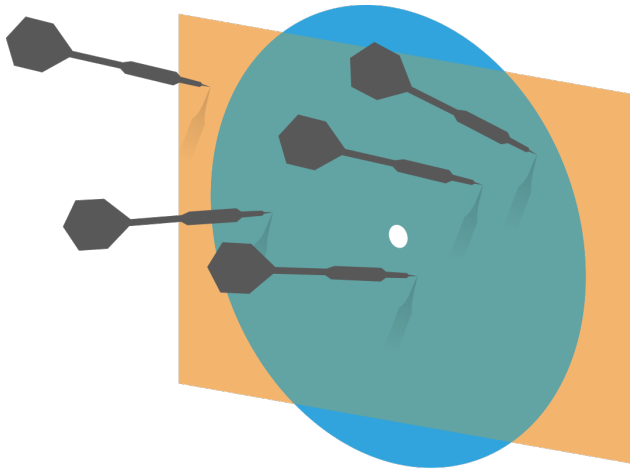


Probability and Statistics for Computer Science

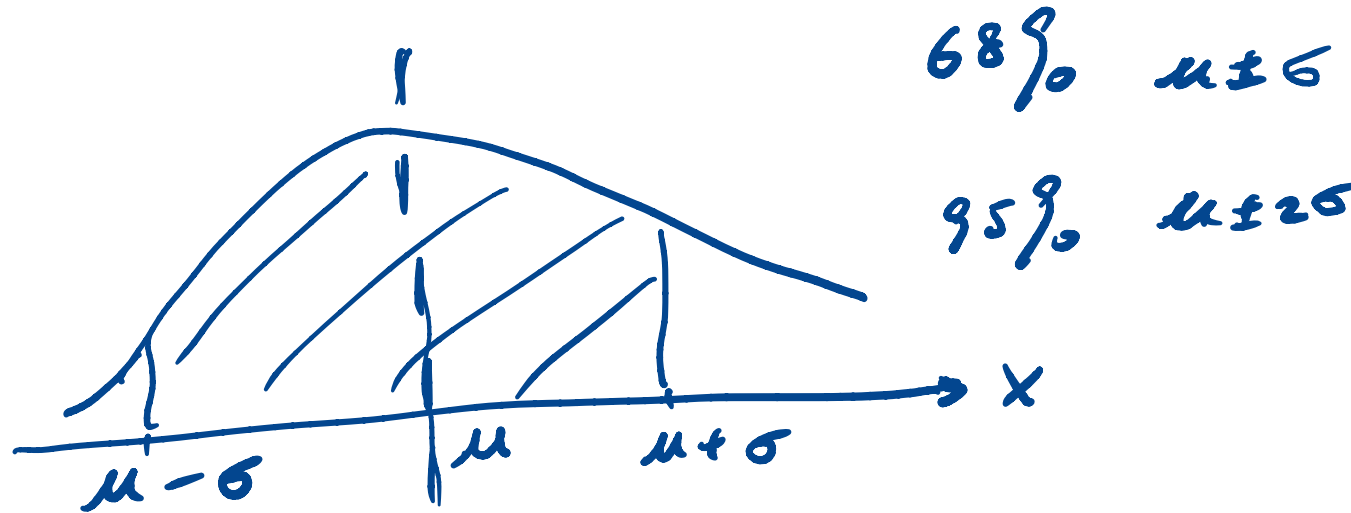


“In statistics we apply probability
to draw conclusions from data.”
---Prof. J. Orloff

Credit: wikipedia

Last time

- ✱ Exponential distribution
- ✱ Normal (Gaussian) distribution



Objectives


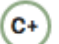
- ✱ Sample mean
- ✱ confidence interval
- ✱ t-distribution

Motivation for drawing conclusion from samples

- ✿ In a study of new-born babies' health, random samples from different time, places and different groups of people will be collected to see how the overall health of the babies is like.



Motivation of sampling: the poll example

	DATES	POLLSTER	SAMPLE	RESULT	NET RESULT
 U.S. Senate	Miss. NOV 25, 2018	 Change Research	1,211 LV	Espy 46% 51% Hyde-Smith	Hyde-Smith +5

Source: [FiveThirtyEight.com](https://www.forty8.com)

- ✱ This senate election poll tells us:
 - ✱ The sample has 1211 likely voters
 - ✱ Ms. Hyde-Smith has realized sample mean equal to 51%
- ✱ What is the estimate of the percentage of votes for Hyde-smith?
- ✱ How confident is that estimate?

Population

- ✱ What is a population?
 - ✱ It's the entire possible data set $\{X\}$
 - ✱ It has a countable size N_p
 - ✱ The population mean $popmean(\{X\})$ is a number
 - ✱ The population standard deviation is $popstd(\{X\})$ and is also a number
- ✱ The population mean and standard deviation are the same as defined previously in chapter 1

Population

$$\{X\} = \{1, 2, 3, \dots, 12\} \quad N_p = 12$$

$$\text{popmean}(\{X\}) = ? \quad 6.5$$

$$\text{popstd}(\{X\}) = ?$$

$$\sqrt{\frac{\sum (x_i - 6.5)^2}{12}}$$

3.4...

Sample

- ✱ The sample is a random subset of the population and is denoted as $\{x\}$, where sampling is done with **replacement**
- ✱ The sample size N is assumed to be much less than population size N_p
- ✱ The **sample mean of a population** is $\bar{X}^{(N)}$ and is a **random variable**

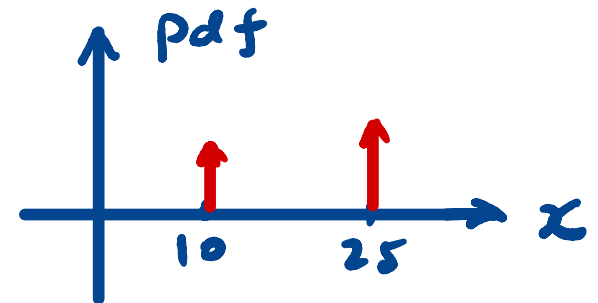
Sample mean example



Money box

* shake and take one and put back.

10¢ dime
25¢ quarter



X_1 takes $x_1=10$ $E[X]=?$

X_2 takes $x_2=10$

X_3 takes $x_3=25$

\vdots
 X_N takes $x_N=10$

$$\boxed{\bar{X}} = \frac{\sum X_i}{N}$$

\parallel
 $X^{(N)}$

Sample $\{x\}$ and Sample Mean $X^{(N)}$

$$\{X\} = \{1, 2, 3, \dots, 12\} \quad N < N_p$$

One random sample $\rightarrow \{x\} = \{1, 1, 2, 3, 3\} \quad N = 5$

$X^{(N)}$ RV takes a value?

$$\frac{10}{5} = 2$$

$$X^{(N)} = \frac{x_1 + x_2 + \dots + x_N}{N}$$

Another random sample $\rightarrow \{1, 1, 1, 1, 1\} \Rightarrow X^{(N)} = 1$

Sample mean of a population

- ✱ The sample mean is the average of **IID** samples

$$\underline{X^{(N)}} = \frac{1}{N} (\mathbf{x}_1 + \mathbf{x}_2 + \dots + \mathbf{x}_N) = \text{mean}(\{x\})$$

- ✱ By linearity of the expectation and the fact the sample items are identically drawn from the same population with replacement

$$E[X^{(N)}] = \frac{1}{N} (E[X^{(1)}] + E[X^{(1)}] \dots + E[X^{(1)}]) = E[X^{(1)}]$$

Expected value of one random sample is the population mean

- ✱ Since each sample is drawn uniformly from the population

$$\begin{aligned} E[X^{(1)}] &= \text{popmean}(\{X\}) \\ &= \sum x_i \frac{1}{N_p} = \text{popmean} \end{aligned}$$

therefore $E[X^{(N)}] = \text{popmean}(\{X\})$

- ✱ We say that $X^{(N)}$ is an unbiased estimator of the population mean.

Standard deviation of the sample mean

- ✱ We can also rewrite another result from the lecture on the weak law of large numbers

$$\text{var}[X^{(N)}] = \frac{\text{popvar}(\{X\})}{N} = \frac{\text{var}(\sum RV_i)}{N}$$

if RV_i are indpt.

- ✱ The standard deviation of the sample mean

$$\text{std}[X^{(N)}] = \frac{\text{popstd}(\{X\})}{\sqrt{N}} = \frac{\text{var}(\frac{1}{N} \cdot \sum x_i)}{(\frac{1}{N})^2 \text{var}(\sum x_i)}$$

- ✱ But we need the population standard deviation in order to calculate the $\text{std}[X^{(N)}]$!

$$\frac{1}{N^2} \cdot N \cdot \text{var}(x_i)$$
$$\frac{1}{N} \cdot \text{var}(x_i)$$

Unbiased estimate of population standard deviation & Stderr

- ✱ The unbiased estimate of $popstd(\{X\})$ is defined as

$$stdunbiased(\{x\}) = \sqrt{\frac{1}{N-1} \sum_{x_i \in \text{sample}} (x_i - \text{mean}(\{x_i\}))^2}$$

std({x}) → popstd

- ✱ So the **standard error** is an estimate of

$$std[X^{(N)}] \quad \text{approximation} \quad std[X^{(N)}] = \frac{popstd(\{X\})}{\sqrt{N}}$$

approx. of popstd

$$\frac{popstd(\{X\})}{\sqrt{N}} \stackrel{\bullet}{=} \frac{stdunbiased(\{x\})}{\sqrt{N}} = \boxed{stderr(\{x\})}$$

Standard error: election poll

	DATES	POLLSTER	SAMPLE	RESULT	NET RESULT
U.S. Senate	Miss. NOV 25, 2018	C+ Change Research	1,211 LV	Espy 46% 51% Hyde-Smith	Hyde-Smith +5

What is the estimate of the percentage of votes for Hyde-Smith? 51%

$$X_i = \begin{cases} 1 & \text{for Smith} \\ 0 & \text{otherwise} \end{cases}$$

Number of sampled voters who selected Ms. Smith is:

$$1211(0.51) \approx 618$$

$$\sum X_i = \text{total votes for Smith}$$

Number of sampled voters who didn't selected Ms. Smith was

$$1211(0.49) \approx 593$$

$$\frac{\sum X_i}{N} \quad N = 1211$$

Standard error: election poll

✱ $stdunbiased(\{x\})$

1 1 0 0 ... 1
N = 1211

$$= \sqrt{\frac{1}{1211 - 1} (618(1 - 0.51)^2 + 593(0 - 0.51)^2)} = 0.5001001$$

$$X^{(N)} = \frac{x_1 + \dots + x_N}{N}$$

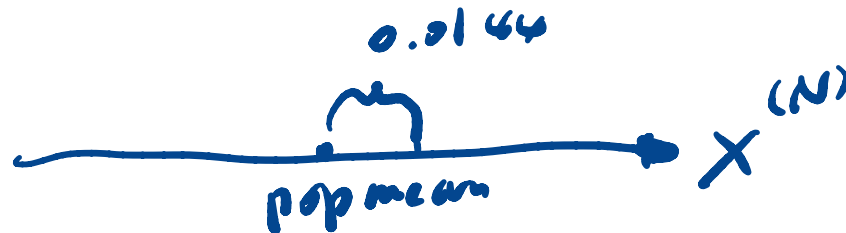
618 "1"

✱ $stderr(\{x\})$

for Smith

$$= \frac{0.5001001}{\sqrt{1211}} \simeq 0.0144$$

593 "0"
not for her



Interpreting the standard error

- ✱ **Sample mean** is a random variable and has its own probability distribution, `stderr` is an estimate of sample mean's standard deviation $X^{(N)} \rightarrow \text{Normal}$ if $N \rightarrow \infty$
- ✱ When N is very large, according to the **Central Limit Theorem**, sample mean is approaching a normal distribution with μ σ

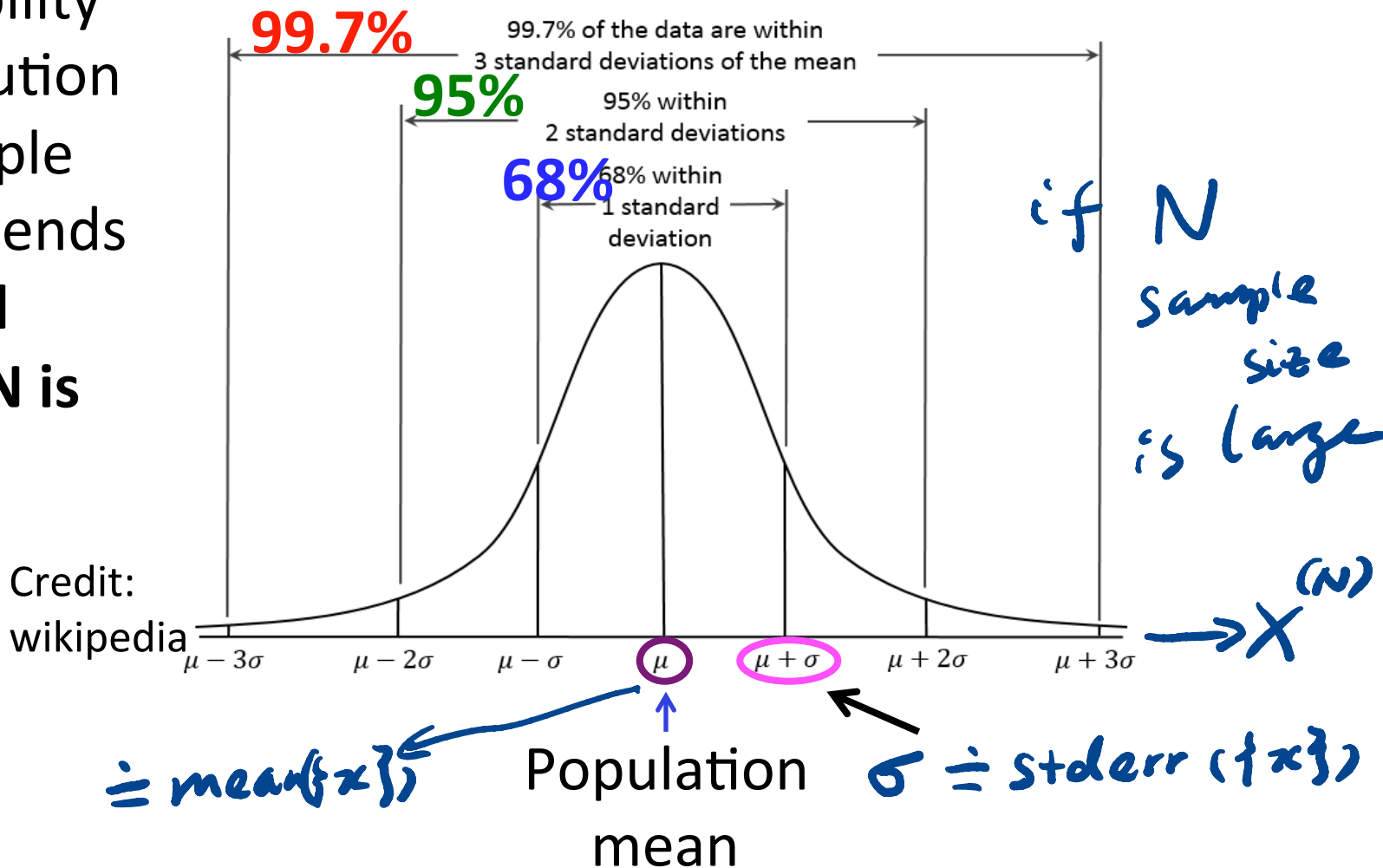
$$\mu = \text{popmean}(\{X\}) ; \sigma = \frac{\text{popstd}(\{X\})}{\sqrt{N}} \doteq \text{stderr}(\{x\})$$

$$\text{stderr}(\{x\}) = \frac{\text{stdunbiased}(\{x\})}{\sqrt{N}}$$

$$\text{pdf} = P(x) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

Interpreting the standard error

Probability distribution of sample mean tends normal when N is large



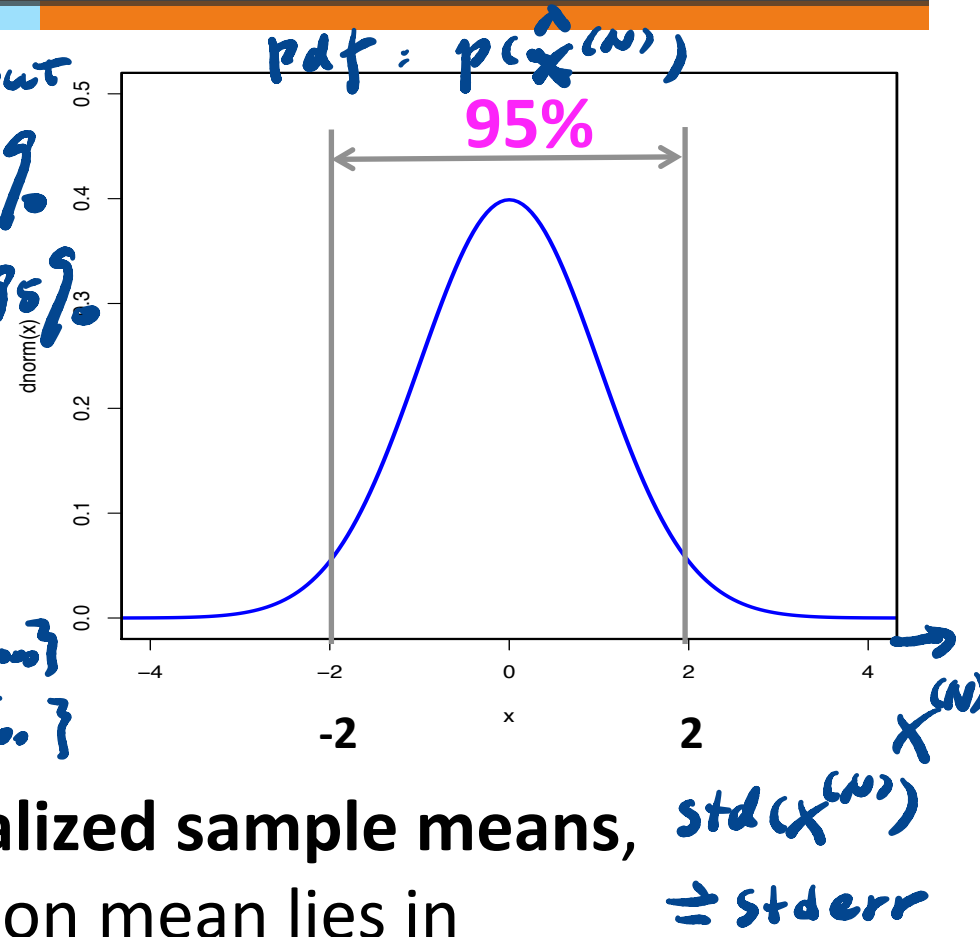
Confidence intervals

✱ Confidence interval for a population mean is defined by fraction

input
c%
c% = 95%

✱ Given a percentage, find how many units of strerr it covers.

$N=1000$ { x_1, \dots, x_{1000} }
 $N=1000$ { x'_1, \dots, x'_{1000} }



For **95%** of the **realized sample means**,
 the population mean lies in
 [sample mean - 2 stderr, sample mean + 2 stderr]

Confidence intervals when N is large

- ✱ For about 68% of realized sample means

$$\text{mean}(\{x\}) - \text{stderr}(\{x\}) \leq \text{popmean}(\{X\}) \leq \text{mean}(\{x\}) + \text{stderr}(\{x\})$$

- ✱ For about 95% of realized sample means

$$\text{mean}(\{x\}) - 2\text{stderr}(\{x\}) \leq \text{popmean}(\{X\}) \leq \text{mean}(\{x\}) + 2\text{stderr}(\{x\})$$


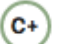
- ✱ For about 99.7% of realized sample means

$$\text{mean}(\{x\}) - 3\text{stderr}(\{x\}) \leq \text{popmean}(\{X\}) \leq \text{mean}(\{x\}) + 3\text{stderr}(\{x\})$$

Q. Confidence intervals

- ✱ What is the 68% confidence interval for a population mean?
 - A. [sample mean-2stderr, sample mean+2stderr]
 - B. [sample mean-stderr, sample mean+stderr]
 - C. [sample mean-std, sample mean+std]

Standard error: election poll



	DATES	POLLSTER	SAMPLE	RESULT	NET RESULT
 U.S. Senate	Miss. NOV 25, 2018	 Change Research	1,211 LV	Espy 46% 51% Hyde-Smith	Hyde-Smith +5

51%

✱ We estimate the population mean as 51% with stderr 1.44%

✱ The 95% confidence interval is
[51%-2×1.44%, 51%+2×1.44%]= [48.12%, 53.88%]

Q.

✱ A store staff mixed their fuji  and gala  apples and they were individually wrapped, so they are indistinguishable. if I pick 30 apples and found 21 fuji , what is my 95% confidence interval to estimate the popmean is 70% for fuji? (hint: $\text{strerr} > 0.05$)

A. $[0.7-0.17, 0.7+0.17]$

B. $[0.7-0.056, 0.7+0.056]$

What if N is small? When is N large enough?

- * If samples are taken from normal distributed population, the following variable is a random variable whose distribution is Student's t-distribution with **N-1** degree of freedom.

$$T = \frac{\text{mean}(\{x\}) - \text{popmean}(\{X\})}{\text{stderr}(\{x\})}$$

$E[T] = 0$
 $\text{var}[T] = 1$

$\rightarrow X^{(N)}$ $E[X^{(N)}]$ $X^{(N)}$ $x_1 \dots x_N$

$\approx \rightarrow \text{std}(X^{(N)})$ N is the sample size

Degree of freedom is **N-1** due

to this constraint:
$$\sum_i (x_i - \text{mean}(\{x\})) = 0$$

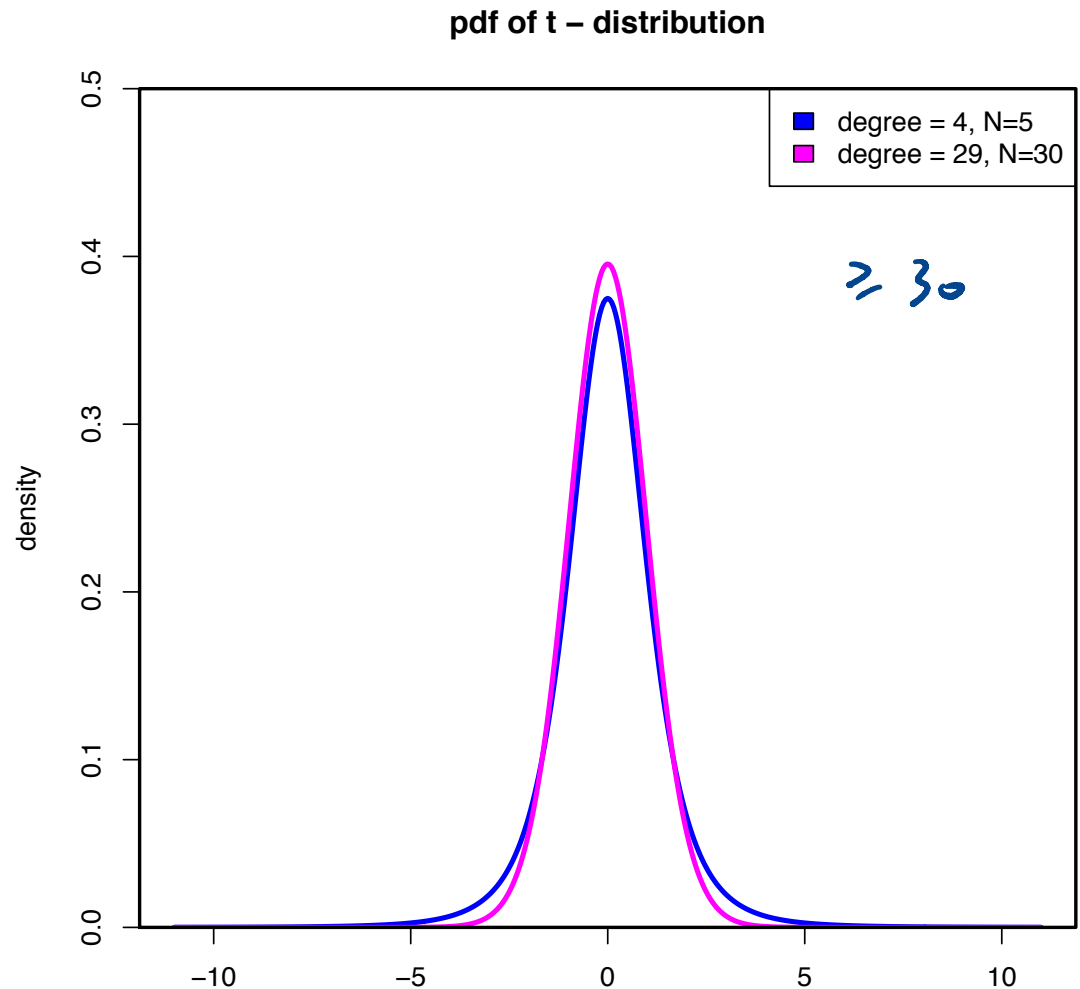
t-distribution is a family of distri. with different degrees of freedom

t-distribution with $N=5$
and $N=30$



Credit :
wikipedia

William Sealy Gosset 1876-1937

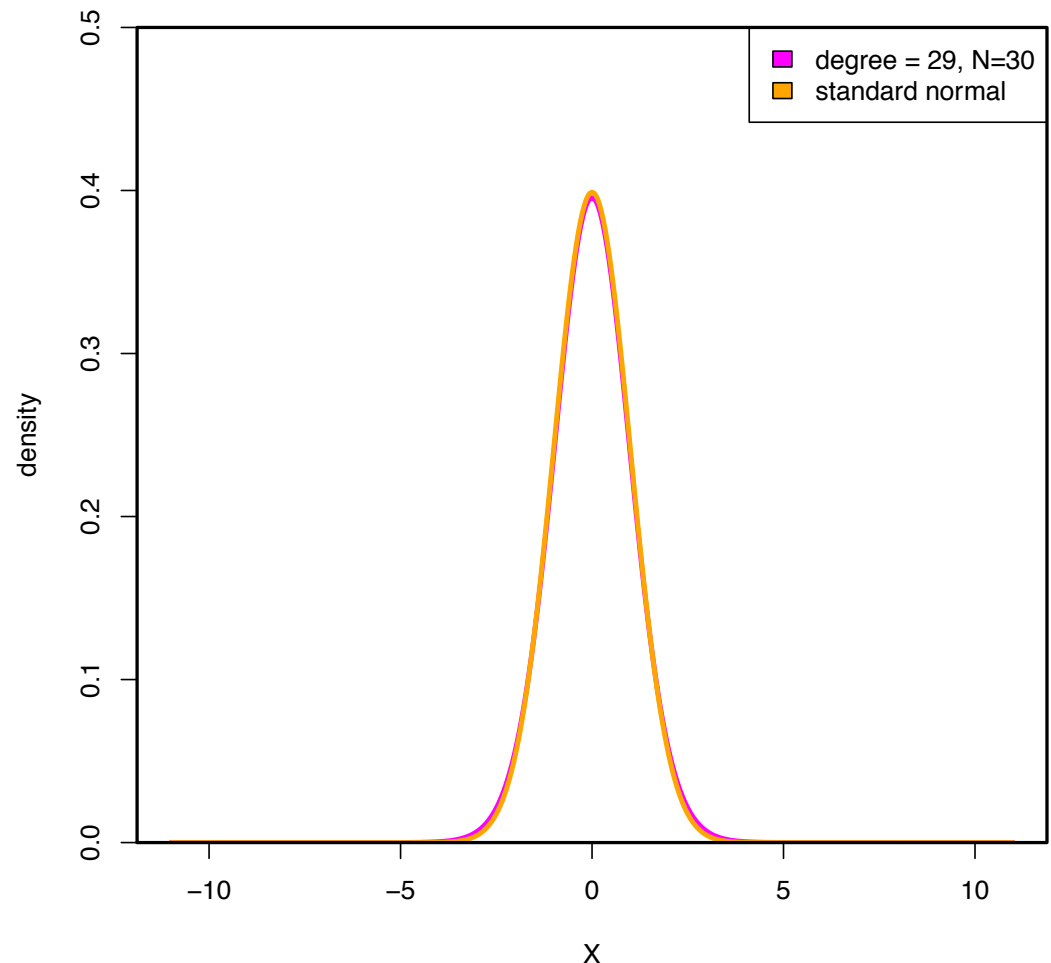


When $N=30$, t-distribution is almost Normal

t-distribution looks very similar to normal when $N=30$.

So $N=30$ is a rule of thumb to decide N is large or not

pdf of t (n=30) and normal distribution



Assignments

- ✱ Read Chapter 7 of the textbook
- ✱ Next time: Bootstrap, Hypothesis tests
- ✱ Prepare for Midterm1

Additional References

- ✱ Charles M. Grinstead and J. Laurie Snell
"Introduction to Probability"
- ✱ Morris H. Degroot and Mark J. Schervish
"Probability and Statistics"

See you next time

*See
you!*

