

WEEK 11 EXERCISES

You are expected to do all these problems, but for **Homework 10** please turn in **only Problems 2, 5, and 6** on **Thursday November 8 at the start of lecture**.

1. Baseline and Post-Treatment

Weight loss is an undesirable side effect of a new medical treatment. Patients are weighed, in pounds, at the start and end of the treatment. Let X be the baseline (pre-treatment) weight of a randomly picked patient and let Y be that patient's post-treatment weight. Suppose X and Y have the bivariate normal distribution with parameters $(160, 140, 15, 20, 0.6)$.

(a) Find $P((X + Y)/2 > 150)$. Did you use the bivariate normal assumption?

$$\mathbb{P}((X + Y)/2 > 150) = \mathbb{P}((X + Y)/2 - 150 > 0) = 0.5,$$

since $(X + Y)/2 - 150$ is a normal variable with mean $(E(X) + E(Y))/2 - 150 = 0$. We are using the bivariate normal assumption here to conclude that this linear combination is also normal.

(b) Find $P(Y < X)$.

$Y - X$ is a normal variable with mean $E(Y) - E(X) = 140 - 160 = -20$ and variance

$$\text{Var}(Y) + \text{Var}(X) - 2 \cdot \text{Cov}(Y, X) = 20 + 15 - 2 \cdot 0.6 \cdot \sqrt{20}\sqrt{15} \approx 14.22 \approx 3.77^2.$$

Thus, $\frac{Y - X + 20}{3.77}$ is a standard normal variable, and

$$\mathbb{P}(Y < X) = \mathbb{P}(Y - X + 20 < 20) = \mathbb{P}\left(\frac{Y - X + 20}{3.77} < 5.31\right) = \Phi(5.31) \approx 1.$$

2. Slices of a Normal Cake

You don't have to integrate in this exercise. In fact, integration is not recommended. Instead, draw some straight lines, think about angles, and remember the shape of the joint density of two i.i.d. standard normal variables.

(a) Let V and W be i.i.d. standard normal variables. Find $P(V > 0, W > \sqrt{3}V)$.

We are looking at the integral of the joint density of V and W over the region bordered by the half lines $x = 0, y > 0$ and $y = \sqrt{3}x, x > 0$. Since the joint density is rotationally symmetric, this is equal to the angle of these two half-lines divided by 2π , which is $\frac{\pi/2 - \arctan \sqrt{3}}{2\pi} = \frac{\pi/2 - \pi/3}{2\pi} = \frac{1}{12}$.

(b) In a population of adult mother-daughter pairs, the heights of the mothers and the daughters have correlation 0.6 and an approximately bivariate normal joint distribution. In roughly what proportion of mother-daughter pairs are both women above average in height?

If X is the height of a mother and Y is the height of her daughter, then with the usual notation, the quantity that we are looking for is

$$\mathbb{P}(X > \mu_X, Y > \mu_Y) = \mathbb{P}\left(\frac{X - \mu_X}{\sigma_X} > 0, \frac{Y - \mu_Y}{\sigma_Y} > 0\right) = \mathbb{P}(X^* > 0, Y^* > 0).$$

For the standardized variables $Y^* = 0.6 \cdot X^* + \sqrt{1 - 0.6^2} \cdot Z$ where $Z \sim N(0, 1)$ is independent of X . Hence,

$$\mathbb{P}(X^* > 0, Y^* > 0) = \mathbb{P}(X^* > 0, 0.6 \cdot X^* + \sqrt{1 - 0.6^2} \cdot Z > 0) = \mathbb{P}\left(X^* > 0, Z > -\frac{0.6}{\sqrt{1 - 0.6^2}}X^*\right).$$

By the same reasoning as in part (a), this is equal to

$$\frac{\pi/2 + \arctan \frac{0.6}{\sqrt{1 - 0.6^2}}}{2\pi} = \frac{\pi/2 + \arctan 0.75}{2\pi} = 0.3524$$

3. Permutations and Subsets

Let $\mathbf{X} = [X_1 \ X_2 \ \dots \ X_n]^T$ be multivariate normal with mean vector μ_X and covariance matrix Σ_X .

(a) Let \mathbf{Y} be a permutation of \mathbf{X} . Show that \mathbf{Y} is multivariate normal.

It is enough to show that every linear combination of Y_1, \dots, Y_n is normally distributed. But for any constants a_1, \dots, a_n , the linear combination $\sum_{i=1}^n a_i Y_i$ can be written as a linear combination of X_1, \dots, X_n with the same coefficients in some different order (following the given permutation), and that is normal.

Alternatively, we can notice that \mathbf{Y} is a linear transformation of \mathbf{X} by $\mathbf{Y} = \mathbf{P} \cdot \mathbf{X}$ where \mathbf{P} is the $n \times n$ permutation matrix of the given permutation π , that is, $\mathbf{P}_{ij} = 1$ if $j = \pi(i)$ and 0 otherwise.

(b) For $m < n$, let W_1, W_2, \dots, W_m be any subset of X_1, X_2, \dots, X_n and let $\mathbf{W} = [W_1 \ W_2 \ \dots \ W_m]^T$. Show that \mathbf{W} is multivariate normal.

Similarly to part (a), a linear combination $\sum_{i=1}^m a_i W_i$ can be written as a linear combination of a subset of X_1, \dots, X_n (with the exact same coefficients), and that is normal.

Alternatively, $\mathbf{W} = \mathbf{A} \cdot \mathbf{X}$, where \mathbf{A} is the $m \times n$ matrix that has $\mathbf{A}_{ij} = 1$ if $W_i = X_j$ and 0 otherwise.

4. Normal Sample Mean and Sample Variance, Part 1

(a) Let R have the chi-squared distribution with n degrees of freedom. What is the mgf of R ?

$R \stackrel{d}{=} Z_1^2 + \dots + Z_n^2$ where Z_i are i.i.d. standard normal variables. Hence, for the mgf's: $M_R(t) = (M_{Z_1^2}(t))^n$. We have

$$M_{Z_1^2}(t) = \int_{-\infty}^{\infty} e^{tx^2} \frac{1}{\sqrt{2\pi}} e^{-x^2/2} dx = \frac{1}{\sqrt{1-2t}} \int_{-\infty}^{\infty} \frac{1}{\sqrt{2\pi}} \sqrt{1-2t} \cdot e^{-\frac{x^2}{2}(1-2t)} dx = \frac{1}{\sqrt{1-2t}},$$

so

$$M_R(t) = \left(\frac{1}{1-2t} \right)^{n/2}$$

for $t < 1/2$.

(b) For R as in Part (a), suppose $R = V + W$ where V and W are independent and V has the chi-squared distribution with $m < n$ degrees of freedom. Can you identify the distribution of W ? Justify your answer.

For the mgf's:

$$\begin{aligned} M_R(t) &= M_V(t)M_W(t) \\ \left(\frac{1}{1-2t} \right)^{n/2} &= \left(\frac{1}{1-2t} \right)^{m/2} M_W(t) \\ M_W(t) &= \left(\frac{1}{1-2t} \right)^{(n-m)/2}, \end{aligned}$$

which is the mgf of the $\chi^2(n-m)$ distribution, so $W \sim \chi^2(n-m)$.

(c) Let X_1, X_2, \dots, X_n be any sequence of random variables and let $\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i$. Let α be any constant. Prove the **sum of squares decomposition**

$$\sum_{i=1}^n (X_i - \alpha)^2 = \sum_{i=1}^n (X_i - \bar{X})^2 + n(\bar{X} - \alpha)^2$$

$$\begin{aligned} \sum_{i=1}^n (X_i - \alpha)^2 &= \sum_{i=1}^n ((X_i - \bar{X}) + (\bar{X} - \alpha))^2 \\ &= \sum_{i=1}^n (X_i - \bar{X})^2 + \sum_{i=1}^n (\bar{X} - \alpha)^2 + 2 \sum_{i=1}^n (X_i - \bar{X})(\bar{X} - \alpha) \\ &= \sum_{i=1}^n (X_i - \bar{X})^2 + n(\bar{X} - \alpha)^2 + 2(\bar{X} - \alpha) \sum_{i=1}^n (X_i - \bar{X}) \\ &= \sum_{i=1}^n (X_i - \bar{X})^2 + n(\bar{X} - \alpha)^2 \end{aligned}$$

(d) Now let X_1, X_2, \dots, X_n be i.i.d. normal with mean μ and variance $\sigma^2 > 0$. Let S^2 be the “sample variance” defined by

$$S^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2$$

Use Parts (b) and (c) with an appropriate value of α to find a constant c such that cS^2 has a chi-squared distribution. Provide the degrees of freedom.

Using part (c) with $\alpha = \mu$ and dividing through by σ^2 we get

$$\sum_{i=1}^n \left(\frac{X_i - \mu}{\sigma} \right)^2 = \frac{1}{\sigma^2} \sum_{i=1}^n (X_i - \bar{X})^2 + n \left(\frac{\bar{X} - \mu}{\sigma} \right)^2.$$

Now since the $\frac{X_i - \mu}{\sigma}$ are i.i.d standard normal variables, the left-hand side above has the $\chi^2(n)$ distribution. $n \left(\frac{\bar{X} - \mu}{\sigma} \right)^2 = \left(\frac{\bar{X} - \mu}{\sigma/\sqrt{n}} \right)^2$ is also the square of a standard normal, so it has the $\chi^2(1)$ distribution. The other term on the right-hand side is equal to $\frac{n-1}{\sigma^2} S^2$. Using part (b) along with the independence of \bar{X} and S proven in 5 (c), we get that $\frac{n-1}{\sigma^2} S^2 \sim \chi^2(n-1)$.

5. Normal Sample Mean and Sample Variance, Part 2

Required reading: Carefully go through Exercise 6 in Homework 5. In Part (e) of that exercise you defined the sample variance of an i.i.d. sample. The same definition was used in the exercise above. Note that the Homework 5 exercise made no assumption about the underlying distribution of the elements of the sample, other than the existence of the expectation μ and variance σ^2 .

Let X_1, X_2, \dots, X_n be i.i.d. with mean μ and variance σ^2 . Define the sample mean \bar{X} and the sample variance S^2 as in Exercise 4.

(a) For $1 \leq i \leq n$ let $D_i = X_i - \bar{X}$. Find $Cov(D_i, \bar{X})$.

$$\begin{aligned}
\text{Cov}(D_i, \bar{X}) &= \text{Cov}(X_i - \bar{X}, \bar{X}) = \text{Cov}(X_i, \bar{X}) - \text{Cov}(\bar{X}, \bar{X}) \\
&= \text{Cov}\left(X_i, \frac{1}{n}X_i + \dots + \frac{1}{n}X_n\right) - \text{Var}(\bar{X}) \\
&= \text{Cov}\left(X_i, \frac{1}{n}X_i\right) - \frac{\sigma^2}{n} = \frac{\sigma^2}{n} - \frac{\sigma^2}{n} = 0
\end{aligned}$$

(b) Now assume in addition that X_1, X_2, \dots, X_n are i.i.d. normal (μ, σ^2) . What is the joint distribution of $\bar{X}, D_1, D_2, \dots, D_{n-1}$? Explain why D_n isn't on the list.

$\bar{X}, D_1, D_2, \dots, D_{n-1}$ are jointly normal since every linear combination of them is some linear combination of X_1, \dots, X_n , which is normal. The parameters are the following:

$$\begin{aligned}
E[\bar{X}] &= \mu \\
E[D_i] &= 0 \\
\text{Var}[\bar{X}] &= \frac{\sigma^2}{n} \\
\text{Var}[D_i] &= \text{Var}[X_i - \bar{X}] = \text{Var}[X_i] + \text{Var}[\bar{X}] - 2 \cdot \text{Cov}(X_i, \bar{X}) = \left(1 - \frac{1}{n}\right) \sigma^2 \\
\text{Cov}(\bar{X}, D_i) &= 0 \\
\text{Cov}(D_i, D_j) &= \text{Cov}(X_i - \bar{X}, X_j - \bar{X}) = \text{Cov}(X_i, X_j - \bar{X}) = -\text{Cov}(X_i, \bar{X}) = -\frac{\sigma^2}{n}
\end{aligned}$$

D_n is not on the list, since $D_1 + \dots + D_n = 0$, and this linear dependence makes the joint distribution of \bar{X}, D_1, \dots, D_n degenerate.

(c) True or false (justify your answer): The sample mean and sample variance of an i.i.d. normal sample are independent of each other.

We saw that for any i , \bar{X} and D_i are uncorrelated and they are jointly normal, which imply that they are independent. Therefore, \bar{X} and $S^2 = \frac{1}{n-1} \sum_{i=1}^n D_i^2$ are also independent.

6. The F Distribution

The F distribution is named for Sir Ronald Fisher, the creator of much of modern statistical theory. The distribution is on the positive real numbers and has two positive integer parameters which we will call n and d for “numerator” and “denominator”.

A random variable X has the $F_{n,d}$ distribution if $X = \frac{N/n}{D/d}$ where N and D are independent random variables such that N has the chi-squared (n) distribution and D has the chi-squared (d) distribution.

In statistics, the distribution arises for example in tests for whether two normal populations have the same variance. The variables N and D both arise as sums of squares of centered normal random variables.

The goal of this exercise is for you to derive the following impressive formula for the $F_{n,d}$ density:

$$f_X(x) = \frac{\Gamma\left(\frac{n}{2} + \frac{d}{2}\right)}{\Gamma\left(\frac{n}{2}\right)\Gamma\left(\frac{d}{2}\right)} \left(\frac{n}{d}\right)^{\frac{n}{2}} x^{\frac{n}{2}-1} \left(1 + \frac{n}{d}x\right)^{-\frac{n+d}{2}}, \quad x > 0$$

That looks horrendous but it really isn't. As always, the most impressive part is the constant of integration. The functional form is in fact rather straightforward. For example, if $n = 50$ and $d = 100$ then the density

becomes

$$f_X(x) = Cx^{24}(1 + 0.5x)^{-75}$$

which doesn't seem so bad after all.

(a) Find $E(N/n)$ and $E(D/d)$. This will help explain why F statistics are often compared to 1.

$$E(N/n) = E(D/d) = 1$$

(b) Since $X = \frac{d}{n}Y$ where $Y = \frac{N}{D}$, you will first find the density of Y and then the density of X . Start by writing Y in terms of $R = \frac{N}{N+D}$.

$$Y = \frac{N}{D} = \frac{N/(N+D)}{D/(N+D)} = \frac{R}{1-R}$$

(c) Go back to Lecture 11 given on 9/27 and review the beta-gamma algebra. Find the density of Y .

$N \sim \Gamma(n/2, 1/2)$ and $D \sim \Gamma(d/2, 1/2)$, so $R = \frac{N}{N+D} \sim \beta(n/2, d/2)$, hence, the density of R is

$$f_R(r) = \frac{\Gamma\left(\frac{n+d}{2}\right)}{\Gamma\left(\frac{n}{2}\right)\Gamma\left(\frac{d}{2}\right)} r^{n/2-1} (1-r)^{d/2-1}$$

for $r \in (0, 1)$. With a change of variable, the density of $Y = \frac{R}{1-R}$ is (with $y = \frac{r}{1-r}$, $r = \frac{y}{1+y}$):

$$\begin{aligned} f_Y(y) &= f_R(r) \cdot \left| \frac{dr}{dy} \right| = f_R\left(\frac{y}{1+y}\right) \cdot \frac{1}{(1+y)^2} \\ &= \frac{\Gamma\left(\frac{n+d}{2}\right)}{\Gamma\left(\frac{n}{2}\right)\Gamma\left(\frac{d}{2}\right)} \left(\frac{y}{1+y}\right)^{n/2-1} \left(\frac{1}{1+y}\right)^{d/2-1} \cdot \frac{1}{(1+y)^2} \end{aligned}$$

for $y > 0$.

(d) Now use Part (b) to find the density of X .

With another change of variable, the density of $X = \frac{d}{n}Y$ is (with $x = \frac{d}{n}y$)

$$\begin{aligned} f_X(x) &= f_Y(y) \cdot \left| \frac{dy}{dx} \right| = f_Y\left(\frac{n}{d} \cdot x\right) \cdot \frac{n}{d} \\ &= \frac{\Gamma\left(\frac{n+d}{2}\right)}{\Gamma\left(\frac{n}{2}\right)\Gamma\left(\frac{d}{2}\right)} \left(\frac{\frac{n}{d}x}{1+\frac{n}{d}x}\right)^{n/2-1} \left(\frac{1}{1+\frac{n}{d}x}\right)^{d/2-1} \cdot \frac{1}{(1+\frac{n}{d}x)^2} \cdot \frac{n}{d} \\ &= \frac{\Gamma\left(\frac{n}{2} + \frac{d}{2}\right)}{\Gamma\left(\frac{n}{2}\right)\Gamma\left(\frac{d}{2}\right)} \left(\frac{n}{d}\right)^{\frac{n}{2}} x^{\frac{n}{2}-1} \left(1 + \frac{n}{d}x\right)^{-\frac{n+d}{2}} \end{aligned}$$

for $x > 0$.

7. The t Distribution

This distribution is much used in tests for the mean of a normal population. It is sometimes called the *Student's t* distribution, and using related constructions is sometimes called Studentizing. That's in honor of the originator of the distribution, William Sealy Gosset. He worked in the Guinness Brewery in Ireland early in the 20th century and also spent time publishing statistical papers under the pseudonym Student.

A random variable T has the t distribution with d degrees of freedom if $T = \frac{Z}{\sqrt{D/d}}$ where Z is standard normal, D has the chi-squared (d) distribution, and Z and D are independent.

The goal of this problem is for you to find the density of T and also to identify a commonly used statistic that has a t distribution.

(a) Show that $T \stackrel{d}{=} -T$.

$-T = \frac{-Z}{\sqrt{D/d}}$ where $-Z$ is a standard normal variable, which is independent of D . Thus, $-T$ also has the t distribution with d degrees of freedom.

(b) Show that T^2 has an F distribution and hence use Exercise 6 to find the density of T^2 .

$T^2 = \frac{Z^2}{D/d}$, hence, T^2 has the $F_{1,d}$ distribution with density

$$f_{T^2}(x) = \frac{\Gamma(\frac{1}{2} + \frac{d}{2})}{\Gamma(\frac{1}{2})\Gamma(\frac{d}{2})} \left(\frac{1}{d}\right)^{\frac{1}{2}} x^{\frac{1}{2}-1} \left(1 + \frac{1}{d}x\right)^{-\frac{1+d}{2}}.$$

for $x > 0$.

(c) Find the density of $\sqrt{T^2}$, the positive square root of T^2 .

By a change of variable, the density of $\sqrt{T^2} = |T|$ is (with $y = \sqrt{x}$):

$$\begin{aligned} f_{|T|}(y) &= f_{T^2}(x) \cdot \left| \frac{dx}{dy} \right| = f_{T^2}(y^2) \cdot 2y \\ &= \frac{\Gamma(\frac{1}{2} + \frac{d}{2})}{\Gamma(\frac{1}{2})\Gamma(\frac{d}{2})} \left(\frac{1}{d}\right)^{\frac{1}{2}} y^{-1} \left(1 + \frac{1}{d}y^2\right)^{-\frac{1+d}{2}} \cdot 2y \\ &= 2 \cdot \frac{\Gamma(\frac{1}{2} + \frac{d}{2})}{\Gamma(\frac{1}{2})\Gamma(\frac{d}{2})} \left(\frac{1}{d}\right)^{\frac{1}{2}} \left(1 + \frac{1}{d}y^2\right)^{-\frac{1+d}{2}}. \end{aligned}$$

(d) Find the density of T . At this point this should require almost no calculation.

Since $T \stackrel{d}{=} -T$, for the density of T it holds that $f_T(t) = f_T(-t)$ for every t . Thus,

$$\begin{aligned} f_T(t) &= \frac{f_T(t) + f_T(-t)}{2} = \frac{f_{|T|}(|t|)}{2} \\ &= \frac{\Gamma(\frac{1}{2} + \frac{d}{2})}{\Gamma(\frac{1}{2})\Gamma(\frac{d}{2})} \left(\frac{1}{d}\right)^{\frac{1}{2}} \left(1 + \frac{1}{d}t^2\right)^{-\frac{1+d}{2}} \\ &= \frac{\Gamma(\frac{d+1}{2})}{\sqrt{d\pi}\Gamma(\frac{d}{2})} \left(1 + \frac{t^2}{d}\right)^{-\frac{1+d}{2}} \end{aligned}$$

for $t \in \mathbb{R}$.

(e) Let X_1, X_2, \dots, X_n be i.i.d. normal (μ, σ^2) random variables. As in Exercise 5, let \bar{X} be the sample mean and S^2 the sample variance. Let $S = \sqrt{S^2}$ be the positive square root of S^2 . Show that the random variable

$$\frac{\bar{X} - \mu}{S/\sqrt{n}}$$

has a t distribution and find the degrees of freedom.

$$\frac{\bar{X} - \mu}{S/\sqrt{n}} = \frac{(\bar{X} - \mu)/(\sigma/\sqrt{n})}{S/\sigma} = \frac{(\bar{X} - \mu)/(\sigma/\sqrt{n})}{\sqrt{S^2/\sigma^2}} = \frac{(\bar{X} - \mu)/(\sigma/\sqrt{n})}{\sqrt{[(n-1)S^2/\sigma^2]/(n-1)}},$$

where the nominator has the standard normal distribution, it is independent of the denominator by 5(c), and $\frac{(n-1)S^2}{\sigma^2}$ in the denominator has the $\chi^2(n-1)$ distribution by 4(d). Hence, $\frac{\bar{X}-\mu}{S/\sqrt{n}}$ has the t distribution with $n-1$ degrees of freedom.