



# STATISTICS FOR DATA SCIENCE

## Confidence Intervals for Small Samples

---

**Prof. Uma D**

**Prof. Suganthi S**

**Prof. Silviya Nancy J**

Department of Computer Science and Engineering

# STATISTICS FOR DATA SCIENCE

---

## Confidence Intervals for Small Samples

Prof. Uma D

Prof. Suganthi S

Prof. Silviya Nancy J

## Topics to be covered...

---

- **Confidence Intervals for population mean of small samples**
- **Student's t Distribution**
- **Confidence Intervals using t Distribution**
- **Student's t Distribution Is Appropriate?**
- **One-Sided CI for Small Samples**

- If the sample size is small, standard deviation ( $s$ ) of the sample may not be close to  $\sigma$  (population standard deviation). Hence  $\bar{X}$  (sample\_mean) may not be approximately normal.
- However, if the population from which the sample is drawn is known to be approximately normal (can be confirmed using normal probability plot).

- It turns out that we can still use the quantity.

- $(\bar{X} - \mu) / (s/\sqrt{n}),$

but since  $s$  is not necessarily close to  $\sigma$ , the quantity will not have a normal distribution.

- Instead it has Student's  $t$  distribution with  $n - 1$  degrees of freedom, denoted as  $t_{n-1}$ .

- The t distribution is a theoretical probability distribution.
- It is symmetrical, bell-shaped, and similar to the standard normal curve.
- It differs from the standard normal curve, however, in that it has an additional parameter, called **degrees of freedom**, which changes its shape.

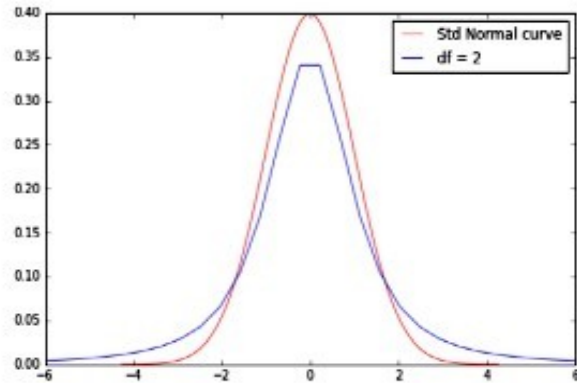
$$\text{df} = \text{sample size} - 1$$

- Setting the value of df defines a particular member of the family of t distributions. ( $\text{df} > 0 \Rightarrow \text{Sample Size} > 1$ )

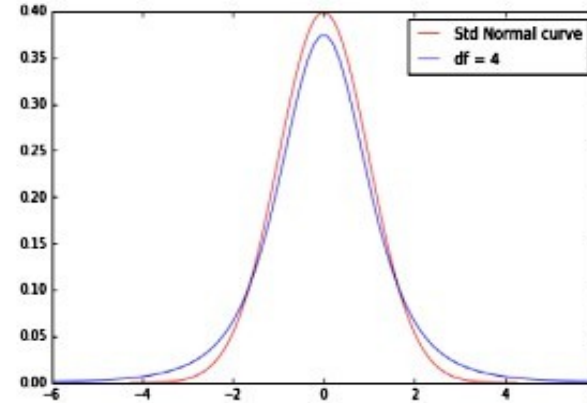
# STATISTICS FOR DATA SCIENCE

## Students t Distribution

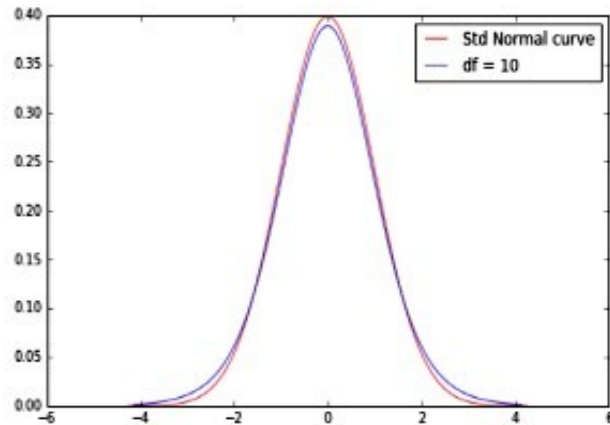
1)  $df = 2$



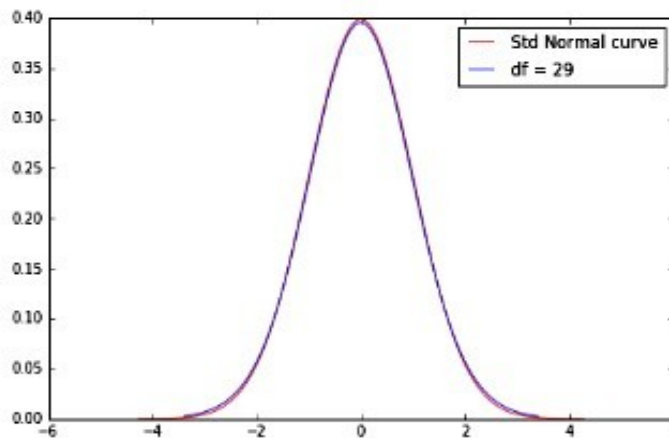
2)  $df = 4$



3)  $df = 10$



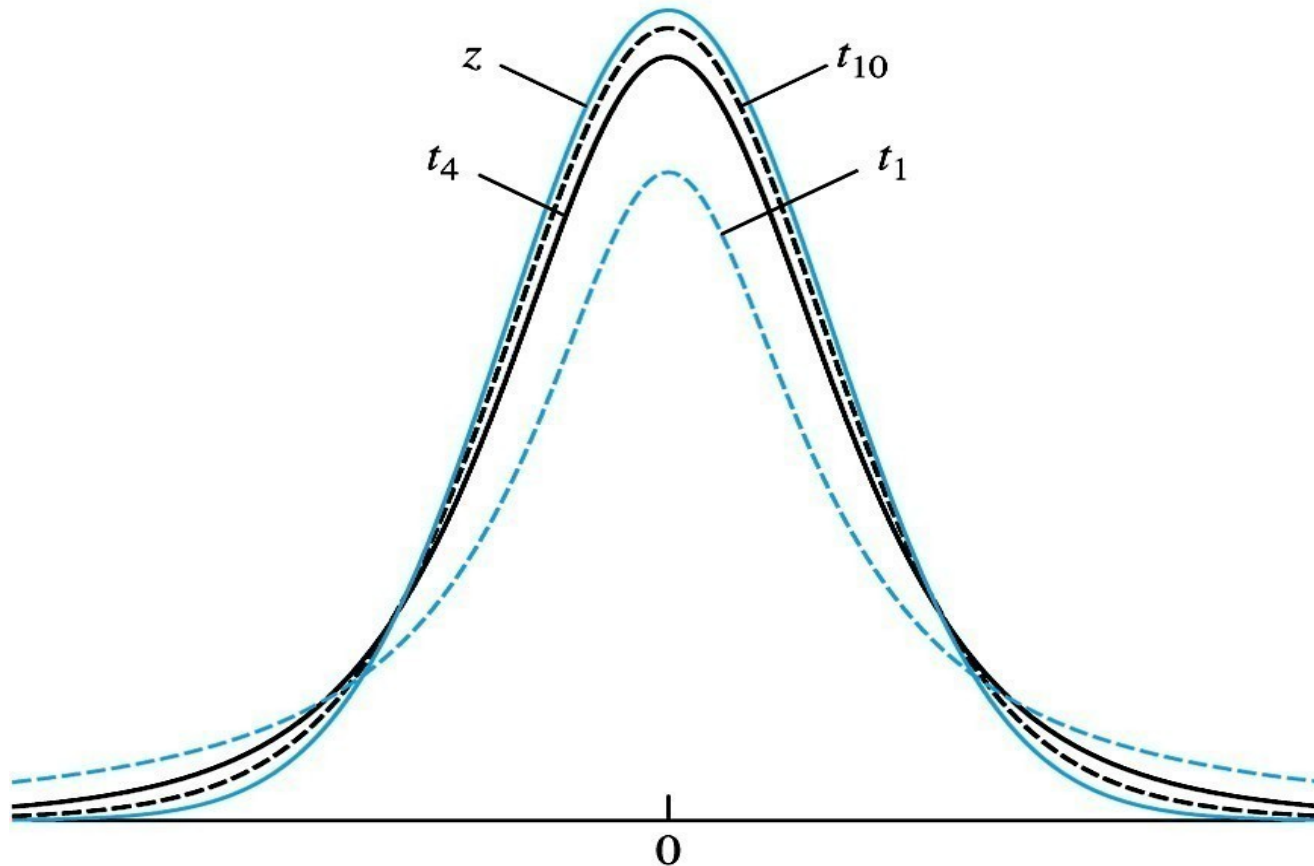
4)  $df = 30$



# STATISTICS FOR DATA SCIENCE

## PDF for Students t curve

Note that the smaller the distribution function, the flatter the shape of the distribution, resulting in greater area in the tails of the distribution.





As the df increase, the t distribution approaches the standard normal distribution ( $\mu=0.0$ ,  $\sigma=1.0$ ).

The standard normal curve is a special case of the t distribution when  $df = \text{infinity}$ .

For practical purposes, the t distribution approaches the standard normal distribution relatively quickly, such that when  $df=30$  the two are almost identical.

- We use t table to find probabilities associated with t distribution.
- Row headings – denotes degree of freedom
- Column headings – denotes the area to the right(probabilities)
- The value in particular row and column specifies the t-score where,

$$P(t > t\text{-score}) = \text{col\_heading}$$

1) A random sample of size 10 is drawn from a normal distribution with mean 4.

a) Find  $P(t > 1.833)$

b) Find  $P(t > 1.5)$

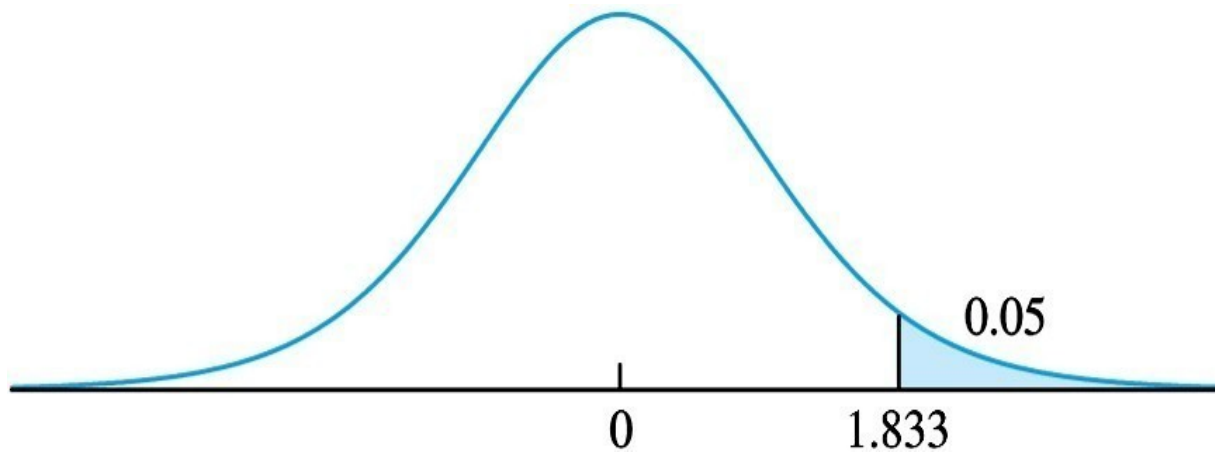
a) Find  $P(t > 1.833)$

df = 9 (row\_heading)

t-score = 1.833

corresponding col\_heading = 0.05

**$P(t > 1.833) = 0.05$**





### b) Find $P(t > 1.5)$

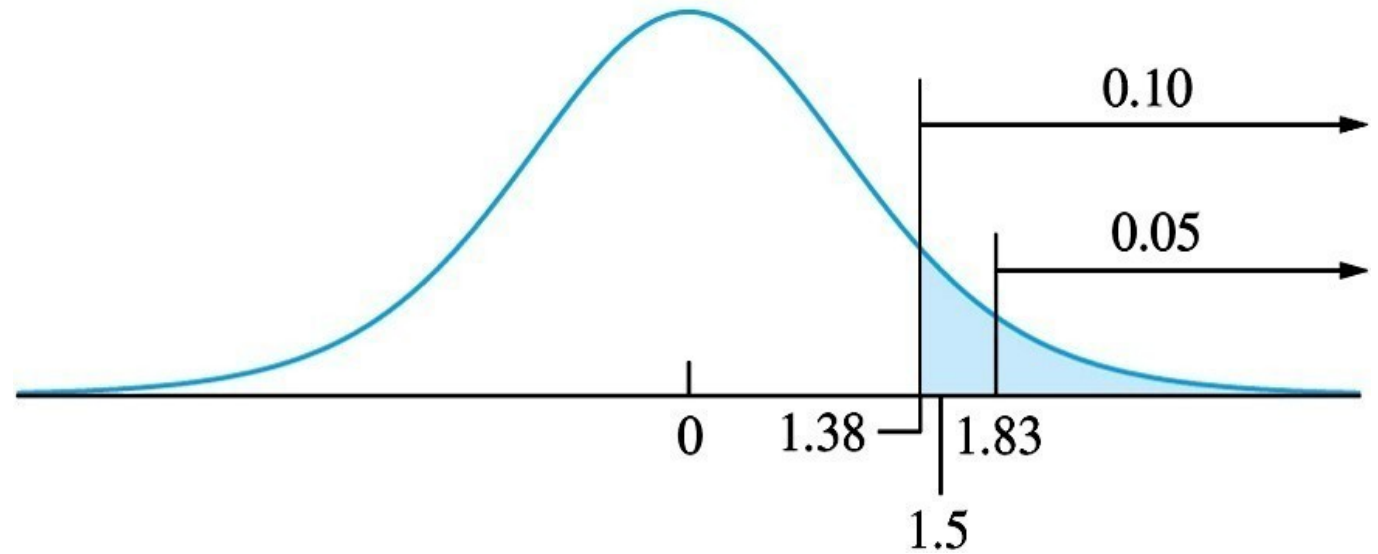
df = 9 (row\_heading)

t-score = 1.5 [does not correspond to any of the values in that row]

but we do have t-scores 1.383, 1.833 corresponding to upper tail probabilities 0.10 and 0.05 respectively. That is,

$P(t > 1.383) = 0.10$  and  $P(t > 1.833) = 0.05$

Since  $1.383 < 1.5 < 1.833 \Rightarrow 0.05 < P(t > 1.5) < 0.10$



2) Find the value of  $t_{12}$  distribution where upper-tail probability is 0.025.

**Solution:**

row\_head = 12

col\_head = 0.025

=> t-score = 2.179

The quantity,

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

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

We can generate a  $(1 - \alpha)$  100% Confidence Interval for  $\mu$  as

$$\bar{X} \pm t_{n-1, \alpha/2} \frac{s}{\sqrt{n}}$$

## Student's $t$ Distribution is Appropriate when

---



- Sample size is small ( $n < 30$ )
- Sample comes from a population that is approximately normal.
- In many cases, we must examine the sample for normality, by constructing a box plot or normal probability plot.
- Unfortunately, when the sample size is small, departures from normality may be hard to detect.
- If these plots do not reveal a strong asymmetry or any outliers, then in most cases the Student's  $t$  distribution will be reliable.



We can generate a  $(1 - \alpha)$  100% Upper Confidence bound for  $\mu$  as:

$$\bar{X} + t_{n-1, \alpha} * s/\sqrt{n}$$

We can generate a  $(1 - \alpha)$  100% Lower Confidence bound for  $\mu$  as:

$$\bar{X} - t_{n-1, \alpha} * s/\sqrt{n}$$

## Example1

---



Find the value of  $t_{n-1, \alpha/2}$  needed to construct a two-sided confidence interval of the given level with the given sample size:

- a) 90% with sample size 12
- b) 95% with sample size 7

**a) 90% with sample size 12**

$$df = 11$$

$$\alpha = 0.10 \quad \Rightarrow \alpha/2 = 0.05$$

$$\Rightarrow \text{in t table : row\_heading} = 11, \text{col\_heading} = 0.05 \Rightarrow t_{11, 0.05} = 1.796$$

**b) 95% with sample size 7**

$$df = 6$$

$$\alpha = 0.05 \quad \Rightarrow \alpha/2 = 0.025$$

$$\Rightarrow \text{in t table : row\_heading} = 6, \text{col\_heading} = 0.025 \Rightarrow t_{6, 0.025} = 2.447$$

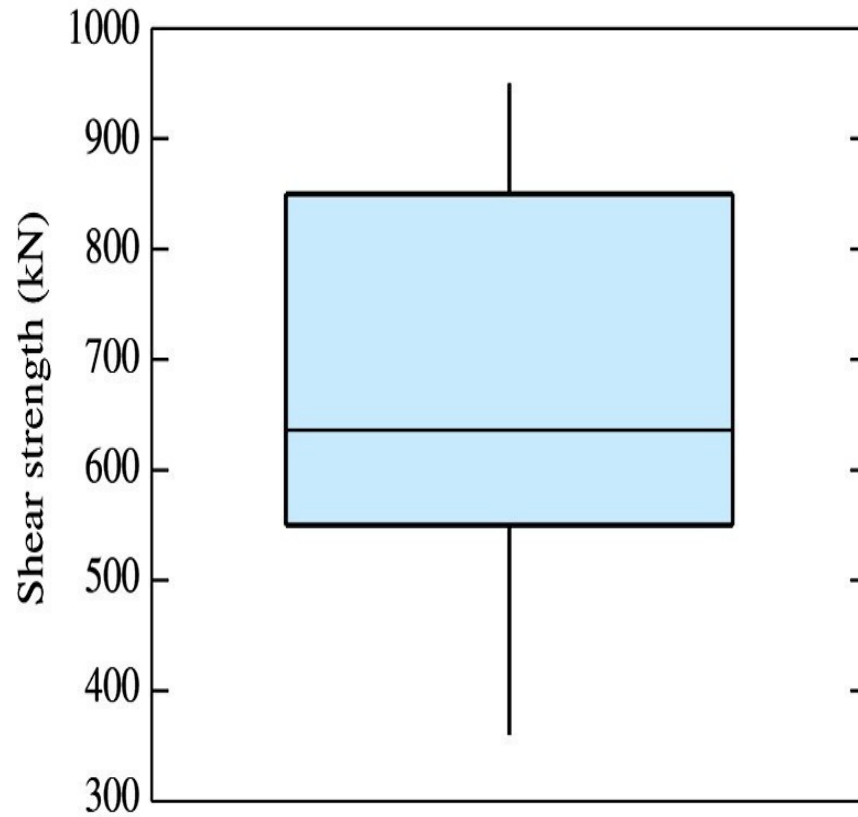
Example2



Following represents the measurements of the nominal shear strength (in kN) for a sample of 15 prestressed concrete beams:

|     |     |     |     |     |     |     |     |
|-----|-----|-----|-----|-----|-----|-----|-----|
| 580 | 400 | 428 | 825 | 850 | 875 | 920 | 550 |
| 575 | 750 | 636 | 360 | 590 | 735 | 950 |     |

- a)Is it appropriate to use the Student's t statistic to construct a 99% confidence interval for the mean shear strength?
- b)If so, construct the confidence interval. If not, explain why not.



Since there are no outliers in the data set, Student's  $t$  statistic can be used to construct 99% CI.

Sample mean =  $\bar{X}$  = 668.27

Sample standard deviation =  $s$  = 192.089

$t_{n-1, \alpha/2} = t_{15-1, 0.005} = 2.977$

99% CI :

$$668.27 \pm 2.977 * 192.089/\sqrt{15}$$

$$=(520.62, 815.92)$$

If it is known that the sample indeed was drawn from a **normal population**, also the **standard deviation of the population is known**, use z not t distribution to find out the confidence interval irrespective of the sample size.

### Summary

Let  $X_1, \dots, X_n$  be a random sample (of any size) from a *normal* population with mean  $\mu$ . If the standard deviation  $\sigma$  is known, then a level  $100(1 - \alpha)\%$  confidence interval for  $\mu$  is

$$\bar{X} \pm z_{\alpha/2} \frac{\sigma}{\sqrt{n}} \quad (5.12)$$



**THANK YOU**

---

**Prof. Uma D**

**Prof. Suganthi S**

**Prof. Silviya Nancy J**

Department of Computer Science and Engineering