

STAT 345/445 Lecture 18

Random sample from a normal population – Section 5.3

1 Sampling from the normal distribution

- Distribution of sample mean and sample variance
- Correlation and independence for normals
- The t -distribution
- The F distribution

Sampling from the Normal Distribution

Theorem 5.3.1: Distributions of \bar{X} and S^2

Let X_1, \dots, X_n be a random sample from $N(\mu, \sigma^2)$ and let

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

then

(a) \bar{X} and S^2 are independent

(b) $\bar{X} \sim N\left(\mu, \frac{\sigma^2}{n}\right)$

(c) $\frac{(n-1)S^2}{\sigma^2} \sim \chi_{n-1}^2$

! *surprising* Since \bar{X} and S^2

saw this already one both functions of

X_1, \dots, X_n .

Proof of Theorem 5.3.1 (a)

(a) \bar{X} and S^2 are independent

- Note first that S^2 is “over determined” with $n+1$ terms

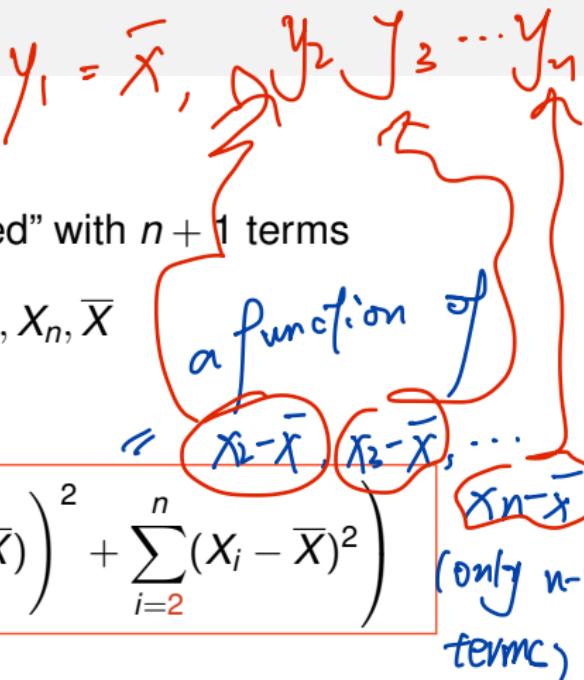
$X_1, X_2, \dots, X_n, \bar{X}$

Can write S^2 without X_1 :

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

$$(n-1)S^2 = \sum_{i=1}^n (X_i - \bar{X})^2 = (X_1 - \bar{X})^2 + \sum_{i=2}^n (X_i - \bar{X})^2$$

$= A$



$$\text{note: } n \cdot \bar{x} = x_1 + x_2 + \dots + x_n \Rightarrow x_1 = n\bar{x} - (x_2 + x_3 + \dots + x_n)$$

$$= (\bar{x} - x_2) + (\bar{x} - x_3) + \dots$$

$$+ (\bar{x} - x_n) + \bar{x}$$

$$= \bar{x} + \sum_{i=2}^n (\bar{x} - x_i)$$

$$A = \left(\bar{x} + \sum_{i=2}^n (\bar{x} - x_i) - \bar{x} \right)^2$$

$$= \left(\sum_{i=2}^n (\bar{x} - x_i) \right)^2 = \sum_{i=1}^n (x_i - \bar{x})^2$$

Proof of Theorem 5.3.1 (a) - continued

- We have written S^2 as a function of $n - 1$ terms:

$$Y_2, \dots, Y_n$$

$$(X_2 - \bar{X}), (X_3 - \bar{X}), \dots, (X_n - \bar{X})$$

- Show that \bar{X} and the random vector $(X_2 - \bar{X}, X_3 - \bar{X}, \dots, X_n - \bar{X})$ are independent

 - then we have shown that \bar{X} and S^2 are independent

$(X_1, \dots, X_n) \rightarrow (Y_1, \dots, Y_n)$ chw: Y_1 und

- Define an n dimensional transformation:

$$Y_1 = \bar{X}, \quad Y_2 = X_2 - \bar{X}, \quad Y_3 = X_3 - \bar{X}, \dots, \quad Y_n = X_n - \bar{X}$$

- Want to show that

S^2 is a function of these
 $f(y_1, y_2, y_3, \dots, y_n) = g(y_1) h(y_2, y_3, \dots, y_n)$

for some functions $g(\cdot)$ and $h(\cdot)$

Proof of Theorem 5.3.1 (a) - continued

- Find inverse functions $h_i(\mathbf{y})$ and Jacobian and then

$$f(\mathbf{y}) = f_{\mathbf{X}}(h_1(\mathbf{y}), \dots, h_n(\mathbf{y}) | J|)$$

- Assuming X_1, \dots, X_n are i.i.d. $N(0, 1)$ we have

$$f_{\mathbf{X}} = \prod_{i=1}^n \frac{1}{\sqrt{2\pi}} \exp(-x_i^2/2) = (2\pi)^{-n/2} \exp\left(-\frac{1}{2} \sum_{i=1}^n x_i^2\right)$$

- Therefore

$$f(\mathbf{y}) = \text{c } \sim - - \cdot - \left(\cancel{\times} \right)$$

$$y_1 = \bar{x}$$

$$\begin{aligned} y_2 &= x_2 - \bar{x} \\ &\vdots \\ y_n &= x_n - \bar{x} \end{aligned} \Rightarrow \begin{aligned} x_2 &= y_2 + \bar{x} \\ &\vdots \\ x_n &= y_n + \bar{x} \end{aligned} = y_2 + y_1$$

$$x_n = y_n + \bar{x}$$

$$\text{and } x_1 = n\bar{x} - (x_2 + x_3 + \dots + x_n)$$

$$\begin{aligned} &= ny_1 - [(y_2 + y_1) + (y_3 + y_1) + \dots + (y_n + y_1)] \end{aligned}$$

$$= y_1 - y_2 - y_3 - \dots - y_n$$

$$\frac{\partial x_1}{\partial y_1} = 1, \quad \frac{\partial x_1}{\partial y_2} = -1, \quad \dots, \quad \frac{\partial x_1}{\partial y_n} = -1$$

$$\frac{\partial x_2}{\partial y_1} = 1, \quad \frac{\partial x_2}{\partial y_2} = 1, \quad \frac{\partial x_2}{\partial y_3} = 0, \dots, \quad \frac{\partial x_L}{\partial y_n} = 0$$

for $k=2, \dots, n$

$$\frac{\partial x_k}{\partial y_1} = 1; \quad \frac{\partial x_k}{\partial y_k} = 1, \text{ otherwise } \frac{\partial x_k}{\partial y_j} = 0, j > k$$

$$= \left[\begin{array}{ccccccc|c}
1 & -1 & -1 & 0 & 0 & \dots & 0 & -1 \\
1 & 1 & 0 & - & \hline
0 & 1 & 0 & 0 & - & \dots & 0 & 0 \\
\vdots & \vdots & \vdots & \vdots & \vdots & \ddots & \vdots & \vdots \\
1 & 0 & 0 & 0 & \dots & \dots & 0 & 0
\end{array} \right]$$

Gauss-Elimination | odd all rows 2 from to the 1st row

def of a diagonal

matrix = multiple
of diagonal values

$$\left[\begin{array}{cccc} n & 0 & 0 & \cdots & 0 \\ 1 & 1 & 0 & \cdots & 0 \\ 1 & 0 & 1 & \cdots & 0 \\ \vdots & & & & \\ 1 & 0 & 0 & \cdots & 1 \end{array} \right]$$

$$= n [I_{n-1} = n \cdot \begin{bmatrix} 1 & & & \\ & 1 & & \\ & & \ddots & \\ & & & 1 \end{bmatrix}] = n$$

$$(*) = \int [y_1, \dots, y_n] = (2\pi)^{-n/2} \exp\left(-\frac{1}{2} \underbrace{(y_1 - y_2 - \dots - y_n)^2}_{\sum_{i=2}^n y_i}\right)$$

$$-\frac{1}{2} \sum_{i=2}^n (y_i - y_1) \Big) \cdot n$$

$$= n (2\pi)^{-\frac{n}{2}} \exp\left(-\frac{1}{2} \left(y_1^2 + \sum_{i=2}^n y_i^2 - 2y_1 \sum_{i=2}^n y_i + \dots\right)\right)$$

$$= \dots = \int [y_1] h(y_2, \dots, y_n) \text{ berücksichtigt}$$

About the χ_p^2 distribution

*(*1) or can show using mgf*

$$\begin{aligned} &= \text{Gamma}\left(\frac{\sum_i p_i}{2}, 2\right) \\ &= \chi_{p_1 + \dots + p_n}^2. \end{aligned}$$

Lemma 5.3.2

(a) If $Z \sim N(0, 1)$ then $Z^2 \sim \chi_1^2$

(b) If V_1, \dots, V_n are independent and $V_i \sim \chi_{p_i}^2$ then

$$V_1 + \dots + V_n \sim \chi_{p_1 + \dots + p_n}^2$$

(a): Seen before . just a univariate transformation
of $Z \sim N(0, 1)$

(b): Remember: χ_p^2 distribution. = Gamma $(\frac{p}{2}, 1)$ distr.

Seen before : Sum of independent Gammas with same beta β

is a Gamma with that same β and d

= sum of the original d 's $\Rightarrow \gamma_1 + \dots + \gamma_n \sim$

Gamma($\sum_{i=1}^n \frac{p_i}{2}$, 2)

Note: if x_1, \dots, x_n are iid $N(\mu, \sigma^2)$,

then $\frac{x_i - \mu}{\sigma} \stackrel{\text{iid}}{\sim} N(0, 1)$

$$\Rightarrow \frac{(x_i - \mu)^2}{\sigma^2} \sim \chi^2.$$

$$\Rightarrow \sum_{i=1}^n \frac{(x_i - \mu)^2}{\sigma^2} \sim \chi_n^2$$

$$= \frac{\sum_{i=1}^n (x_i - \mu)^2}{\sigma^2} \sim \chi_n^2$$

but $(n-1)S^2 = \sum_{i=1}^n (x_i - \bar{x})^2$

"lose one degree
of freedom"

Proof of Theorem 5.3.1 (c)

$$(c) \frac{(n-1)S^2}{\sigma^2} \sim \chi_{n-1}^2$$

$X_1, X_2, \dots, X_n \sim i.i.d N(0, 1)$

- Proof by induction. Set

$$\bar{X}_k = \frac{1}{k} \sum_{i=1}^k X_i \quad \text{and} \quad S_k^2 = \frac{1}{k-1} \sum_{i=1}^k (X_i - \bar{X}_k)^2$$

and note that

recursive
formulas

$$\bar{X}_{k+1} = \frac{k\bar{X}_k + X_{k+1}}{k+1}$$

$$\text{and } kS_{k+1}^2 = (k-1)S_k^2 + \frac{k}{k+1}(X_{k+1} - \bar{X}_k)^2$$

Can show
on your own

- Start with $k = 2 \dots$

X_1, \dots, X_n iid $N(\mu, \sigma^2)$, show $\frac{(n-1)S_n^2}{\sigma^2} \sim \chi_{n-1}^2$

(k=2) Show that $\frac{S_2^2}{\sigma^2} \sim \chi_1^2$

$$\frac{S_2^2}{\sigma^2} = \frac{1}{2!} \sum_{i=1}^2 (X_i - \bar{X}_2)^2 = \frac{1}{2!} \sum_{i=1}^2 \left[(X_i - \frac{X_1+X_2}{2})^2 + (X_i - \frac{X_1+X_2}{2})^2 \right]$$

$$= \frac{1}{2!} \left(\left(\frac{X_1-X_2}{2} \right)^2 + \left(\frac{X_2-X_1}{2} \right)^2 \right) = \frac{2}{2!} \left(\frac{X_1-X_2}{2} \right)^2 \\ = \frac{(X_1-X_2)^2}{2! \sigma^2}$$

$$X_1 - X_2 \sim N(0, 2\sigma^2) \Rightarrow \frac{X_1-X_2}{\sqrt{2\sigma^2}} \sim N(0, 1)$$

$$\Rightarrow \frac{(X_1-X_2)^2}{2\sigma^2} \sim \chi_1^2$$

Assume that $\frac{(k-1)S_k^2}{\chi^2} \sim \chi_{k-1}^2$ and Show that

$$\frac{KS_{k+1}^2}{\chi^2} \sim \chi_k^2$$

$$\frac{KS_{k+1}^2}{\chi^2} = \underbrace{\frac{k-1}{\chi^2} S_k^2}_{\sim \chi_{k-1}^2 \text{ by Assume}} + \underbrace{\frac{k}{\chi^2(k+1)} (\bar{X}_{k+1} - \bar{X}_k)^2}_{\sim \chi_1^2 \text{ indep to first term}}$$

$\sim \chi_{k-1}^2$ by Assume

by the two $\xrightarrow{*} \chi_{k-1+1}^2 = \chi_k^2$

S_k^2 and \bar{X}_{k+1} are indep. (Theorem)

S_k^2 and \bar{X}_{k+1} are indep.

and X_{k+1} and \bar{X}_k are indep \Rightarrow

Two terms are indep $\xrightarrow{\text{(*)}}$ and \bar{X}_k

$$X_{k+1} - \bar{X}_k \sim N(0, \sigma^2 + \sigma^2/k)$$

$$N(\mu, \sigma^2) \quad \backslash \quad N(\mu, \sigma^2/k) \quad \frac{\cancel{\mu} \sigma^2 + \sigma^2}{k} = \frac{k+1}{k} \sigma^2$$

$$\Rightarrow \frac{X_{k+1} - \bar{X}_k}{\sqrt{\frac{k+1}{k} \sigma^2}} \sim N(0, 1) \Rightarrow \frac{(X_{k+1} - \bar{X}_k)^2}{\frac{k+1}{k} \sigma^2} \sim \chi^2_1$$

$$\Rightarrow \frac{k}{(k+1)\sigma^2} (X_{k+1} - \bar{X}_k)^2 \sim \chi^2_1 *$$

Correlation and independence

- The normal distribution has a very special property

Independence \Leftrightarrow correlation = 0

Let (X_1, X_2) be a bivariate normal random vector. Then X_1 and X_2 are independent if and only if $\rho = \text{Cor}(X_1, X_2) = 0$

Recall the pdf of a bivariate normal pdf:

$$f(x_1, x_2) = \frac{1}{2\pi\sigma_1\sigma_2\sqrt{1-\rho^2}} \times \\ \exp \left\{ -\frac{1}{2(1-\rho^2)} \left[\frac{(x_1 - \mu_1)^2}{\sigma_1^2} - 2\rho \frac{(x_1 - \mu_1)(x_2 - \mu_2)}{\sigma_1\sigma_2} + \frac{(x_2 - \mu_2)^2}{\sigma_2^2} \right] \right\}$$

The student's t distribution

- Students- t distribution was developed because we want to know the distribution of the following statistic

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

- Used in inference about the mean
- Introduced by W. S. Gosset, who wrote under the alias “Student”
- The legend:
 - Gosset derived the t_m distribution while working for the Guinness Brewery in Dublin. In fear of competition he was forbidden to publish his analysis of brewery data and hence he wrote under the pseudonym Student.

The student's t distribution

Definition: The t distribution

Let U and V be independent random variables and $U \sim N(0, 1)$ and $V \sim \chi_p^2$. Then the distribution of

$$T = \frac{U}{\sqrt{V/p}}$$

$\xrightarrow{\text{N}(0, 1)}$
 $\xleftarrow{\sqrt{\chi_p^2 / p}}$

is called the **t distribution with p degrees of freedom, or t_p**

What is the pdf for T ?

Deriving the pdf for the student's t distribution

Let's do the derivation that Gosset did 120 years ago!

- Let $U \sim N(0, 1)$ and $V \sim \chi_p^2$ be independent
- Want the pdf of

$$T = \frac{U}{\sqrt{V/p}}$$

- Strategy: Do a bi-variate transformation

$$(U, V) \mapsto (T, W)$$

and then find the marginal of T

- Set $W = V$

Deriving the pdf for the student's t distribution - cont

R slide

The student's t distribution

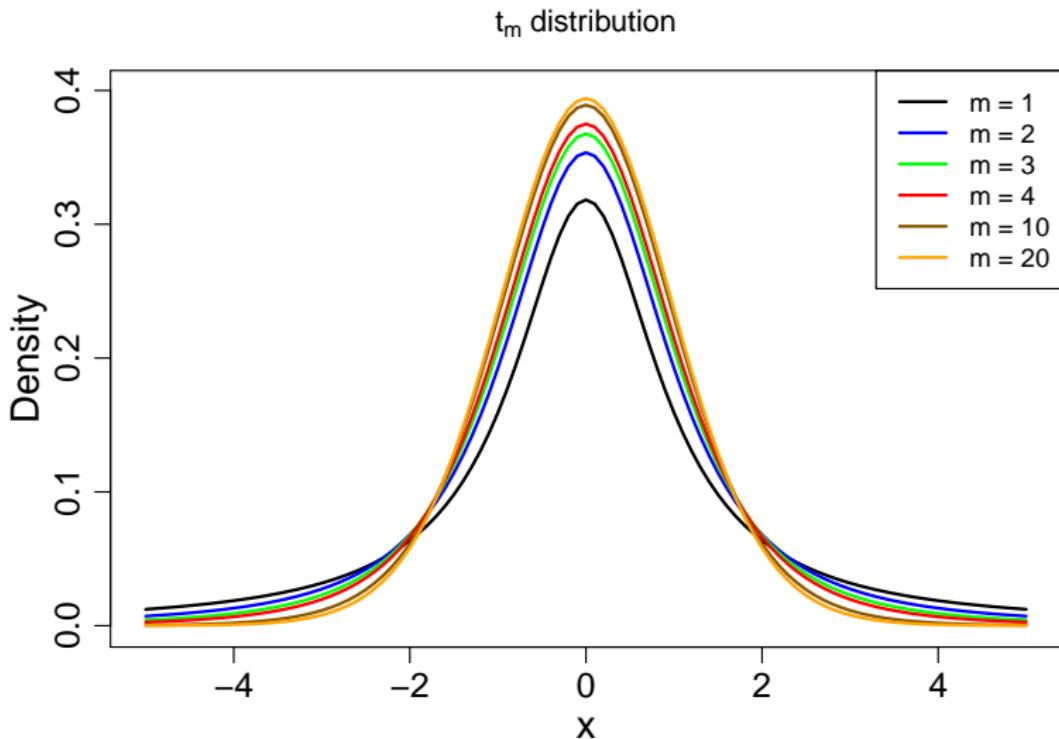
- The pdf of the t_p distribution is

$$f(x) = \frac{\Gamma\left(\frac{p+1}{2}\right)}{\Gamma\left(\frac{p}{2}\right)} \frac{1}{(p\pi)^{1/2}} \frac{1}{(1+x^2/p)^{(p+1)/2}}, \quad -\infty < x < \infty$$

The parameter p ("degrees of freedom") is an integer

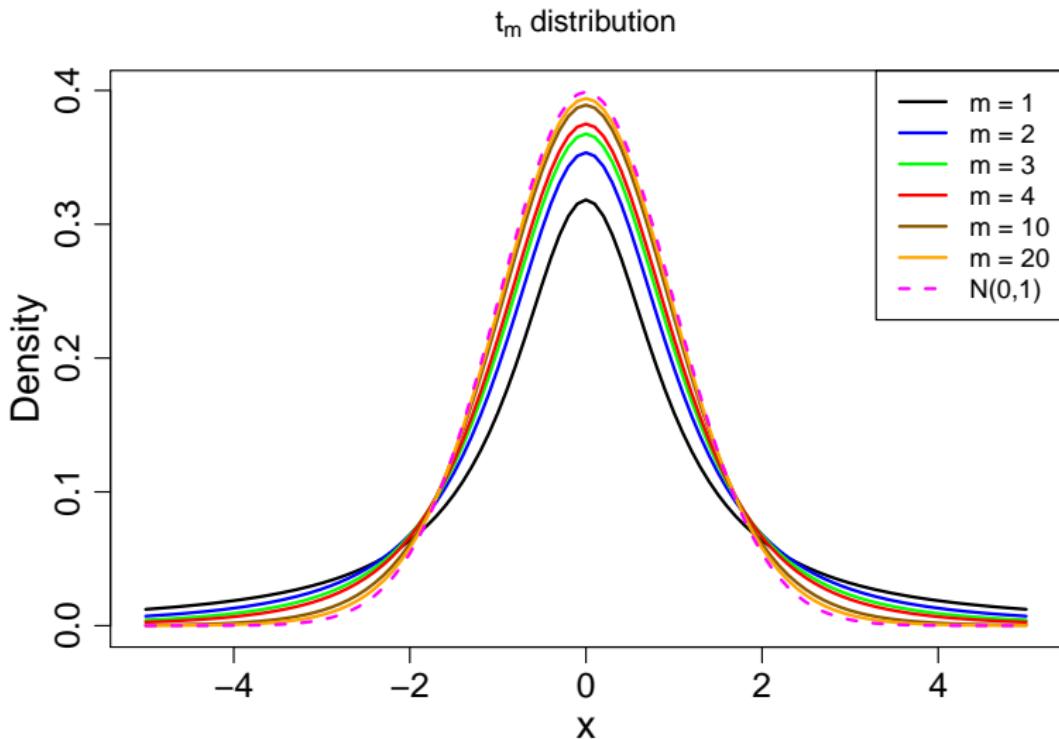
The t_p pdfs and the standard normal pdf

As $p \rightarrow \infty$ the t_p approaches $N(0, 1)$



The t_p pdfs and the standard normal pdf

As $p \rightarrow \infty$ the t_p approaches $N(0, 1)$



The T statistic

Right side

The T statistic

Let X_1, X_2, \dots, X_n be a random sample from a **normal distribution** with mean μ and variance σ^2 . Then the statistic

$$T = \frac{\bar{X} - \mu}{S/\sqrt{n}} \sim t_{n-1}$$

has **t distribution with $n - 1$ degrees of freedom, or t_{n-1}**

proof ...

know: $\frac{\bar{X} - \mu}{\sqrt{n}} \sim N(0, 1)$ and $\frac{(n-1)S^2}{\sigma^2} \sim \chi_{n-1}^2$

\bar{X} and S^2 are independent $\Rightarrow \frac{\bar{X} - \mu}{\sigma} \text{ and } \frac{(n-1)S^2}{\sigma^2}$ are indep

$$\Rightarrow \frac{\bar{x} - \mu}{\frac{s}{\sqrt{n}}} \sim \mathcal{E}_{n-1}$$

(n-1)s²
 s²(n-1)

→

$$= \frac{(\bar{x} - \mu)\sqrt{n}}{s} \Rightarrow \frac{\bar{x} - \mu}{\frac{s}{\sqrt{n}}} \sim \mathcal{E}_{n-1}$$

Snedecor's F distribution¹

- Comparing two variances: S_1^2/S_2^2
- Linear models: Ratio of sums-of-squares

Def: $F_{u,\nu}$ -distribution

Let $X \sim \chi_p^2$ and $Y \sim \chi_q^2$ be independent random variables. The distribution of

$$U = \frac{X/p}{Y/q} \sim F_{p,q}$$

is called the *F distribution with p and q degrees of freedom*

¹ *Trivia:* George W. Snedecor founded the first academic department of statistics in the United States, at Iowa State University in 1947.

F distribution $\mathcal{N} \sim \chi_p^2$, $\mathcal{Y} \sim \chi_q^2$ $u = \frac{\chi_p}{\chi_q}$

Define $V = Y$ (for example)

$f(x, y) = f(x) \cdot f(y)$ find $f_{(u, v)}$ and

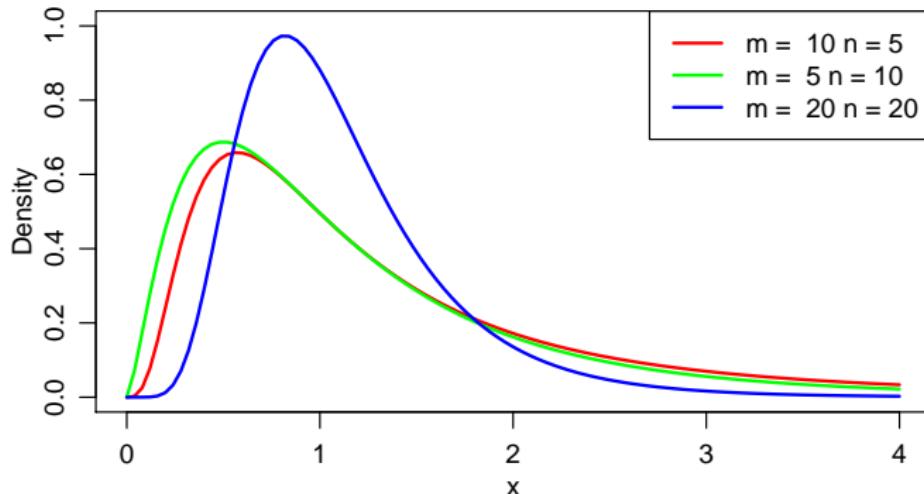
$$f(u) = \int f_{(u, v)} dy$$

F -distribution –pdf

F -distribution - pdf

$$f(x) = \frac{\Gamma\left(\frac{p+q}{2}\right)}{\Gamma\left(\frac{p}{2}\right)\Gamma\left(\frac{q}{2}\right)} \left(\frac{p}{q}\right)^{p/2} \frac{x^{(p/2)-1}}{(1 + (p/q)x)^{(p+q)/2}}, \quad 0 < x < \infty$$

F_{mn} distributions



The F statistic

$$\frac{(n_x - 1) S_x^2}{S_x^2} \sim \chi_{n_x - 1}^2 \quad \frac{(n_y - 1) S_y^2}{S_y^2} \sim \chi_{n_y - 1}^2$$

The F statistic

- Let X_1, X_2, \dots, X_n be a random sample from a normal distribution with mean μ_X and variance σ_X^2 .
- Let Y_1, Y_2, \dots, Y_m be a random sample from a normal distribution with mean μ_Y and variance σ_Y^2 .

Then the statistic

$$F = \frac{S_x^2 / \sigma_X^2}{S_y^2 / \sigma_Y^2} \sim F_{n-1, m-1}$$

$$n_x = n$$

$$n_y = m$$

has F distribution with $n - 1$ and $m - 1$ degrees of freedom, or

$$F_{n-1, m-1}$$

$$\Rightarrow \frac{(n_x - 1) S_x^2 / S_x^2}{(n_y - 1) S_y^2 / S_y^2} \sim \chi_{n_x - 1, n_y - 1}^2$$

F -distribution – properties

Some useful (univariate) transformations of the F distribution

Theorem 5.3.8

- (a) If $X \sim F_{p,q}$ then $1/X \sim F_{q,p} \Rightarrow$ tables in their test books only know
- (b) If $X \sim t_q$ then $X^2 \sim F_{1,q}$
- (c) If $X \sim F_{p,q}$ then $\frac{(p/q)X}{1+(p/q)X} \sim \text{Beta}(q/2, p/2)$ upper quantiles

\Rightarrow In STAT 325: Marginal t -test in the summary table is equivalent to the last F -test in the (sequential) Anova Table.

$$\frac{\int_{1-\alpha/2}^{\infty} f_{F_{1,q}}(x) dx}{\int_0^{\alpha/2} f_{F_{1,q}}(x) dx} = \frac{1}{F_{\alpha/2, q, 1}}$$

T and F statistics – uses

- The T statistic is used in inference of the mean of a normal distribution when variance is unknown, e.g.
 - Population mean of one or two populations
 - t -test, two-sample t -test, paired t -test ...
 - Regression coefficients in a normal linear regression model
- The F statistic is used in many situations, e.g.
 - to test for equality of variances from two independent populations
 - Analysis of Variance (ANOVA)
 - comparing nested models in normal linear regression

Linear model example

Coefficients:

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	1.6100262	7.8418734	0.205	0.837
PercPoverty	4.0099847	0.2807055	14.285	< 2e-16 ***
PerclUnemployed	-2.2020717	0.4689496	-4.696	3.57e-06 ***
IncomePerCapita	0.0018748	0.0003045	6.157	1.69e-09 ***

Signif. codes: 0 ‘***’ 0.001 ‘**’ 0.01 ‘*’ 0.05 ‘.’ 0.1 ‘ ’ 1

Residual standard error: 20.55 on 433 degrees of freedom
 Multiple R-squared: 0.3246, Adjusted R-squared: 0.3199
 F-statistic: 69.36 on 3 and 433 DF, p-value: < 2.2e-16

> anova(Fit)

Analysis of Variance Table

Response: CrimeRate

	Df	Sum Sq	Mean Sq	F value	Pr(>F)
PercPoverty	1	60376	60376	142.989	< 2.2e-16 ***
PerclUnemployed	1	11476	11476	27.179	2.878e-07 ***
IncomePerCapita	1	16008	16008	37.912	1.690e-09 ***
Residuals	433	182832	422		

Signif. codes: 0 ‘***’ 0.001 ‘**’ 0.01 ‘*’ 0.05 ‘.’ 0.1 ‘ ’ 1

Summary table

t-test for

$$H_0: \beta_k = 0 \text{ vs}$$

$$H_1: \beta_k \neq 0$$

$$\sqrt{37.912}$$

Same test