

STAT 24400 Lecture 13

Section 6.2  $\chi^2$ , t, and F Distributions

Section 6.3 Sample Mean & Sample Variance

Yibi Huang  
Department of Statistics  
University of Chicago

## Section 6.2 $\chi^2$ , t, and F Distributions

## Chapter 6 Distributions Derived from Normal

There are 3 distributions derived from the normal distributions that occur many statistical problems

- ▶ Chi-Squared ( $\chi^2$ ) distributions
  - ▶ “Chi-squared” is read “kai-squared”
- ▶ t distributions
- ▶ F distributions

## Definitions: Chi-Squared Distributions

Let  $Z_1, Z_2, \dots, Z_n$  be i.i.d.  $\sim N(0, 1)$ . The random variable

$$T_n = \sum_{i=1}^n Z_i^2$$

is said to be a **chi-squared distribution** with  **$n$  degrees of freedom**, denoted as

$$T_n \sim \chi_n^2.$$

In HW9, we show using MGF that chi-squared distributions are special Gamma distributions that

$$\chi_n^2 = \text{Gamma}(\alpha = n/2, \lambda = 1/2)$$

and the corresponding PDF is

$$f_{T_n}(t) = \frac{1}{2^{n/2} \Gamma(\frac{n}{2})} t^{(n/2)-1} e^{-t/2}, \quad t > 0.$$

# Properties of Chi-Squared Distributions

If  $Y \sim \chi_n^2$ , then its MGF is

$$M(t) = (1 - 2t)^{-n/2},$$

from which we can derive its expected value and variance

- ▶  $E[Y] = n$
- ▶  $\text{Var}(Y) = 2n$
- ▶ If  $U \sim \chi_n^2$  and  $V \sim \chi_m^2$  are independent, then  $U + V \sim \chi_{m+n}^2$ 
  - ▶ The proof is straight forward using MGF

## Definition: (Student's) $t$ -Distributions

If  $Z \sim N(0, 1)$  and  $U \sim \chi_n^2$  and  $Z$  and  $U$  are independent, then the distribution of

$$T = \frac{Z}{\sqrt{U/n}}$$

is called the **(Student's)  $t$ -distribution** with  **$n$  degrees of freedom**, denoted as

$$T \sim t_n.$$

The PDF is given by

$$f(t) = \frac{\Gamma(\frac{n+1}{2})}{\sqrt{n\pi} \Gamma(\frac{n}{2})} \left(1 + \frac{t^2}{n}\right)^{-(n+1)/2}, \quad -\infty < t < \infty.$$

## Proof for the PDF of the $t$ -Distribution

By the independence of  $Z$  and  $U$ , their joint PDF is given by

$$f_{ZU}(z, u) = \frac{1}{\sqrt{2\pi}} e^{-z^2/2} \cdot \frac{1}{2^{n/2} \Gamma(\frac{n}{2})} u^{\frac{n}{2}-1} e^{-u/2} = \frac{u^{\frac{n}{2}-1} \exp(-\frac{1}{2}(z^2 + u))}{\sqrt{\pi} 2^{\frac{n+1}{2}} \Gamma(\frac{n}{2})}, \quad \begin{matrix} -\infty < z < \infty, \\ u > 0. \end{matrix}$$

Consider the transformation  $W = \frac{Z}{\sqrt{U}}$ ,  $Y = U$ , with inverse transformation

$$\begin{matrix} Z = W\sqrt{Y}, \\ U = Y \end{matrix} \Rightarrow \text{Jacobian} = \begin{vmatrix} \frac{\partial z}{\partial w} & \frac{\partial z}{\partial y} \\ \frac{\partial u}{\partial w} & \frac{\partial u}{\partial y} \end{vmatrix} = \begin{vmatrix} \sqrt{y} & \frac{w}{2\sqrt{y}} \\ 0 & 1 \end{vmatrix} = \sqrt{y}.$$

The joint PDF for  $(W, Y)$  is

$$\begin{aligned} f_{WY}(w, y) &= f_{ZU}(w\sqrt{y}, y) \cdot \sqrt{y} \\ &= \frac{y^{\frac{n}{2}-1} \exp(-\frac{1}{2}(w^2 y + y))}{\sqrt{\pi} 2^{\frac{n+1}{2}} \Gamma(\frac{n}{2})} \sqrt{y} = \frac{y^{\frac{n+1}{2}-1} \exp(-\frac{y}{2}(1 + w^2))}{\sqrt{\pi} 2^{\frac{n+1}{2}} \Gamma(\frac{n}{2})} \end{aligned}$$

The marginal PDF for  $W$  can be obtained by integrating  $f_{WY}(w, y)$  over  $y$ .

$$f_W(w) = \int_0^\infty f_{WY}(w, y) dy = \frac{1}{\sqrt{\pi} 2^{\frac{n+1}{2}} \Gamma(\frac{n}{2})} \int_0^\infty y^{\frac{n+1}{2}-1} e^{-\frac{y}{2}(1+w^2)} dy.$$

Let

$$x = \frac{y}{2}(1+w^2) \Rightarrow y = \frac{2x}{1+w^2}, \quad dy = \frac{2}{1+w^2} dx.$$

Then,

$$\begin{aligned} f_W(w) &= \frac{1}{\sqrt{\pi} 2^{\frac{n+1}{2}} \Gamma(\frac{n}{2})} \int_0^\infty \left( \frac{2x}{1+w^2} \right)^{\frac{n+1}{2}-1} e^{-x} \frac{2}{1+w^2} dx \\ &= \frac{1}{\sqrt{\pi} \Gamma(\frac{n}{2}) (1+w^2)^{\frac{n+1}{2}}} \underbrace{\int_0^\infty x^{\frac{n+1}{2}-1} e^{-x} dx}_{=\Gamma(\frac{n+1}{2})} \\ &= \frac{\Gamma(\frac{n+1}{2})}{\sqrt{\pi} \Gamma(\frac{n}{2})} (1+w^2)^{-\frac{n+1}{2}}, \quad -\infty < w < \infty. \end{aligned}$$



The PDF for  $T = \frac{Z}{\sqrt{U/n}} = \sqrt{n} W$  is

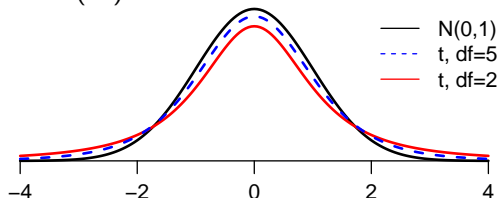
$$f_T(t) = \frac{1}{\sqrt{n}} f_W\left(\frac{t}{\sqrt{n}}\right) = \frac{\Gamma(\frac{n+1}{2})}{\sqrt{n\pi}\Gamma(\frac{n}{2})} \left(1 + \frac{t^2}{n}\right)^{-\frac{n+1}{2}}, \quad -\infty < t < \infty.$$

# Properties of $t$ -Distributions

For  $T \sim t_n$  with the PDF

$$f(t) = \frac{\Gamma(\frac{n+1}{2})}{\sqrt{n\pi} \Gamma(\frac{n}{2})} \left(1 + \frac{t^2}{n}\right)^{-(n+1)/2}, \quad -\infty < t < \infty.$$

- ▶ Bell-shaped, symmetric about 0
- ▶ With 1 degree of freedom,  $t_1 = \text{Cauchy}$
- ▶  $E[T] = 0$  if  $\text{df} > 1$
- ▶ For large  $t$ , the  $t$ -density with  $n$  df is  $\approx \frac{\text{constant}}{t^{n+1}} \Rightarrow$  **heavier tail** than normal
- ▶  $E[T^k]$  doesn't exist if  $k \geq \text{degrees of freedom (df)}$
- ▶ higher df  $\Rightarrow$  lighter tails
- ▶ As  $\text{df} \rightarrow \infty$ ,  $t \rightarrow N(0, 1)$



## Definition: $F$ -Distributions

Let  $U$  and  $V$  be independent chi-square random variables with  $m$  and  $n$  degrees of freedom, respectively. The distribution of

$$X = \frac{U/m}{V/n}$$

is called the  *$F$ -distribution with  $m$  and  $n$  degrees of freedom*, denoted by

$$X \sim F_{m,n}.$$

The PDF is given by

$$f(x) = \frac{\Gamma(\frac{m+n}{2})}{\Gamma(\frac{m}{2})\Gamma(\frac{n}{2})} \left(\frac{m}{n}\right)^{m/2} x^{\frac{m}{2}-1} \left(1 + \frac{m}{n}x\right)^{-\frac{m+n}{2}}, \quad x > 0.$$

The PDF can be obtained similarly as for the  $t$ -distribution.

# Properties of $F$ -Distributions

For  $X \sim F_{m,n}$

- ▶  $E(X) = \frac{n}{n-2}$  if  $n > 2$
- ▶  $E(X^k)$  exists only if  $k < n/2$
- ▶ If  $T \sim t_n$ , then  $T^2 \sim F_{1,n}$
- ▶ asymmetric PDF
- ▶  $F$ -distribution can be transformed to Beta distribution

$$X \sim F_{m,n} \quad \Rightarrow \quad Y = \frac{(m/n)X}{1 + (m/n)X} \sim \text{Beta} \left( a = \frac{n}{2}, b = \frac{m}{2} \right)$$

## Section 6.3 Sample Mean & Sample Variance

## First Statistics Question in STAT 24400

If we observed

$X_1, X_2, \dots, X_n$ , i.i.d.  $\sim N(\mu, \sigma^2)$ , but  $\mu$  and  $\sigma^2$  are UNKNOWN.

How to use the observed values of  $X_1, X_2, \dots, X_n$  to estimate the unknown  $\mu$  and  $\sigma^2$ ?

# First Statistics Question in STAT 24400

If we observed

$X_1, X_2, \dots, X_n$ , i.i.d.  $\sim N(\mu, \sigma^2)$ , but  $\mu$  and  $\sigma^2$  are UNKNOWN.

How to use the observed values of  $X_1, X_2, \dots, X_n$  to estimate the unknown  $\mu$  and  $\sigma^2$ ?

- ▶  $X_1, X_2, \dots, X_n$  are sometimes called the *sample*,  $n$  is called the *sample size*
- ▶ Usually estimate  $\mu$  by  $\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i$ , called the *sample mean*.
- ▶ As  $\sigma^2 = E[(X_i - \mu)^2]$ , one might attempt to estimate it by

$$\frac{\sum_{i=1}^n (X_i - \mu)^2}{n}.$$

However,  $\mu$  is unknown. We thus estimate  $\sigma^2$  by

$$S^2 = \frac{\sum_{i=1}^n (X_i - \bar{X})^2}{n - 1}, \quad \text{called the } \textit{sample variance}.$$

- ▶ Why divide by  $n - 1$ , not  $n$ ?
- ▶ We will discuss estimation problems in Chapter 8 in detail

# Population Mean/Variance v.s. Sample Mean/Variance

If  $X_1, \dots, X_n$  are i.i.d.  $\sim N(\mu, \sigma^2)$ ,

►  $\mu$  is called the *population mean*

►  $\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i$  is called the *sample mean*

►  $\sigma^2$  is called the *population variance*

►  $S^2 = \frac{\sum_{i=1}^n (X_i - \bar{X})^2}{n-1}$  is called the *sample variance*



# Sample Mean

For the sample mean

$$\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i$$

we have shown earlier that

$$E(\bar{X}) = \mu \quad \text{and} \quad \text{Var}(\bar{X}) = \frac{\sigma^2}{n}$$

and

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

## A Useful Identity

The following identity always holds for any value of  $c$ .

$$\sum_{i=1}^n (X_i - c)^2 = \sum_{i=1}^n (X_i - \bar{X})^2 + n(\bar{X} - c)^2, \quad \text{where } \bar{X} = \frac{1}{n} \sum_{i=1}^n X_i.$$

*Proof.*

$$\begin{aligned} \sum_{i=1}^n (X_i - c)^2 &= \sum_{i=1}^n (X_i - \bar{X} + \bar{X} - c)^2 \\ &= \sum_{i=1}^n (X_i - \bar{X})^2 + 2 \sum_{i=1}^n (X_i - \bar{X}) \underbrace{(\bar{X} - c)}_{\text{constant}} + \sum_{i=1}^n \underbrace{(\bar{X} - c)^2}_{\text{constant}} \\ &= \sum_{i=1}^n (X_i - \bar{X})^2 + 2(\bar{X} - c) \underbrace{\sum_{i=1}^n (X_i - \bar{X})}_{=0, \text{ see below}} + n(\bar{X} - c)^2 \end{aligned}$$

$$\text{where } \sum_{i=1}^n (X_i - \bar{X}) = \sum_{i=1}^n X_i - \sum_{i=1}^n \bar{X} = n\bar{X} - n\bar{X} = 0.$$

## Corollary of the Useful Identity

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

- The case  $c = 0$  gives the shortcut formula for the sample variance

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

- The value  $c$  that minimizes  $\sum_{i=1}^n (X_i - c)^2$  is  $c = \bar{X}$ .

## Expectation of Sample Variance (Why Divide by $n - 1$ , not $n$ ?)

Letting  $c = \mu = E[X_i]$  in the useful identity

$$\sum_{i=1}^n (X_i - \mu)^2 = \underbrace{\sum_{i=1}^n (X_i - \bar{X})^2}_{=(n-1)S^2} + n(\bar{X} - \mu)^2.$$

gives the following expression for  $S^2$

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

Taking expected values on both sides, we get

$$\begin{aligned} E[S^2] &= \frac{1}{n-1} \left( \sum_{i=1}^n \underbrace{E[(X_i - \mu)^2]}_{=\text{Var}(X_i)=\sigma^2} - n \underbrace{E[(\bar{X} - \mu)^2]}_{=\text{Var}(\bar{X})=\sigma^2/n} \right) \\ &= \frac{1}{n-1} \left( n\sigma^2 - n \cdot \frac{\sigma^2}{n} \right) = \sigma^2. \end{aligned}$$

$\overline{X}$  is Independent of  $S^2$

## $\bar{X}$ Is Independent of $S^2$ (★★)

We will first prove that

$$\bar{X} \text{ is indep. of } \overbrace{(X_2 - \bar{X}, X_3 - \bar{X}, \dots, X_n - \bar{X})}^{\text{no } X_1 - \bar{X}}.$$

This would imply  $\bar{X}$  is independent of  $S^2$  since  $(n-1)S^2$  can be written as a function of  $(X_2 - \bar{X}, X_3 - \bar{X}, \dots, X_n - \bar{X})$  as follows

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

where  $X_1 - \bar{X} = -\sum_{i=2}^n (X_i - \bar{X})$  since  $\sum_{i=1}^n (X_i - \bar{X}) = 0$ .

## $\bar{X}$ Is Independent of $S^2$ (★★)

We will first prove that

$$\bar{X} \text{ is indep. of } \overbrace{(X_2 - \bar{X}, X_3 - \bar{X}, \dots, X_n - \bar{X})}^{\text{no } X_1 - \bar{X}}.$$

This would imply  $\bar{X}$  is independent of  $S^2$  since  $(n-1)S^2$  can be written as a function of  $(X_2 - \bar{X}, X_3 - \bar{X}, \dots, X_n - \bar{X})$  as follows

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

where  $X_1 - \bar{X} = -\sum_{i=2}^n (X_i - \bar{X})$  since  $\sum_{i=1}^n (X_i - \bar{X}) = 0$ .

Steps of the proof:

1. find the joint PDF  $f_{\mathbf{Y}}(y_1, y_2, \dots, y_n)$  of  $Y_1 = \bar{X}$ ,  $Y_i = X_i - \bar{X}$  for  $i = 2, \dots, n$ .
2. show that the joint PDF  $f_{\mathbf{Y}}(y_1, y_2, \dots, y_n)$  can factor as the product of a function of  $y_1$  and a function of  $(y_2, \dots, y_n)$ .

$$f(y_1, y_2, \dots, y_n) = g(y_1)h(y_2, \dots, y_n), \quad \text{for all } y_1, y_2, \dots, y_n.$$

## Multivariate Transformation

Suppose  $(X_1, \dots, X_n)$  are continuous r.v.'s with joint PDF

$$f_{\mathbf{X}}(x_1, \dots, x_n).$$

They are mapped onto  $(Y_1, \dots, Y_n)$  by a 1-to-1 transformation

$$\begin{aligned} y_1 &= g_1(x_1, \dots, x_n) \\ &\vdots \\ y_n &= g_n(x_1, \dots, x_n) \end{aligned}$$

and the transformation can be inverted to obtain

$$\begin{aligned} x_1 &= h_1(y_1, \dots, y_n) \\ &\vdots \\ x_n &= h_n(y_1, \dots, y_n). \end{aligned}$$



The joint PDF  $f_{\mathbf{Y}}(y_1, \dots, y_n)$  is given by

$$f_{\mathbf{Y}}(y_1, \dots, y_n) = f_{\mathbf{X}}(h_1(y_1, \dots, y_n), \dots, h_n(y_1, \dots, y_n)) \left| \frac{\partial(x_1, \dots, x_n)}{\partial(y_1, \dots, y_n)} \right|,$$

where  $\left| \frac{\partial(x_1, \dots, x_n)}{\partial(y_1, \dots, y_n)} \right|$  is **absolute value** of the *Jacobian of the transformation*, defined as the determinant of the  $n \times n$  matrix

$$\begin{vmatrix} \frac{\partial x_1}{\partial y_1} & \frac{\partial x_1}{\partial y_2} & \dots & \frac{\partial x_1}{\partial y_n} \\ \frac{\partial x_2}{\partial y_1} & \frac{\partial x_2}{\partial y_2} & \dots & \frac{\partial x_2}{\partial y_n} \\ \vdots & \ddots & \vdots & \vdots \\ \frac{\partial x_n}{\partial y_1} & \frac{\partial x_n}{\partial y_2} & \dots & \frac{\partial x_n}{\partial y_n} \end{vmatrix}$$

that the  $(i, j)$  element is  $\frac{\partial x_i}{\partial y_j}$ .

## Joint PDF of $\bar{X}$ and $(X_2 - \bar{X}, \dots, X_n - \bar{X})$

For  $Y_1 = \bar{X}$ ,  $Y_i = X_i - \bar{X}$ , for  $i = 2, 3, \dots, n$ , the inverse transformation is

$$X_1 = Y_1 - (Y_2 + Y_3 + \dots + Y_n),$$

$$X_i = Y_1 + Y_i, \quad \text{for } i = 2, 3, \dots, n.$$

We see

$$\frac{\partial x_1}{\partial y_j} = \begin{cases} 1 & \text{if } j = 1 \\ -1 & \text{if } j = 2, 3, \dots, n, \end{cases} \quad \text{and} \quad \frac{\partial x_i}{\partial y_j} = \begin{cases} 1 & \text{if } j = 1 \text{ or } i \\ 0 & \text{otherwise.} \end{cases}$$

The Jacobian matrix is

$$\begin{vmatrix} \frac{\partial x_1}{\partial y_1} & \frac{\partial x_1}{\partial y_2} & \frac{\partial x_1}{\partial y_3} & \dots & \frac{\partial x_1}{\partial y_n} \\ \frac{\partial x_2}{\partial y_1} & \frac{\partial x_2}{\partial y_2} & \frac{\partial x_2}{\partial y_3} & \dots & \frac{\partial x_2}{\partial y_n} \\ \frac{\partial x_3}{\partial y_1} & \frac{\partial x_3}{\partial y_2} & \frac{\partial x_3}{\partial y_3} & \dots & \frac{\partial x_3}{\partial y_n} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ \frac{\partial x_n}{\partial y_1} & \frac{\partial x_n}{\partial y_2} & \frac{\partial x_n}{\partial y_3} & \dots & \frac{\partial x_n}{\partial y_n} \end{vmatrix} = \begin{vmatrix} 1 & -1 & -1 & \dots & -1 \\ 1 & 1 & 0 & \dots & 0 \\ 1 & 0 & 1 & \dots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & 0 & 0 & \dots & 1 \end{vmatrix}$$

The determinant can be shown by induction to be  $n$ .

As  $X_i$ 's are independent, their joint PDF is

$$\begin{aligned} f_{\mathbf{X}}(x_1, x_2, \dots, x_n) &= \frac{1}{(2\pi)^{n/2} \sigma^n} \exp \left( \frac{-1}{2\sigma^2} \sum_{i=1}^n (x_i - \mu)^2 \right) \\ &= \frac{1}{(2\pi)^{n/2} \sigma^n} \exp \left( \frac{-1}{2\sigma^2} \left( \sum_{i=1}^n (x_i - \bar{x})^2 + n(\bar{x} - \mu)^2 \right) \right) \end{aligned}$$

in which,  $x_1 - \bar{x} = -\sum_{i=2}^n (x_i - \bar{x}) = -\sum_{i=2}^n y_i$

$$\sum_{i=1}^n (x_i - \bar{x})^2 = (x_1 - \bar{x})^2 + \sum_{i=2}^n (x_i - \bar{x})^2 = \left( \sum_{i=2}^n y_i \right)^2 + \sum_{i=2}^n y_i^2$$

The joint PDF of  $(Y_1, \dots, Y_n)$  is thus

$$f_{\mathbf{Y}}(y_1, y_2, \dots, y_n) = \frac{|J|}{(2\pi)^{n/2} \sigma^n} \exp \left[ \frac{-1}{2\sigma^2} \left( \left( \sum_{i=2}^n y_i \right)^2 + \sum_{i=2}^n y_i^2 + n(y_1 - \mu)^2 \right) \right]$$

where  $|J| = n$  is the Jacobian shown on the previous page.

We can see the joint PDF  $f_{\mathbf{Y}}(y_1, y_2, \dots, y_n)$  can factor into

- ▶ a function  $\exp(-\frac{n}{2\sigma^2}(y_1 - \mu)^2)$  of  $y_1$ , and
- ▶ a function  $\exp[\frac{-1}{2\sigma^2}((\sum_{i=2}^n y_i)^2 + \sum_{i=2}^n y_i^2)]$  of  $y_2, \dots, y_n$ .

This proves the independence of

$$Y_1 = \bar{X} \quad \text{and} \quad (Y_2, \dots, Y_n) = (X_2 - \bar{X}, \dots, X_n - \bar{X}),$$

which implies the independence of  $\bar{X}$  and  $S^2$ .

## Distribution of $S^2$

If  $X_1, X_2, \dots, X_n$  are i.i.d.  $\sim N(\mu, \sigma^2)$ , then

$$\frac{X_1 - \mu}{\sigma}, \frac{X_2 - \mu}{\sigma}, \dots, \frac{X_n - \mu}{\sigma} \text{ are i.i.d. } \sim N(0, 1),$$

which implies

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

has a **chi-squared** distribution with  $n$  degrees of freedom.

If  $X_1, X_2, \dots, X_n$  are i.i.d.  $\sim N(\mu, \sigma^2)$ , then

$$\frac{X_1 - \mu}{\sigma}, \frac{X_2 - \mu}{\sigma}, \dots, \frac{X_n - \mu}{\sigma} \text{ are i.i.d. } \sim N(0, 1),$$

which implies

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

has a **chi-squared** distribution with  $n$  degrees of freedom.

**Question:** What's the distribution of

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

If  $X_1, X_2, \dots, X_n$  are i.i.d.  $\sim N(\mu, \sigma^2)$ , then

$$\frac{X_1 - \mu}{\sigma}, \frac{X_2 - \mu}{\sigma}, \dots, \frac{X_n - \mu}{\sigma} \text{ are i.i.d. } \sim N(0, 1),$$

which implies

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

has a **chi-squared** distribution with  $n$  degrees of freedom.

**Question:** What's the distribution of

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

**Ans:** **chi-squared** distribution with  $n - 1$  degrees of freedom.



## Proof of $(n-1)S^2/\sigma^2 \sim \chi_{n-1}^2$

Define  $V_1, V_2, V_3$  as follows:

$$\underbrace{\frac{\sum_{i=1}^n (X_i - \mu)^2}{\sigma^2}}_{=V_1} = \underbrace{\frac{\sum_{i=1}^n (X_i - \bar{X})^2}{\sigma^2}}_{=V_2} + \underbrace{\frac{n(\bar{X} - \mu)^2}{\sigma^2}}_{=V_3}.$$

## Proof of $(n-1)S^2/\sigma^2 \sim \chi_{n-1}^2$

Define  $V_1, V_2, V_3$  as follows:

$$\underbrace{\frac{\sum_{i=1}^n (X_i - \mu)^2}{\sigma^2}}_{=V_1} = \underbrace{\frac{\sum_{i=1}^n (X_i - \bar{X})^2}{\sigma^2}}_{=V_2} + \underbrace{\frac{n(\bar{X} - \mu)^2}{\sigma^2}}_{=V_3}.$$

► From the previous page,  $V_1 \sim \chi_n^2$  has MGF  $M_{V_1}(t) = (1 - 2t)^{-n/2}$

## Proof of $(n-1)S^2/\sigma^2 \sim \chi_{n-1}^2$

Define  $V_1, V_2, V_3$  as follows:

$$\underbrace{\frac{\sum_{i=1}^n (X_i - \mu)^2}{\sigma^2}}_{=V_1} = \underbrace{\frac{\sum_{i=1}^n (X_i - \bar{X})^2}{\sigma^2}}_{=V_2} + \underbrace{\frac{n(\bar{X} - \mu)^2}{\sigma^2}}_{=V_3}.$$

- ▶ From the previous page,  $V_1 \sim \chi_n^2$  has MGF  $M_{V_1}(t) = (1 - 2t)^{-n/2}$
- ▶  $\sqrt{n}(\bar{X} - \mu)/\sigma \sim N(0, 1) \Rightarrow V_3 \sim \chi_1^2$  with MGF  $M_{V_3}(t) = (1 - 2t)^{-1/2}$

## Proof of $(n-1)S^2/\sigma^2 \sim \chi_{n-1}^2$

Define  $V_1, V_2, V_3$  as follows:

$$\underbrace{\frac{\sum_{i=1}^n (X_i - \mu)^2}{\sigma^2}}_{=V_1} = \underbrace{\frac{\sum_{i=1}^n (X_i - \bar{X})^2}{\sigma^2}}_{=V_2} + \underbrace{\frac{n(\bar{X} - \mu)^2}{\sigma^2}}_{=V_3}.$$

- ▶ From the previous page,  $V_1 \sim \chi_n^2$  has MGF  $M_{V_1}(t) = (1 - 2t)^{-n/2}$
- ▶  $\sqrt{n}(\bar{X} - \mu)/\sigma \sim N(0, 1) \Rightarrow V_3 \sim \chi_1^2$  with MGF  $M_{V_3}(t) = (1 - 2t)^{-1/2}$
- ▶ Indep of  $V_2$  and  $V_3$  comes from the indep of  $S^2$  and  $\bar{X}$ .  
The MGF of  $V_1 = V_2 + V_3$  is thus

$$M_{V_1}(t) = M_{V_2}(t)M_{V_3}(t)$$

## Proof of $(n-1)S^2/\sigma^2 \sim \chi_{n-1}^2$

Define  $V_1, V_2, V_3$  as follows:

$$\underbrace{\frac{\sum_{i=1}^n (X_i - \mu)^2}{\sigma^2}}_{=V_1} = \underbrace{\frac{\sum_{i=1}^n (X_i - \bar{X})^2}{\sigma^2}}_{=V_2} + \underbrace{\frac{n(\bar{X} - \mu)^2}{\sigma^2}}_{=V_3}.$$

- ▶ From the previous page,  $V_1 \sim \chi_n^2$  has MGF  $M_{V_1}(t) = (1 - 2t)^{-n/2}$
- ▶  $\sqrt{n}(\bar{X} - \mu)/\sigma \sim N(0, 1) \Rightarrow V_3 \sim \chi_1^2$  with MGF  $M_{V_3}(t) = (1 - 2t)^{-1/2}$
- ▶ Indep of  $V_2$  and  $V_3$  comes from the indep of  $S^2$  and  $\bar{X}$ .

The MGF of  $V_1 = V_2 + V_3$  is thus

$$M_{V_1}(t) = M_{V_2}(t)M_{V_3}(t) \Rightarrow M_{V_2}(t) = \frac{M_{V_1}(t)}{M_{V_3}(t)} = \frac{(1 - 2t)^{-n/2}}{(1 - 2t)^{-1/2}} = (1 - 2t)^{-\frac{n-1}{2}},$$

which is the MGF for  $\chi_{n-1}^2$ . By the uniqueness of MGFs, this proves

$$V_2 = \frac{\sum_{i=1}^n (X_i - \bar{X})^2}{\sigma^2} = \frac{(n-1)S^2}{\sigma^2} \sim \chi_{n-1}^2.$$

## Back to Statistics

Recall our goal is to **estimate the unknown mean  $\mu$**  using the observed values of  $X_1, X_2, \dots, X_n$  that are i.i.d.  $\sim N(\mu, \sigma^2)$ .

## Back to Statistics

Recall our goal is to **estimate the unknown mean  $\mu$**  using the observed values of  $X_1, X_2, \dots, X_n$  that are i.i.d.  $\sim N(\mu, \sigma^2)$ .

For  $Z \sim N(0, 1)$ , using the normal CDF we know

$$P(-1.96 \leq Z \leq 1.96) = 0.95.$$

As  $\bar{X} \sim N(\mu, \sigma^2/n)$ , which implies  $\frac{\bar{X}-\mu}{\sigma/\sqrt{n}} \sim N(0, 1)$ , we have

$$P\left(-1.96 \leq \frac{\bar{X} - \mu}{\sigma/\sqrt{n}} \leq 1.96\right) = 0.95,$$

or equivalently

$$P\left(\bar{X} - 1.96 \frac{\sigma}{\sqrt{n}} \leq \mu \leq \bar{X} + 1.96 \frac{\sigma}{\sqrt{n}}\right) = 0.95.$$

This means, for 95% of the time, the sample mean  $\bar{X}$  is within  $1.96\sigma/\sqrt{n}$  from the true value of  $\mu$ , but  **$\sigma$  is UNKNOWN**.

## t-Statistic

The result on the previous page relies on the fact that

$$\frac{\bar{X} - \mu}{\sigma/\sqrt{n}} = \frac{\bar{X} - \mu}{\sqrt{\sigma^2/n}} \sim N(0, 1), \quad \text{but } \sigma^2 \text{ is UNKNOWN.}$$

If we replace  $\sigma^2$  by  $S^2$ , what's the distribution of

$$T = \frac{\bar{X} - \mu}{\sqrt{S^2/n}}?$$

The random variable  $T$  defined above is called the *t-statistic*.



## t-Statistic (2)

Dividing both the numerator and denominator of  $T$  by  $\sqrt{\sigma^2/n}$ , we can rewrite  $T$  as

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

where

1.  $Z = \frac{(\bar{X} - \mu)}{\sqrt{\sigma^2/n}} \sim N(0, 1)$
2.  $U = \frac{(n-1)S^2}{\sigma^2} \sim \chi_{n-1}^2$ , and
3.  $Z$  and  $U$  are independent (from the indep of  $\bar{X}$  and  $S^2$ ).

From the definition of  $t$ -distribution, we know

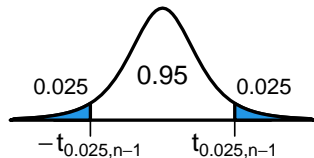
$$T = \frac{\bar{X} - \mu}{\sqrt{S^2/n}} \sim t_{n-1}$$

has a  $t$ -**distribution** with  $n - 1$  degrees of freedom.

## 95% One-Sample $t$ -Confidence Interval

If  $T \sim t_{n-1}$ , let  $t_{0.025, n-1}$  be the value so that

$$P(-t_{0.025, n-1} \leq T \leq t_{0.025, n-1}) = 0.95$$



This means

$$P\left(-t_{0.025, n-1} \leq T = \frac{\bar{X} - \mu}{\sqrt{S^2/n}} \leq t_{0.025, n-1}\right) = 0.95,$$

or equivalently

$$P\left(\bar{X} - t_{0.025, n-1} \sqrt{\frac{S^2}{n}} \leq \mu \leq \bar{X} + t_{0.025, n-1} \sqrt{\frac{S^2}{n}}\right) = 0.95.$$

meaning, for 95% of the time, the sample mean  $\bar{X}$  is within  $t_{0.025, n-1} \sqrt{\frac{S^2}{n}}$  from the true value of  $\mu$ .

## 95% One-Sample $t$ -Confidence Interval

The interval

$$\left( \bar{X} - t_{0.025, n-1} \sqrt{\frac{S^2}{n}}, \bar{X} + t_{0.025, n-1} \sqrt{\frac{S^2}{n}} \right).$$

is thus call the *95% one-sample  $t$ -confidence interval* for  $\mu$ .

For example, with  $n = 16$  observations,  $t_{0.025, 16-1} \approx 2.131$ , the 95% confidence interval for  $\mu$  is

$$\left( \bar{X} - 2.131 \sqrt{\frac{S^2}{16}}, \bar{X} + 2.131 \sqrt{\frac{S^2}{16}} \right).$$