R}} e^{-\frac{u^2}{2}} du \quad (21)$$

$$= \mu \quad (22)$$

here we use the property that $ue^{-\frac{u^2}{2}}$ is an odd function and that $\int_{\mathbb{R}} e^{-\frac{u^2}{2}} du = \sqrt{2\pi}$ (this property can be deduced from the standard normal density, an easy way to remember). The calculation of variance is left to the reader.

The **standard normal CDF** is one of the most frequently used notations in statistics. The definition is

$$\Phi(x) = \int_{-\infty}^x \frac{1}{\sqrt{2\pi}} e^{-\frac{t^2}{2}} dt = \int_{-\infty}^x \frac{1}{\sqrt{2\pi}} e^{-\frac{t^2}{2}} dt = \mathbb{P}(G \leq x) \quad (G \sim N(0, 1)) \quad (23)$$

a property of Φ is that

$$\forall x \in \mathbb{R}, \Phi(x) + \Phi(-x) = 1 \quad (24)$$

to see this, notice that $\frac{1}{\sqrt{2\pi}} e^{-\frac{t^2}{2}}$ is an even function in t , so

$$\Phi(x) = \int_{-\infty}^x \frac{1}{\sqrt{2\pi}} e^{-\frac{t^2}{2}} dt = \int_{-x}^{+\infty} \frac{1}{\sqrt{2\pi}} e^{-\frac{u^2}{2}} du \quad (u = -t) \quad (25)$$

$$= \int_{-\infty}^{+\infty} \frac{1}{\sqrt{2\pi}} e^{-\frac{u^2}{2}} du - \int_{-\infty}^{-x} \frac{1}{\sqrt{2\pi}} e^{-\frac{u^2}{2}} du \quad (26)$$

$$= 1 - \Phi(-x) \quad (27)$$

Week 2

HW 0

For the problems in HW 0, let's look at problem 6 and 7 briefly. The important fact used in problem 6 is that for independent r.v. X, Y , it's true that $\mathbb{E}XY = \mathbb{E}X \cdot \mathbb{E}Y$. Let's prove this property for continuous random variables.

If X, Y are independent with density f, g , the joint density is $h(x, y) = f(x)g(y)$

$$\mathbb{E}XY = \int_{\mathbb{R}^2} xyh(x, y) dx dy \quad (28)$$

$$= \int_{\mathbb{R}^2} xyf(x)g(y) dx dy \quad (29)$$

$$= \int_{\mathbb{R}} xf(x) dx \cdot \int_{\mathbb{R}} yg(y) dy \quad (30)$$

$$= \mathbb{E}X \cdot \mathbb{E}Y \quad (31)$$

one can also try to prove the property in the discrete case.

For problem 7, the main idea is to tell you that often it's the case that you can greatly simplify the calculations by applying the properties of expectation or variance. For $X \sim N(\mu_x, \sigma_x^2), Y \sim N(\mu_y, \sigma_y^2)$ independent, by independence, the joint density is

$$f(x, y) = f_X(x)f_Y(y) = \frac{1}{\sqrt{2\pi}\sigma_x} e^{-\frac{(x-\mu_x)^2}{2\sigma_x^2}} \frac{1}{\sqrt{2\pi}\sigma_y} e^{-\frac{(y-\mu_y)^2}{2\sigma_y^2}} \quad (32)$$

the expectation can be calculated by linearity

$$\mathbb{E}(aX + bY + c) = a\mathbb{E}X + b\mathbb{E}Y + c = a\mu_x + b\mu_y + c \quad (33)$$

the second moment computed with the variance identity

$$\mathbb{E}X^2 = \text{Var}(X) + (\mathbb{E}X)^2 = \sigma_x^2 + \mu_x^2 \quad (34)$$

and the variance of linear combination is

$$\text{Var}(aX + bY + c) = \text{Var}(aX) + \text{Var}(bY) = a^2\sigma_x^2 + b^2\sigma_y^2 \quad (35)$$

note that generally the variance of sum does not equal the sum of variance, here it holds because of independence (actually this property holds if and only if X, Y are uncorrelated by the conclusion of problem 5).

HW 1

Let's talk about calculating the distribution of the transformation of a random variable. The most important idea comes from **the CDF method** that focuses on deriving the CDF of the transformed r.v.

To see how this method works, let's first look at some examples and then build up the theory for this method.

Now $X \sim N(0, 1)$, and we want to derive the PDF of $Y = |X|$ and to calculate $\mathbb{E}|X|$. The first step is to set up the CDF of Y , denoted $F_Y(y) = \mathbb{P}(Y \leq y)$, it's obvious that when $y < 0$ the CDF always has value 0 so we only have to consider the non-trivial case where $y \geq 0$. Denote $f_X(x)$ as the density of X so $f_X(x) = \frac{1}{\sqrt{2\pi}}e^{-\frac{x^2}{2}}$

$$F_Y(y) = \mathbb{P}(|X| \leq y) = \mathbb{P}(-y \leq X \leq y) = \int_{-y}^y f_X(x) dx = 2 \int_0^y f_X(x) dx \quad (36)$$

since $f_X(x)$ is an even function. Actually one does not have to calculate this integral, but to notice that PDF is the derivative of CDF, so taking derivative w.r.t. y on both sides gives

$$f_Y(y) = \frac{d}{dy} F_Y(y) = 2 \frac{d}{dy} \int_0^y f_X(x) dx = 2f_X(y) = \sqrt{\frac{2}{\pi}} e^{-\frac{y^2}{2}} \quad (y \geq 0) \quad (37)$$

the calculation of expectation follows

$$\mathbb{E}Y = \int_0^\infty y f_Y(y) dy \quad (38)$$

$$= \sqrt{\frac{2}{\pi}} \int_0^\infty y e^{-\frac{y^2}{2}} dy \quad (39)$$

$$= \sqrt{\frac{2}{\pi}} \int_0^\infty e^{-\frac{y^2}{2}} d\frac{y^2}{2} \quad (40)$$

$$= \sqrt{\frac{2}{\pi}} \int_0^\infty e^{-u} du \quad (41)$$

$$= \sqrt{\frac{2}{\pi}} \quad (42)$$

Remark. Do not forget that the density for Y only works on $[0, \infty)$ so it's necessary to label out $y \geq 0$.

Remark. One might have to take the derivative of an integral with variables in the integration region a lot when calculating the distribution of the transformation of r.v. As a result, one might find the following property from calculus useful:

$$\frac{d}{dx} \int_{f(x)}^{g(x)} h(t) dt = \frac{d}{dx} \int_0^{g(x)} h(t) dt - \frac{d}{dx} \int_0^{f(x)} h(t) dt \quad (43)$$

$$= h(g(x))g'(x) - h(f(x))f'(x) \quad (44)$$

to see this, one can consider $p(x) = \int_0^x h(t) dt$ and $\int_0^{g(x)} h(t) dt = p(g(x))$, so

$$\frac{d}{dx} \int_0^{g(x)} h(t) dt = \frac{d}{dx} p(g(x)) \quad (45)$$

$$= p'(g(x))g'(x) \quad (46)$$

$$= h(g(x))g'(x) \quad (47)$$

since $p'(x) = h(x)$ by Newton-Lebniz formula. So this is actually just an application of **the chain rule**.

The example above tells us the way to apply the CDF method, now let's build up the method in theory. Let's assume that we already know the PDF of X and want to get the PDF of $Y = h(X)$ with h to be strictly monotone increasing (this assumption is made to simplify the proof but not necessary).

$$F_Y(y) = \mathbb{P}(Y \leq y) = \mathbb{P}(h(X) \leq y) \quad (48)$$

$$= \mathbb{P}(X \leq h^{-1}(y)) \quad (49)$$

$$= \int_{-\infty}^{h^{-1}(y)} f_X(x) dx \quad (50)$$

take derivative w.r.t. y on both sides to get

$$f_Y(y) = \frac{d}{dy} F_Y(y) \quad (51)$$

$$= \frac{d}{dy} \int_{-\infty}^{h^{-1}(y)} f_X(x) dx \quad (52)$$

$$= f_X(h^{-1}(y)) \cdot \frac{d}{dy} h^{-1}(y) \quad (53)$$

by the calculations we have already made in the remark above. (This is part of the homework problem 8, please try to prove the other half when h is strictly decreasing on your own) Notice that the density has to be non-negative and here since h is increasing, $\frac{d}{dy} h^{-1}(y)$ has to be non-negative, making the density f_Y non-negative. For the case where h is decreasing, there is a slight difference in the sign that you have to notice. In all, the general formula is given by

$$f_Y(y) = f_X(h^{-1}(y)) \cdot \left| \frac{d}{dy} h^{-1}(y) \right| \quad (54)$$

for **any strictly monotone** h and is called **the transformation method**.

Remark. Although this method directly comes from the CDF method, one will see that in multi-dimensional case this method is much easier to generalize and to apply.

Now let me raise an example to show you how to apply this method. Consider $X \sim N(\mu, \sigma^2)$ and we want to find the PDF of $Y = \frac{X-\mu}{\sigma}$. It's immediate that $h(x) = \frac{x-\mu}{\sigma}$ is a linear function so it's strictly monotone, $h^{-1}(y) = \sigma y + \mu$

and $\frac{d}{dy}h^{-1}(y) = \frac{1}{\frac{dh(y)}{dy}} = \sigma$. Now it's clear that $f_X(x) = \frac{1}{\sqrt{2\pi}\sigma}e^{-\frac{(x-\mu)^2}{2\sigma^2}}$, plug in the formula to see

$$f_Y(y) = f_X(h^{-1}(y)) \cdot \left| \frac{d}{dy}h^{-1}(y) \right| \quad (55)$$

$$= \sigma \cdot f_X(\sigma y + \mu) \quad (56)$$

$$= \frac{1}{\sqrt{2\pi}}e^{-\frac{y^2}{2}} \quad (57)$$

as a result, we see that $Y \sim N(0, 1)$ is standard Gaussian. So we proved a very important **scaling property of Gaussian random variable** that $X \sim N(\mu, \sigma^2)$ if and only if $\frac{X-\mu}{\sigma} \sim N(0, 1)$.

The last method to talk about is **the method of moment generating function (MGF)**. This method is based on two properties of MGF defined as $M_X(t) = \mathbb{E}e^{tX}$. The first one is that MGF characterizes the distribution, so two random variables have the same MGF if and only if they have the same distribution. The second one is that for independent X, Y , $M_{X+Y}(t) = M_X(t)M_Y(t)$ the MGF of the sum is the product of respective MGF. **The MGF method is especially effective for dealing with Gaussian random variables.**

An example is that for $Y_1, Y_2, \dots, Y_n \sim N(0, 1)$ *i.i.d.*, let's calculate the distribution of $Z = a_1Y_1 + a_2Y_2 + \dots + a_nY_n$. One has to know that the MGF for $N(\mu, \sigma^2)$ Gaussian r.v. is $M(t) = e^{\mu t + \frac{\sigma^2}{2}t^2}$ (refer to the remark is you are not familiar with this conclusion).

$$M_Z(t) = M_{a_1Y_1}(t)M_{a_2Y_2}(t)\dots M_{a_nY_n}(t) \quad (58)$$

$$= \mathbb{E}e^{ta_1Y_1}\mathbb{E}e^{ta_2Y_2}\dots\mathbb{E}e^{ta_nY_n} \quad (59)$$

$$= M_{Y_1}(ta_1)\dots M_{Y_n}(ta_n) \quad (60)$$

$$= e^{\frac{a_1^2}{2}t^2}\dots e^{\frac{a_n^2}{2}t^2} \quad (61)$$

$$= e^{\frac{\sum_{i=1}^n a_i^2}{2}t^2} \quad (62)$$

comparing with the MGF for $N(\mu, \sigma^2)$, one immediately find that $Z \sim N(0, \sum_{i=1}^n a_i^2)$. This is telling us that **the linear combination of independent Gaussian r.v. must still be Gaussian**. (Try to do the same problem for $Y_i \sim N(\mu_i, \sigma_i^2)$ independent but not *i.i.d.* to see the conclusion that the linear combination is still Gaussian)

Remark. Let's calculate the MGF for $X \sim N(\mu, \sigma^2)$

$$M_X(t) = \mathbb{E}e^{tX} \quad (63)$$

$$= \int_{\mathbb{R}} e^{tx} f_X(x) dx \quad (64)$$

$$= \frac{1}{\sqrt{2\pi}\sigma} \int_{\mathbb{R}} e^{tx - \frac{(x-\mu)^2}{2\sigma^2}} dx \quad \left(u = \frac{x-\mu}{\sigma}\right) \quad (65)$$

$$= \frac{1}{\sqrt{2\pi}} \int_{\mathbb{R}} e^{t(\sigma u + \mu) - \frac{u^2}{2}} du \quad (66)$$

extract the constant term $e^{t\mu}$ to continue

$$M_X(t) = \frac{1}{\sqrt{2\pi}} e^{t\mu} \int_{\mathbb{R}} e^{t\sigma u - \frac{v^2}{2}} du \quad (67)$$

$$= \frac{1}{\sqrt{2\pi}} e^{t\mu} \int_{\mathbb{R}} e^{-\frac{1}{2}(u^2 - 2t\sigma u + t^2\sigma^2) + \frac{\sigma^2}{2}t^2} du \quad (68)$$

$$= \frac{1}{\sqrt{2\pi}} e^{\mu t + \frac{\sigma^2}{2}t^2} \int_{\mathbb{R}} e^{-\frac{(u-t\sigma)^2}{2}} du \quad (v = u - t\sigma) \quad (69)$$

$$= \frac{1}{\sqrt{2\pi}} e^{\mu t + \frac{\sigma^2}{2}t^2} \int_{\mathbb{R}} e^{-\frac{v^2}{2}} dv \quad (70)$$

$$= e^{\mu t + \frac{\sigma^2}{2}t^2} \quad (71)$$

Week3

It's always important that when calculating the distribution of the transformation of random variables, one choose the method that fits the best with the problem. Now we have three methods: the CDF method, the transformation method and the MGF method.

Generally, when dealing with the distribution of the sum of independent random variables, always use MGF method. The reason is that the MGF of the sum of independent r.v. is always the product of respective MGF. From the one-to-one correspondence between MGF and distribution, one would always find out the distribution of the independent sum easily.

For transformation method, it's always easy to apply when we see a transformation from $\mathbb{R}^n \rightarrow \mathbb{R}^n$, mapping a random vector of length n to another random vector of length n . The key point is that the dimension of the domain and image space of the transformation should be the same (because it depends on the determinant of the Jacobian as we will see later). As a result, transformation method won't be applied for problem like deriving the distribution of the sum of random variables, since it's actually mapping $(X_1, \dots, X_n) \in \mathbb{R}^n$ to $X_1 + \dots + X_n \in \mathbb{R}$. Moreover, there's some restrictions on the 'invertible' property of the transformation. For example, $Y = X^2$ has transformation $h(x) = x^2$ which is not invertible, so the transformation method will fail.

For distribution method, it's the most general method but also the method with the most calculations involved. It can be applied in all circumstances, regardless of the transformation function and the random variables one is using. Problems like $Y = |X|$, $Y = X^2$ can only be dealt with using the CDF method.

Let's look at some problems that consider the distribution of the average of *i.i.d.* random variables $\frac{S_n}{n} = \frac{X_1 + \dots + X_n}{n}$ to get familiar with the MGF method.

Let's first take X_1 following the Bernoulli distribution $B(1, p)$. Let's first calculate the MGF of X_1

$$M_{X_1}(t) = \mathbb{E}e^{tX_1} = 1 - p + pe^t \quad (72)$$

so now we see that

$$M_{\frac{S_n}{n}}(t) = \mathbb{E}e^{\frac{S_n}{n}t} = \mathbb{E}e^{\frac{t}{n}S_n} = M_{S_n}\left(\frac{t}{n}\right) \quad (73)$$

since S_n is *i.i.d.* sum, its MGF is the product of the respective MGF, so

$$M_{S_n}\left(\frac{t}{n}\right) = \left[M_{X_1}\left(\frac{t}{n}\right)\right]^n = (1 - p + pe^{\frac{t}{n}})^n \quad (74)$$

notice the trick to put the denominator n of $\frac{S_n}{n}$ into the variable in the MGF.

Similarly, we can compute the example for Poisson distribution $X_1 \sim P(\lambda)$

$$M_{X_1}(t) = \mathbb{E}e^{tX_1} = \sum_{k=0}^{\infty} e^{tk} \frac{\lambda^k}{k!} e^{-\lambda} = e^{\lambda e^t - \lambda} \quad (75)$$

so now we see that

$$M_{\frac{S_n}{n}}(t) = \mathbb{E}e^{\frac{S_n}{n}t} = \mathbb{E}e^{\frac{t}{n}S_n} = M_{S_n}\left(\frac{t}{n}\right) \quad (76)$$

since S_n is *i.i.d.* sum, its MGF is the product of the respective MGF, so

$$M_{S_n}\left(\frac{t}{n}\right) = \left[M_{X_1}\left(\frac{t}{n}\right)\right]^n = e^{n\lambda e^{\frac{t}{n}} - n\lambda} \quad (77)$$

one might try to prove the additivity of Poisson distribution as an exercise (if $\forall i = 1, 2, \dots, n, X_i \sim P(\lambda_i)$ are independent random variables, then $X_1 + \dots + X_n \sim P(\lambda_1 + \dots + \lambda_n)$) using MGF.

Week 4

Quiz Answer

1. We are playing a game of darts, where every throw results in the dart landing randomly somewhere on the dartboard, which has a radius of 1. If we say that X is the horizontal coordinate and Y is the vertical coordinate (both measured from the center/bullseye), then each throw results in a random pair (X, Y) with a joint density of

$$f_{X,Y}(x,y) = \begin{cases} \frac{1}{\pi} & x^2 + y^2 \leq 1 \\ 0 & \text{else} \end{cases} \quad (78)$$

(a) (6 points) What is the probability we are closer to the bullseye in the horizontal direction than in the vertical direction?

solution. The probability is $\mathbb{P}(|X| < |Y|)$. Representing as the integral of joint density

$$\mathbb{P}(|X| < |Y|) = \frac{1}{\pi} \int \int_{x^2 + y^2 \leq 1, |x| < |y|} dx dy \quad (79)$$

$$= \frac{1}{\pi} \cdot \text{area}(x^2 + y^2 \leq 1, |x| < |y|) \quad (80)$$

$$= \frac{1}{\pi} \cdot \frac{\pi}{2} = \frac{1}{2} \quad (81)$$

□

(b) (4 points) What is the marginal density of the vertical component, $f_Y(y)$

solution.

$$f_Y(y) = \int_{-\sqrt{1-y^2}}^{\sqrt{1-y^2}} f_{X,Y}(x,y) dx \quad (82)$$

$$= \frac{2}{\pi} \sqrt{1-y^2}, \quad y \in (-1, 1) \quad (83)$$

□

2. Suppose that at the start of the week, the amount of gasoline (measured in 1K gallons) in the holding tank of our neighborhood gas station can be represented by a random variable, G , with a uniform distribution over $[0, 1]$ (so the maximum possible is 1000 gallons of gasoline at the start of the week). Further, suppose that the amount of gasoline sold by the gas station during the week can also be modeled with uniform random variable, Y , over the interval $[0, g_1]$, where g_1 is the particular amount at the start of the week.

(a) (5 points) Find the joint density function for the amount at the start of the week and the amount sold during the week.

solution. Now we know that $G \sim U(0, 1)$, $Y|_{G=g_1} \sim U(0, g_1)$ so

$$f_{Y,G}(y, g) = f_G(g) \cdot f_{Y|G}(y|g) = \frac{1}{g}, \quad 0 < y < g < 1 \quad (84)$$

□

(b) (5 points) If the station stocks 600 gallons at the start of the week, what is the probability they sell more than 250 gallons?

solution. Condition on $G = 0.6$, $Y \sim U(0, 0.6)$, so

$$\mathbb{P}(Y \geq 0.25 | G = 0.6) = \int_{0.25}^{0.6} \frac{1}{0.6} dy = \frac{0.35}{0.6} = \frac{7}{12} \quad (85)$$

□

(c) (5 points) Find the marginal density of Y .

solution.

$$f_Y(y) = \int_y^1 f_{Y,G}(y, g) dg = \int_y^1 \frac{1}{g} dg = -\log y, \quad y \in (0, 1) \quad (86)$$

□

(d) (5 points) If we know that the station sold 250 gallons of gasoline, what is the probability that they had more than 500 gallons at the start of the week?

solution. In order to get $\mathbb{P}(G \geq 0.5 | Y = 0.25)$, let's first derive the conditional density $f_{G|Y}(g|y)$

$$f_{G|Y}(g|y) = \frac{f_{Y,G}(y, g)}{f_Y(y)} = \frac{1}{-\log y \cdot g}, \quad 0 < y < g < 1 \quad (87)$$

so now

$$\mathbb{P}(G \geq 0.5 | Y = 0.25) = \int_{0.5}^1 f_{G|Y}(g|0.25) dg \quad (88)$$

$$= \int_{0.5}^1 \frac{1}{-\log 0.25 \cdot g} dg \quad (89)$$

$$= \frac{\log 0.5}{\log 0.25} = \frac{1}{2} \quad (90)$$

□

(e) (5 points) What is the expected amount of gasoline we will have left at the end of the week?

solution.

$$\mathbb{E}(G - Y) = \mathbb{E}G - \mathbb{E}Y \quad (91)$$

where $\mathbb{E}G = \frac{1}{2}$ and $\mathbb{E}Y = \mathbb{E}[\mathbb{E}(Y|G)] = \mathbb{E}\frac{G}{2} = \frac{1}{4}$, so $\mathbb{E}(G - Y) = \frac{1}{4}$. \square

3. The hens at Lilly's Cage-free Egg Farm produce eggs according to a Poisson distribution, with a mean of β per day. Because the hens are allowed to roam free, and not kept in cages, they lay their eggs around the barnyard which increases the chance that they will break. The eggs survive to be collected with some probability p .

(a) (5 points) Find the expected number of eggs that will survive to be collected.

solution. Assume there are $N \sim P(\beta)$ eggs laid and C eggs collected. So $C|_{N=n} \sim B(n, p)$.

$$\mathbb{E}C = \mathbb{E}[\mathbb{E}(C|N)] = \mathbb{E}pN = p\mathbb{E}N = p\beta \quad (92)$$

\square

(b) (5 points) Find the variance for the number of eggs that will survive to be collected.

solution.

$$\mathbb{E}C^2 = \mathbb{E}[\mathbb{E}(C^2|N)] \quad (93)$$

now $\mathbb{E}(C^2|N) = \text{Var}(C|N) + [\mathbb{E}(C|N)]^2 = Np(1-p) + N^2p^2$ so

$$\mathbb{E}C^2 = \mathbb{E}[\mathbb{E}(C^2|N)] = \mathbb{E}Np(1-p) + \mathbb{E}N^2p^2 \quad (94)$$

$$= p(1-p)\mathbb{E}N + p^2\mathbb{E}N^2 \quad (95)$$

$$= p(1-p)\beta + p^2(\beta^2 + \beta) \quad (96)$$

since $\mathbb{E}N^2 = \text{Var}(N) + (\mathbb{E}N)^2 = \beta^2 + \beta$, so

$$\text{Var}(C) = \mathbb{E}C^2 - (\mathbb{E}C)^2 = p(1-p)\beta + p^2(\beta^2 + \beta) - p^2\beta^2 \quad (97)$$

$$= p(1-p)\beta + p^2\beta = p\beta \quad (98)$$

\square

Central Limit Theorem

CLT only works for *i.i.d. random variable series that has finite second moment*. The limiting distribution of $\frac{S_n - \mathbb{E}S_n}{\sqrt{\text{Var}(S_n)}}$ is always $N(0, 1)$. When one wants to apply CLT, first verify that the random variable series is *i.i.d.*,

then verify that second moment exists. If those two conditions hold, one only need to calculate the expectation and the variance of the sum $S_n = X_1 + \dots + X_n$ to write out the normal approximation.

For example, if we have $X_1, \dots, X_n \sim P(\lambda)$ *i.i.d.*, since $Var(X_1) = \lambda < \infty$, CLT holds and we compute

$$\mathbb{E}S_n = n\mathbb{E}X_1 = n\lambda \quad (99)$$

$$Var(S_n) = nVar(X_1) = n\lambda \quad (100)$$

to conclude that

$$\frac{S_n - n\lambda}{\sqrt{n\lambda}} \xrightarrow{d} N(0, 1) \quad (n \rightarrow \infty) \quad (101)$$

which means that when n is large enough, the following approximation that

$$\mathbb{P}\left(\frac{S_n - n\lambda}{\sqrt{n\lambda}} \leq x\right) \rightarrow \Phi(x) \quad (n \rightarrow \infty) \quad (102)$$

works for standard Gaussian CDF Φ . If we simplify the expression, we can see that $\sqrt{n}(\bar{X}_n - \lambda) \xrightarrow{d} N(0, \lambda)$.

Now if we have paired observations $X_1, \dots, X_n, Y_1, \dots, Y_n$ to be independent and X_i all have the same distribution $N(\mu_1, \sigma_1^2)$, Y_i all have the same distribution $N(\mu_2, \sigma_2^2)$, then we can consider $S_n = (X_1 - Y_1) + \dots + (X_n - Y_n)$ to see that $X_1 - Y_1, \dots, X_n - Y_n$ are *i.i.d.* with $Var(X_1 - Y_1) = \sigma_1^2 + \sigma_2^2 < \infty$ so CLT holds. Compute

$$\mathbb{E}S_n = n(\mu_1 - \mu_2) \quad (103)$$

$$Var(S_n) = n(\sigma_1^2 + \sigma_2^2) \quad (104)$$

to get the conclusion that

$$\frac{S_n - n(\mu_1 - \mu_2)}{\sqrt{n(\sigma_1^2 + \sigma_2^2)}} \xrightarrow{d} N(0, 1) \quad (n \rightarrow \infty) \quad (105)$$

divide numerator and denominator by n to see

$$\frac{\sqrt{n}[(\bar{X}_n - \bar{Y}_n) - (\mu_1 - \mu_2)]}{\sqrt{(\sigma_1^2 + \sigma_2^2)}} \xrightarrow{d} N(0, 1) \quad (n \rightarrow \infty) \quad (106)$$

Remark. The simplest CLT does not hold for $X_1, \dots, X_n, Y_1, \dots, Y_n$ since they are independent but not identically distributed. However, if we notice the fact that they are paired samples, we can consider the difference $X_i - Y_i$ as a new random variable series so now it's *i.i.d.*.

Week 5

Bias Variance Decomposition

For any estimator $\hat{\theta} = \hat{\theta}(X)$, $X = (X_1, \dots, X_n)$ that estimates the true parameter θ , we hope to set up a criterion for selecting the best estimator. A frequently used criterion is the mean square error

$$MSE(\hat{\theta}) = \mathbb{E}(\hat{\theta} - \theta)^2 \quad (107)$$

and it's important to realize that the mean square error always has the bias variance decomposition

$$MSE(\hat{\theta}) = \mathbb{E}(\hat{\theta} - \theta)^2 \quad (108)$$

$$= \mathbb{E}(\hat{\theta} - \mathbb{E}\hat{\theta} + \mathbb{E}\hat{\theta} - \theta)^2 \quad (109)$$

$$= \mathbb{E}(\hat{\theta} - \mathbb{E}\hat{\theta})^2 + \mathbb{E}(\mathbb{E}\hat{\theta} - \theta)^2 + 2\mathbb{E}[(\hat{\theta} - \mathbb{E}\hat{\theta})(\mathbb{E}\hat{\theta} - \theta)] \quad (110)$$

$$= Var(\hat{\theta}) + Bias^2(\hat{\theta}) \quad (111)$$

since $\mathbb{E}[(\hat{\theta} - \mathbb{E}\hat{\theta})(\mathbb{E}\hat{\theta} - \theta)] = (\mathbb{E}\hat{\theta} - \theta)(\mathbb{E}\hat{\theta} - \mathbb{E}\hat{\theta}) = 0$. This is showing some kind of **trade-off between unbiasedness and efficiency under a fixed MSE**. If one wants to find an unbiased estimator, one always has to sacrifice some kind of efficiency, while if one wants to find an efficient estimator, it may be seriously biased (the trivial estimator $\hat{\theta} = 0$ always has zero variance but may be significantly biased).

Let's see an example in the textbook exercise 8.20 where $Y_1, \dots, Y_4 \sim \mathcal{E}(\frac{1}{\theta})$ and we want to estimate the parameter θ . Now $X = \sqrt{Y_1 Y_2}$ and we want to find a multiple of X which is unbiased.

Let's first calculate $\mathbb{E}X = \mathbb{E}^2\sqrt{Y_1}$. Note that

$$\mathbb{E}\sqrt{Y_1} = \int_0^\infty \sqrt{y} \frac{1}{\theta} e^{-\frac{y}{\theta}} dy \quad (112)$$

$$= \sqrt{\theta} \int_0^\infty \sqrt{u} e^{-u} du \quad \left(u = \frac{y}{\theta}\right) \quad (113)$$

$$= \sqrt{\theta} \cdot \Gamma\left(\frac{3}{2}\right) \quad (114)$$

$$= \sqrt{\theta} \cdot \frac{1}{2} \Gamma\left(\frac{1}{2}\right) \quad (115)$$

$$= \frac{\sqrt{\theta\pi}}{2} \quad (116)$$

as a result, $\mathbb{E}X = \frac{\theta\pi}{4}$ and $\mathbb{E}\frac{4}{\pi}\sqrt{Y_1 Y_2} = \theta$ so $\frac{4}{\pi}\sqrt{Y_1 Y_2}$ is unbiased.

Remark. Recall from calculus the form of Gamma function that

$$\Gamma(s) = \int_0^\infty x^{s-1} e^{-x} dx \quad (s > 0) \quad (117)$$

and the property that $\forall p \in (0, 1), \Gamma(p)\Gamma(1-p) = \frac{\pi}{\sin p\pi}$ so $\Gamma(\frac{1}{2}) = \sqrt{\pi}$ (or directly from the hint).

For another example, refer to exercise 8.36 in the textbook. Y_1, \dots, Y_n follow $\mathcal{E}(\frac{1}{\theta})$ so we now that

$$\mathbb{E}Y_1 = \theta, \text{Var}(Y_1) = \theta^2, \mathbb{E}\bar{Y} = \theta, \text{Var}(\bar{Y}) = \frac{\theta^2}{n} \quad (118)$$

now we want to construct unbiased estimator for θ and provide an estimate for the standard error of this estimator.

Now we want the expectation of the estimator to be θ so it's quite obvious that Y_1 is just an unbiased estimator. To find the standard error

$$se(Y_1) = \theta \quad (119)$$

since θ itself is unknown, it can be estimated in various ways. We can take \bar{Y} as an unbiased estimation of θ so the estimated standard error is

$$\hat{se}(Y_1) = \bar{Y} \quad (120)$$

we can also take sample variance S^2 as an estimation of θ^2 so the estimated standard error is

$$\hat{se}(Y_1) = \sqrt{S^2} \quad (121)$$

there's no fixed way to estimate the standard error so you can put up any reasonable ways to do it! On the other hand, the choice of unbiased estimator is also not unique, Y_1 is unbiased, but \bar{Y} is also unbiased, one may also choose the sample mean as the operator and try to calculate the estimated standard error.

Consistency

Estimator $\hat{\theta}$ is called a consistent estimator if

$$\hat{\theta} \xrightarrow{P} \theta \quad (n \rightarrow \infty) \quad (122)$$

which means that $\forall \varepsilon > 0, \mathbb{P}(|\hat{\theta} - \theta| \geq \varepsilon) \rightarrow 0 \quad (n \rightarrow \infty)$, this convergence is called convergence in probability.

Since convergence in probability is preserved under addition, subtraction, multiplication, division and actually any continuous mapping (continuous mapping theorem), consistency is typically easy to get. In other words, if $\hat{\theta}_1, \hat{\theta}_2$ are consistent estimators of θ , then $\frac{\hat{\theta}_1 + \hat{\theta}_2}{2}$ is still a consistent estimator.

By what we have learnt so far, sample variance S^2 is an unbiased consistent estimator of population variance σ^2 , so S is a **consistent estimator of σ** (apply the continuous function $g(x) = \sqrt{x}$ on both sides). However, note that S is **not an unbiased estimator of σ** ! (since generally $\mathbb{E}\sqrt{X} \neq \sqrt{\mathbb{E}X}$)

We raise an example next to illustrate the way to construct a consistent estimator and to prove its consistency. For sample Y_1, \dots, Y_n with $\mathbb{E}Y_1 = \mu, \text{Var}(Y_1) = \sigma^2$, we want to estimate $\mathbb{E}Y_1^2$ consistently.

A natural estimator comes from the sample mean of second moments

$$T = \frac{\sum_{i=1}^n Y_i^2}{n} \quad (123)$$

and to show its consistency, it directly follows from WLLN that since $\mathbb{E}Y_1^2 < \infty$

$$T \xrightarrow{p} \mathbb{E}Y_1^2 \quad (n \rightarrow \infty) \quad (124)$$

Efficiency

We say an estimator is more efficient than the other estimator if it has lower variance with the relative efficiency defined as

$$\text{eff}(\hat{\theta}_1, \hat{\theta}_2) = \frac{\text{Var}(\hat{\theta}_2)}{\text{Var}(\hat{\theta}_1)} \quad (125)$$

in this course we just need to learn to calculate the relative efficiency between two given estimators.

For example, if Y_1, \dots, Y_n are from a population with mean μ and variance σ^2 , if we have $\hat{\mu}_2 = \frac{1}{4}Y_1 + \frac{Y_2 + \dots + Y_{n-1}}{2(n-2)} + \frac{1}{4}Y_n$ and $\hat{\mu}_3 = \bar{Y}$, then

$$\text{Var}(\hat{\mu}_3) = \frac{\text{Var}(Y_1)}{n} = \frac{\sigma^2}{n} \quad (126)$$

$$\text{Var}(\hat{\mu}_2) = \frac{1}{16}\sigma^2 + \frac{n-2}{4(n-2)^2}\sigma^2 + \frac{1}{16}\sigma^2 = \frac{1}{8}\sigma^2 + \frac{1}{4(n-2)}\sigma^2 \quad (127)$$

so

$$\text{eff}(\hat{\mu}_3, \hat{\mu}_2) = \frac{\text{Var}(\hat{\mu}_2)}{\text{Var}(\hat{\mu}_3)} = \frac{n}{8} + \frac{n}{4(n-2)} \quad (128)$$

Week 6

Quiz 2

1.

Suppose we can model n wait times for the 24X (bus) as Y_1, Y_2, \dots, Y_n *i.i.d.* Exponential random variables with mean β . We want to assess how well the Central Limit Theorem approximates the exponential distribution.

(a) (5 points) First, let's consider the standardized sample mean $\frac{\sqrt{n}(\bar{Y}-\beta)}{\beta}$ as a random variable, write an expression, denoted $\tilde{F}(c)$, to approximate the CDF $F(c) = \mathbb{P}\left(\frac{\sqrt{n}(\bar{Y}-\beta)}{\beta} \leq c\right)$ and calculate $\tilde{F}(2.65)$.

solution. Since Y_1, \dots, Y_n are *i.i.d.* with $\text{Var}(Y_1) = \beta^2 < \infty$, we can apply the CLT. Since $\mathbb{E}\bar{Y} = \mathbb{E}Y_1 = \beta$, $\text{Var}(\bar{Y}) = \frac{\text{Var}(Y_1)}{n} = \frac{\beta^2}{n}$, we know that

$$\frac{\sqrt{n}(\bar{Y} - \beta)}{\beta} \xrightarrow{d} N(0, 1) \quad (n \rightarrow \infty) \quad (129)$$

and the approximation holds $\tilde{F}(c) = \Phi(c)$ where Φ is standard Gaussian CDF.

$$\tilde{F}(2.65) = \Phi(2.65) = 0.99598.$$

□

(b) (5 points) Find c^* in terms of c and β such that $F(c) = \mathbb{P}(n\bar{Y} \leq c^*)$. How can you calculate $F(c)$ exactly? (Hint: what is the sampling distribution of $n\bar{Y}$?)

solution.

$$F(c) = \mathbb{P}\left(\frac{\sqrt{n}(\bar{Y} - \beta)}{\beta} \leq c\right) = \mathbb{P}\left(\bar{Y} \leq \beta + \frac{c\beta}{\sqrt{n}}\right) = \mathbb{P}\left(n\bar{Y} \leq n\beta\left(1 + \frac{c}{\sqrt{n}}\right)\right) \quad (130)$$

so $c^* = n\beta\left(1 + \frac{c}{\sqrt{n}}\right)$.

Notice that $Y_1 \sim \mathcal{E}\left(\frac{1}{\beta}\right) = \Gamma\left(1, \frac{1}{\beta}\right)$ and if $Y_1, \dots, Y_n \sim \Gamma\left(1, \frac{1}{\beta}\right)$ *i.i.d.*, then $Y_1 + \dots + Y_n \sim \Gamma\left(n, \frac{1}{\beta}\right)$. (the additivity of Gamma distribution w.r.t. parameter α , one may prove it using MGF)

$$F(c) = F_G(c^*) \quad (131)$$

where F_G is the CDF of distribution $\Gamma\left(n, \frac{1}{\beta}\right)$.

□

(c) (3 points) For each of the values of n in the table below, calculate the exact probability, $F(2.65)$ and fill in the first empty column. Then, find the distance (absolute difference) between that value and the approximation you found above. Calculate for $n = 3, 9, 36, 121, 169$ and $\beta = 1$.

solution. Now

$$F(c) = F_G(c^*) \quad (132)$$

where F_G is the CDF of distribution $\Gamma(n, 1)$ and $c = 2.65$ is fixed, so $c^* = n \left(1 + \frac{2.65}{\sqrt{n}}\right)$.

When $n = 3$, $c^* = 3 \left(1 + \frac{2.65}{\sqrt{3}}\right) = 7.59$, so $F(c) = F_G(7.59) = 0.9811$.

When $n = 9$, $c^* = 9 \left(1 + \frac{2.65}{\sqrt{9}}\right) = 16.95$, so $F(c) = F_G(16.95) = 0.98704$.

When $n = 36$, $c^* = 36 \left(1 + \frac{2.65}{\sqrt{36}}\right) = 51.9$, so $F(c) = F_G(51.9) = 0.9916$.

When $n = 121$, $c^* = 121 \left(1 + \frac{2.65}{\sqrt{121}}\right) = 150.15$, so $F(c) = F_G(150.15) = 0.99366$.

When $n = 169$, $c^* = 169 \left(1 + \frac{2.65}{\sqrt{169}}\right) = 203.45$, so $F(c) = F_G(203.45) = 0.99403$.

So the difference $|F(2.65) - \tilde{F}(2.65)|$ changes like 0.01488, 0.00894, 0.00438, 0.002332, 0.00195. \square

(d) (2 points) I've said in lecture that the central limit theorem approximation is good when n is "large enough". As we increased the sample size, what did you notice about the difference between the approximation and the exact probability? How large of an n value would you think is sufficient to be "good" and why?

solution. By setting the error tolerance limit as 1%, $n = 9$ suffices. (There's no unique criteria, just set up an appropriate one you prefer) \square

2. (10 points) The client manager at the law firm of Dewey, Cheatem and Howe sends out bids to their larger clients for particular legal needs. (E.g., preparing the paperwork for a merger.) Because these bids are worded as "Price not to exceed...", they need to be careful about their estimation for the hours needed to fulfill the particular task so they don't lose money. Suppose we let X be the amount of paralegal hours required for a project, and suppose further, that those hours are normally distributed with a mean of 16 hours and a standard deviation of 4 hours. Now, suppose that Y is the amount of attorney hours required for a job, and that the hours are normally distributed with a mean of 10 and a standard deviation of 6 hours. The hourly rate for the paralegals is \$85 while the hourly rate for the attorneys is \$175. We also must include an overhead charge equal to \$500, to all jobs we bid. Then we can write the cost equation as, $cost = 85X + 175Y + 500$.

(a) (8 points) If X and Y are independent, how much should the client manager bid so that the probability of losing money is 0.05 (we lose money when the costs exceed the amount of the bid).

solution. B denotes the amount of the bid, $X \sim N(16, 16)$, $Y \sim N(10, 36)$, want to find B such that

$$0.05 = \mathbb{P}(85X + 175Y + 500 \geq B) \quad (133)$$

notice that the linear combination of two independent Gaussian random variables $85X + 175Y$ is still Gaussian, so

calculate its expectation and variance to see

$$\mathbb{E}(85X + 175Y) = 85 \times 16 + 175 \times 10 = 3110 \quad (134)$$

and the variance

$$\text{Var}(85X + 175Y) = 85^2 \times 16 + 175^2 \times 36 = 1218100 \quad (135)$$

so $85X + 175Y \sim N(3110, 1218100)$.

$$\mathbb{P}(85X + 175Y + 500 \geq B) = \mathbb{P}(85X + 175Y \geq B - 500) \quad (136)$$

$$= \mathbb{P}\left(\frac{85X + 175Y - 3110}{\sqrt{1218100}} \geq \frac{B - 500 - 3110}{\sqrt{1218100}}\right) \quad (137)$$

$$= 1 - \Phi\left(\frac{B - 500 - 3110}{\sqrt{1218100}}\right) = 0.05 \quad (138)$$

so now $\frac{B-500-3110}{\sqrt{1218100}} = 1.65$ and $B = 5431.065$.

□

(b) (2 points) Does the independence of the paralegal hours and attorney hours seem reasonable? Why or why not? (No calculations are necessary, just explain what your intuition tells you.)

solution. Just make reasonable explanations.

□

3. (15 points) Imagine X_1, X_2, \dots, X_n are a random sample (of size n) of McConnell's ice cream scoop weights, and they are normally distributed with a mean μ_X and variance σ_X^2 . Now, let Y_1, Y_2, \dots, Y_m be a random sample (of size m) of Rory's ice cream scoop weights, and they are also normally distributed, but with a mean of μ_Y and a variance of σ_Y^2 .

(a) (5 points) Find the expected value of the difference between the two sample means.

solution.

$$\mathbb{E}(\bar{X} - \bar{Y}) = \mathbb{E}X_1 - \mathbb{E}Y_1 = \mu_X - \mu_Y \quad (139)$$

□

(b) (5 points) Find the variance of the difference between the two sample means.

solution. By the independence between two samples,

$$\text{Var}(\bar{X} - \bar{Y}) = \text{Var}(\bar{X}) + \text{Var}(\bar{Y}) \quad (140)$$

$$= \frac{\text{Var}(X_1)}{n} + \frac{\text{Var}(Y_1)}{m} = \frac{\sigma_X^2}{n} + \frac{\sigma_Y^2}{m} \quad (141)$$

□

(c) (5 points) Now, suppose that $\sigma_X^2 = 1.21(\text{grams}^2)$ and $\sigma_Y^2 = 0.81(\text{grams}^2)$, and that we have equal sized samples. Find the sample sizes so that the difference in sample means will be within 1 gram of the difference in population means, with a probability of 0.90.

solution. Now $m = n$ and we want to find n such that

$$\mathbb{P}(|(\bar{X} - \bar{Y}) - (\mu_X - \mu_Y)| \leq 1) = 0.9 \quad (142)$$

now notice that $\bar{X} - \bar{Y} \sim N\left(\mu_X - \mu_Y, \frac{\sigma_X^2 + \sigma_Y^2}{n}\right)$ is Gaussian so

$$\mathbb{P}(|(\bar{X} - \bar{Y}) - (\mu_X - \mu_Y)| \leq 1) = \mathbb{P}(-1 \leq (\bar{X} - \bar{Y}) - (\mu_X - \mu_Y) \leq 1) \quad (143)$$

$$= \mathbb{P}\left(-\sqrt{\frac{n}{\sigma_X^2 + \sigma_Y^2}} \leq \frac{(\bar{X} - \bar{Y}) - (\mu_X - \mu_Y)}{\sqrt{\frac{\sigma_X^2 + \sigma_Y^2}{n}}} \leq \sqrt{\frac{n}{\sigma_X^2 + \sigma_Y^2}}\right) \quad (144)$$

$$= \Phi\left(\sqrt{\frac{n}{\sigma_X^2 + \sigma_Y^2}}\right) - \Phi\left(-\sqrt{\frac{n}{\sigma_X^2 + \sigma_Y^2}}\right) \quad (145)$$

$$= 2\Phi\left(\sqrt{\frac{n}{\sigma_X^2 + \sigma_Y^2}}\right) - 1 = 0.9 \quad (146)$$

so $\sqrt{\frac{n}{\sigma_X^2 + \sigma_Y^2}} = \Phi^{-1}(0.95) = 1.65$ and $n = 1.65^2(\sigma_X^2 + \sigma_Y^2) = 5.5$. So we need at least 6 samples.

□

Fisher Information and Cramer-Rao Bound

Those two concepts appear since we want to find the minimum variance unbiased estimator (MVUE).

Consider $f(y; \theta)$ as the joint likelihood, then the **log-likelihood** is $l(y; \theta) = \log f(y; \theta)$ and the **score function** is defined as $S(\theta; y) = \frac{d}{d\theta} l(y; \theta)$ is the changing rate of the log-likelihood. So the variance of the score function is defined as the **Fisher information**

$$I(\theta) = \text{Var}(S(\theta; Y)) = \text{Var}\left(\frac{d}{d\theta} l(Y; \theta)\right) = \mathbb{E}\left(\frac{d}{d\theta} l(Y; \theta)\right)^2 \quad (147)$$

since $\mathbb{E}\left(\frac{d}{d\theta} l(Y; \theta)\right) = 0$. Another formulation of Fisher information is given by

$$-\mathbb{E} \frac{d^2}{d\theta^2} l(Y; \theta) = I(\theta) \quad (148)$$

so we just need to derive log-likelihood, take second-order derivative w.r.t. θ , take expectation and add negative sign to get the Fisher information. Such formula for Fisher information also tells us that if we have Y_1, \dots, Y_n as *i.i.d.*

samples, to calculate the Fisher information

$$I(\theta) = -\mathbb{E} \frac{d^2}{d\theta^2} l(Y_1, \dots, Y_n; \theta) \quad (149)$$

$$= -\mathbb{E} \frac{d^2}{d\theta^2} \sum_{i=1}^n l(Y_i; \theta) \quad (150)$$

$$= \sum_{i=1}^n -\mathbb{E} \frac{d^2}{d\theta^2} l(Y_i; \theta) \quad (151)$$

$$= \sum_{i=1}^n I_i(\theta) = nI_1(\theta) \quad (152)$$

since the Fisher information on seeing each sample is the same. This shows that for *i.i.d.* samples we just need to calculate the Fisher information of one of them and multiply by sample size.

The Cramer-Rao bound connects the variance of the estimator with Fisher information

$$\text{Var}(\hat{\theta}) \geq \frac{\left(\frac{d}{d\theta} \mathbb{E} \hat{\theta} \right)^2}{I(\theta)} \quad (153)$$

so if $\hat{\theta}$ is an unbiased estimator of θ , then $\text{Var}(\hat{\theta}) \geq \frac{1}{I(\theta)}$. So any unbiased estimator that can attain this bound is called MVUE.

Example

Let's consider the Fisher information of exponential distribution with likelihood $f(y; \theta) = \frac{1}{\theta} e^{-\frac{y}{\theta}}$ so the log-likelihood is

$$l(y; \theta) = -\log \theta - \frac{y}{\theta} \quad (154)$$

and take second derivative

$$\frac{d^2}{d\theta^2} l(y; \theta) = \frac{d}{d\theta} \left(-\frac{1}{\theta} + \frac{y}{\theta^2} \right) = \frac{1}{\theta^2} - \frac{2y}{\theta^3} \quad (155)$$

now compute Fisher information for single sample

$$I_1(\theta) = \mathbb{E} \left(-\frac{1}{\theta^2} + \frac{2Y}{\theta^3} \right) \quad (156)$$

$$= -\frac{1}{\theta^2} + \frac{2}{\theta^2} = \frac{1}{\theta^2} \quad (157)$$

so the Fisher information for Y_1, \dots, Y_n is $I(\theta) = \frac{n}{\theta^2}$ and the Cramer-Rao bound for unbiased estimator $\hat{\theta}$ is given by

$$\text{Var}(\hat{\theta}) \geq \frac{1}{I(\theta)} = \frac{\theta^2}{n} \quad (158)$$

we will see that since the sample mean $\text{Var}(\bar{Y}) = \frac{\theta^2}{n}$, it's MVUE.

Another example for Pareto distribution with likelihood $f(y; \theta) = \frac{2\theta^2}{y^3}$ ($y > \theta$), log-likelihood is

$$l(y; \theta) = \log 2 + 2 \log \theta - 3 \log y \quad (159)$$

and take second derivative

$$\frac{d^2}{d\theta^2} l(y; \theta) = \frac{d}{d\theta} \frac{2}{\theta} = -\frac{2}{\theta^2} \quad (160)$$

now compute Fisher information for single sample

$$I_1(\theta) = \mathbb{E} \frac{2}{\theta^2} \quad (161)$$

$$= \frac{2}{\theta^2} \quad (162)$$

so the Fisher information for Y_1, \dots, Y_n is $I(\theta) = \frac{2n}{\theta^2}$ and the Cramer-Rao bound for unbiased estimator $\hat{\theta}$ is given by

$$\text{Var}(\hat{\theta}) \geq \frac{1}{I(\theta)} = \frac{\theta^2}{2n} \quad (163)$$

Sufficiency

A statistic $T = T(X_1, \dots, X_n)$ is **sufficient** if $X|_T$ does not depend on parameter θ , which means that on knowing the information contained in statistic T , the original parameter is no longer useful for figuring out the distribution of the sample, so one can always reduce the sample to the sufficient statistic if the parametric model is known. From this definition, we develop the important factorization theorem that says T is sufficient if and only if

$$f(x; \theta) = g_\theta(T(x)) \cdot h(x) \quad (164)$$

the joint likelihood can be decomposed into the product of g that only contains $\theta, T(x)$ and h that only contains x . As a result, if a function contains x, θ at the same time and cannot be decomposed further into products, it has to be contained in $g_\theta(T(x))$.

Remark. Be careful with the support of random variables when figuring out sufficiency, some random variables' support depends on the parameter!

Example

Consider $X_1, \dots, X_n \sim U(0, \theta)$ so

$$f(x; \theta) = \frac{1}{\theta^n} \mathbb{I}_{0 < x_1, \dots, x_n < \theta} \quad (165)$$

$$= \frac{1}{\theta^n} \mathbb{I}_{\min\{x_1, \dots, x_n\} > 0, \max\{x_1, \dots, x_n\} < \theta} \quad (166)$$

$$= g_\theta(T(x)) \cdot h(x) \quad (167)$$

with $h(x) = \mathbb{I}_{\min\{x_1, \dots, x_n\} > 0}$, $g_\theta(T(x)) = \frac{1}{\theta^n} \mathbb{I}_{\max\{x_1, \dots, x_n\} < \theta}$, so naturally we see that $T(X) = \max\{X_1, \dots, X_n\}$ is the sufficient statistic for this distribution.

Week 7

Minimal Sufficiency

A **minimal sufficient statistic** is a sufficient statistic which is the function of any other sufficient statistics. Note that sufficient statistic already contains all the information of the sample distribution without replying on the value of parameter, so the minimal sufficient statistic is actually the sufficient statistic that contains the "minimal amount of information".

Remark. To understand the definition, for any random sample observation, let's say 1,2,3. A function of this sample (let's say the sum) can always be calculated, which is $1+2+3=6$. However, on knowing the sum of the samples to be 6, one can not figure out which three realizations the sample gives, it can be 1,2,3, it can also be 2,2,2, etc. This explains the fact that a function of a statistic cannot contain more information than that statistic itself.

The formal definition of minimal sufficient statistic T' is that it's sufficient and for any other sufficient statistic T , it's always true that if $T(X) = T(Y)$ then $T'(X) = T'(Y)$. The minimal sufficient statistic is always found by the following theorem.

Theorem 1. (Lehmann-Scheffe) T is a minimal sufficient statistic if and only if $\frac{L(x_1, \dots, x_n; \theta)}{L(y_1, \dots, y_n; \theta)}$ is independent of θ is equivalent to saying $T(x_1, \dots, x_n) = T(y_1, \dots, y_n)$.

Let's illustrate the way to apply this method in practice by looking at examples. Let's consider $X_1, \dots, X_n \sim P(\lambda)$ so the likelihood ratio is

$$\frac{\lambda^{\sum_i x_i} e^{-n\lambda}}{\lambda^{\sum_i y_i} e^{-n\lambda}} = \lambda^{\sum_i x_i - \sum_i y_i} \quad (168)$$

is independent of λ if and only if $\sum_i x_i - \sum_i y_i = 0$, so by taking $T(X) = \sum_i X_i$, it must be minimal sufficient.

Let's consider another example $X_1, \dots, X_n \sim U(\theta - 1, \theta + 1)$ so the likelihood ratio is

$$\frac{\frac{1}{2^n} \mathbb{I}_{x_1, \dots, x_n \in (\theta-1, \theta+1)}}{\frac{1}{2^n} \mathbb{I}_{y_1, \dots, y_n \in (\theta-1, \theta+1)}} = \frac{\mathbb{I}_{\min\{x_1, \dots, x_n\} > \theta-1} \mathbb{I}_{\max\{x_1, \dots, x_n\} < \theta+1}}{\mathbb{I}_{\min\{y_1, \dots, y_n\} > \theta-1} \mathbb{I}_{\max\{y_1, \dots, y_n\} < \theta+1}} \quad (169)$$

is independent of θ if and only if $\min\{x_1, \dots, x_n\} = \min\{y_1, \dots, y_n\}$, $\max\{x_1, \dots, x_n\} = \max\{y_1, \dots, y_n\}$. So $T(X) = (\min\{X_1, \dots, X_n\}, \max\{X_1, \dots, X_n\})$ is minimal sufficient.

Moment Estimator

Moment estimator is simply calculating population moments and match them with sample moments to get the estimators of the parameters. However, some distributions do not even have the first moment (Cauchy), so the application of moment estimator is actually restricted.

To illustrate the moment estimator, consider examples $X_1, \dots, X_n \sim P(\lambda)$, the population first moment is λ and the sample first moment is \bar{X} so $\hat{\lambda} = \bar{X}$ is the moment estimator.

Consider $X_1, \dots, X_n \sim U(0, \theta)$, the population first moment is $\frac{\theta}{2}$ and the sample first moment is \bar{X} so $\hat{\theta} = 2\bar{X}$. However if $X_1, \dots, X_n \sim U(-\theta, \theta)$, then population mean is always 0 so there's no way to match 0 with the sample mean \bar{X} . In this situation, we need to consider the second moment, population second moment is $\frac{4\theta^2}{12} = \frac{\theta^2}{3}$ and sample second moment is $\frac{\sum_i X_i^2}{n}$, so $\hat{\theta} = \sqrt{\frac{3}{n} \sum_{i=1}^n X_i^2}$. As we can see, even if the moments exist, there would be problems matching population and sample moments.

On the other hand, the property of moment estimators can be very bad. As we can see, the estimator above for $U(-\theta, \theta)$ is not unbiased and it enjoys no asymptotic properties.

Maximum Likelihood Estimator

MLE applies for much more general case than the moment estimator since it deals with likelihood which always exists. It picks the estimator that maximizes the joint likelihood.

For example, for $X_1, \dots, X_n \sim P(\lambda)$, we have

$$L(\lambda) = \lambda^{\sum_i x_i} e^{-n\lambda} \quad (170)$$

so $l(\lambda) = \sum_i x_i \log \lambda - n\lambda$, taking derivative gives $l'(\lambda) = \frac{\sum_i x_i}{\lambda} - n$, so $\hat{\lambda} = \bar{X}$ gives the MLE.

For example, for $X_1, \dots, X_n \sim U(-\theta, \theta)$, we have

$$L(\theta) = \frac{1}{(2\theta)^n} \mathbb{I}_{\min\{x_1, \dots, x_n\} > -\theta, \max\{x_1, \dots, x_n\} < \theta} \quad (171)$$

so $l(\theta) = -n \log(2\theta)$, $\max\{x_1, \dots, x_n\} < \theta$, $-\min\{x_1, \dots, x_n\} < \theta$, since the log-likelihood is monotone decreasing in θ , to make it larger, θ has to be as small as possible. Notice the range of θ and it's natural that

$$\hat{\theta} = \max\{\max\{x_1, \dots, x_n\}, -\min\{x_1, \dots, x_n\}\} \quad (172)$$

the direct interpretation is that the MLE of θ is the smallest positive θ such that $\forall i, |x_i| < \theta$ shall always be true.

Notice some properties of MLE that the MLE of the function of parameter is always the function of the MLE of the parameters and that MLE enjoys asymptotic normality for large sample size.