

Class 3b: Review of concepts in Probability and Statistics

Business Forecasting

Confidence Intervals

Confidence Intervals

- We calculated the mean price in our sample
- How confident are we that our estimate is close to the parameter's value?
- Confidence intervals measure uncertainty around the estimate

Confidence Intervals

- Mean price was 1245.43
- Is it reasonable to think true average price in population is 1100? What about 2000?
- Suppose that we calculated the confidence interval to be:

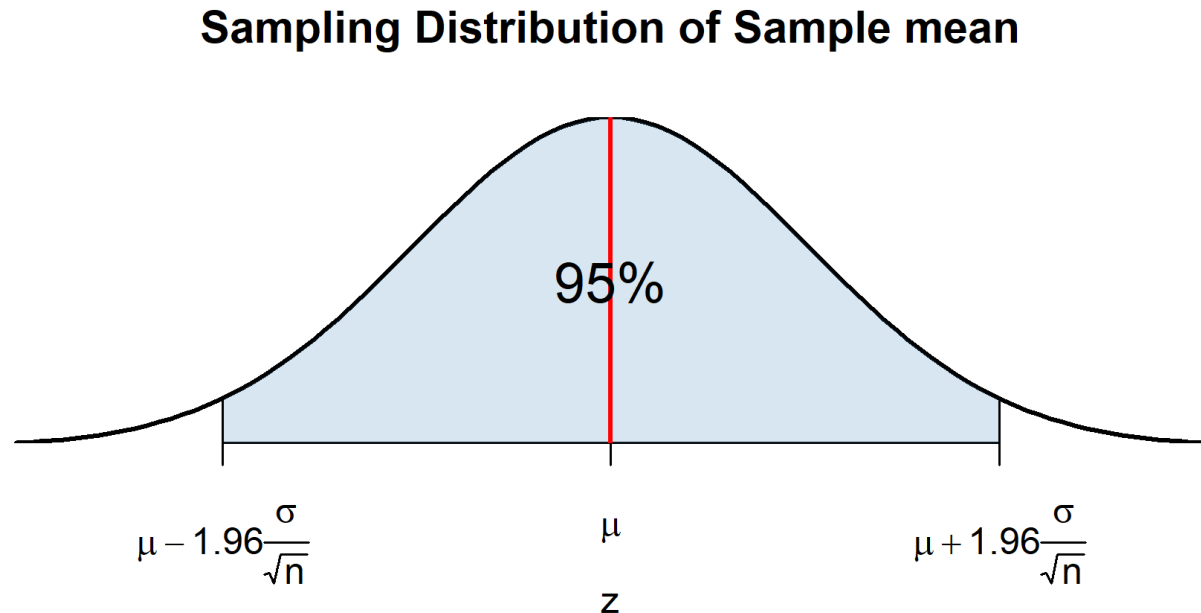
$$\{1086.64, 1404.22\}$$

- Where are these numbers coming from?
1. The sampling distribution of the sample mean tells us how likely we are to get a point estimate which is far away from the true mean
 2. The confidence interval uses this property of the sampling distribution to tell us where the true mean might be
 - Let's go through these statements 1-by-1

Sampling distribution

Q: How likely is it that a sample mean is far away from the true mean?

- Consider a hypothetical sampling distribution of a sample mean
 - Reminder: $\bar{x} \sim \mathcal{N}(\mu, \frac{\sigma}{\sqrt{n}})$
 - If we draw samples repeatedly, 95% of their means will be within the shaded area
- Why 1.96?



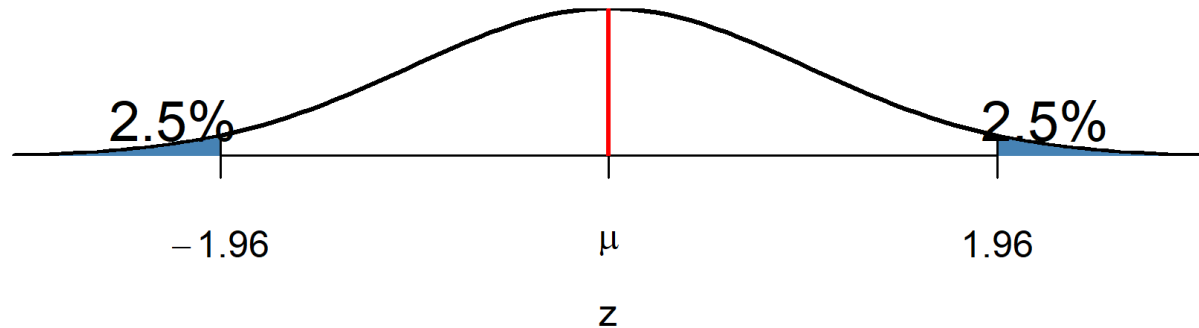
Sampling distribution

- Assume we have $n > 30$
 - CLT applies and $\bar{X} \sim N(\mu_X, \frac{\sigma_X}{\sqrt{n}})$
- We want to find k_1 and k_2 , such that:
 - $P(k_1 < \bar{X} < k_2) = 0.95$, so
 - $P(\bar{X} < k_1) = 0.025$ and $P(\bar{X} > k_2) = 0.025$ (or $P(\bar{X} < k_2) = 0.975$)
 - Trick is to standardize the variable:

$$P(\bar{X} < k_2) = P(\bar{X} - \mu_X < k_2 - \mu_X) = P(\underbrace{\frac{\bar{X} - \mu_X}{\frac{\sigma_X}{\sqrt{n}}}}_{Z \sim N(0,1)} < \underbrace{\frac{k_2 - \mu_X}{\frac{\sigma_X}{\sqrt{n}}}}_{k'_2})$$

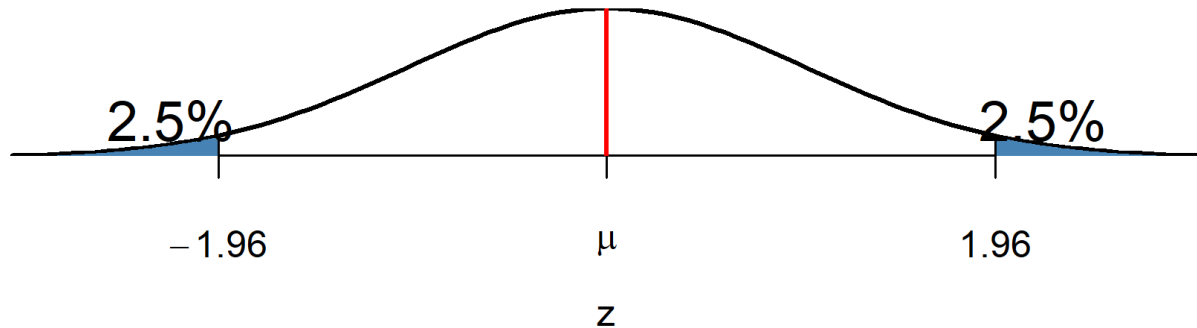
In probability tables, we can find k'_2 , such that $P(Z < k'_2) = 0.975$

Standard normal



- It's 97.5% quantile of standard normal: $k'_2 = v_{0.975}^Z = 1.96$.
 $P(Z < v_{0.975}) = P(Z < 1.96) = 0.975$
- Let's go back to $k'_2 = \frac{k_2 - \mu_X}{\frac{\sigma_X}{\sqrt{n}}}$, from which we can back-out k_2
 - $k_2 = \mu_X + k'_2 \frac{\sigma_X}{\sqrt{n}} = \mu_X + 1.96 \frac{\sigma_X}{\sqrt{n}}$
- By symmetry of normal, $k'_1 = -k'_2$, so $k_1 = \mu_X + \frac{\sigma_X}{\sqrt{n}} k'_1 = \mu_X - 1.96 \frac{\sigma_X}{\sqrt{n}}$

Standard normal



- Another way to see it:
 - Our variable is a linear transformation of a standard normal:
$$\bar{X} = \mu_X + \frac{\sigma_X}{\sqrt{n}} Z$$
 - So its percentiles are linear transformation of st.nor. quantiles:
$$v_{0.975}^{\bar{X}} = \mu_X + \frac{\sigma_X}{\sqrt{n}} v_{0.975}^Z$$

Yet another way to see it

$$\begin{aligned} 0.95 &= P(-1.96 < Z < 1.96) \\ &= P(z_{-\frac{\alpha}{2}} < Z < z_{\frac{\alpha}{2}}) \\ &= P(z_{-\frac{\alpha}{2}} < \frac{\bar{X} - \mu}{\sigma/\sqrt{n}} < z_{\frac{\alpha}{2}}) \\ &= P(z_{-\frac{\alpha}{2}}\sigma/\sqrt{n} < \bar{X} - \mu < z_{\frac{\alpha}{2}}\sigma/\sqrt{n}) \\ &= P(\mu - z_{-\frac{\alpha}{2}}\sigma/\sqrt{n} < \bar{X} < \mu + z_{\frac{\alpha}{2}}\sigma/\sqrt{n}) \end{aligned}$$

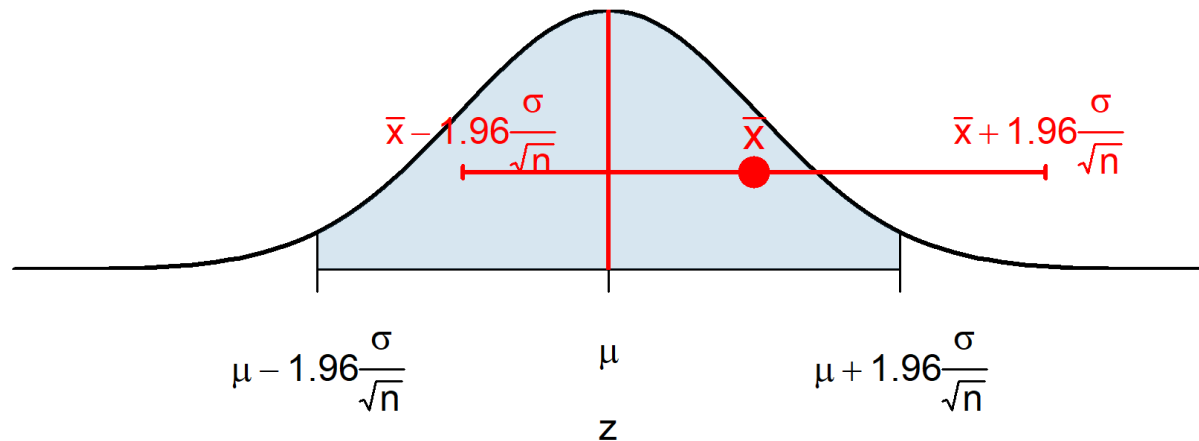
- Theoretically, CLT theorem guarantees that $\frac{\bar{X} - \mu}{\sigma/\sqrt{n}}$ is standard normal
- What happens if you do not know σ ?
- In large sample, $s \rightarrow \sigma$, so $\frac{\bar{X} - \mu}{s/\sqrt{n}} \rightarrow N(0, 1)$
- So in large samples, standardized sample mean (with estimated standard deviation) will also have normal distribution
- You may need a bit higher n to ensure $s \rightarrow \sigma$

Sampling distribution

Q: How far is the sampled mean from the true mean?

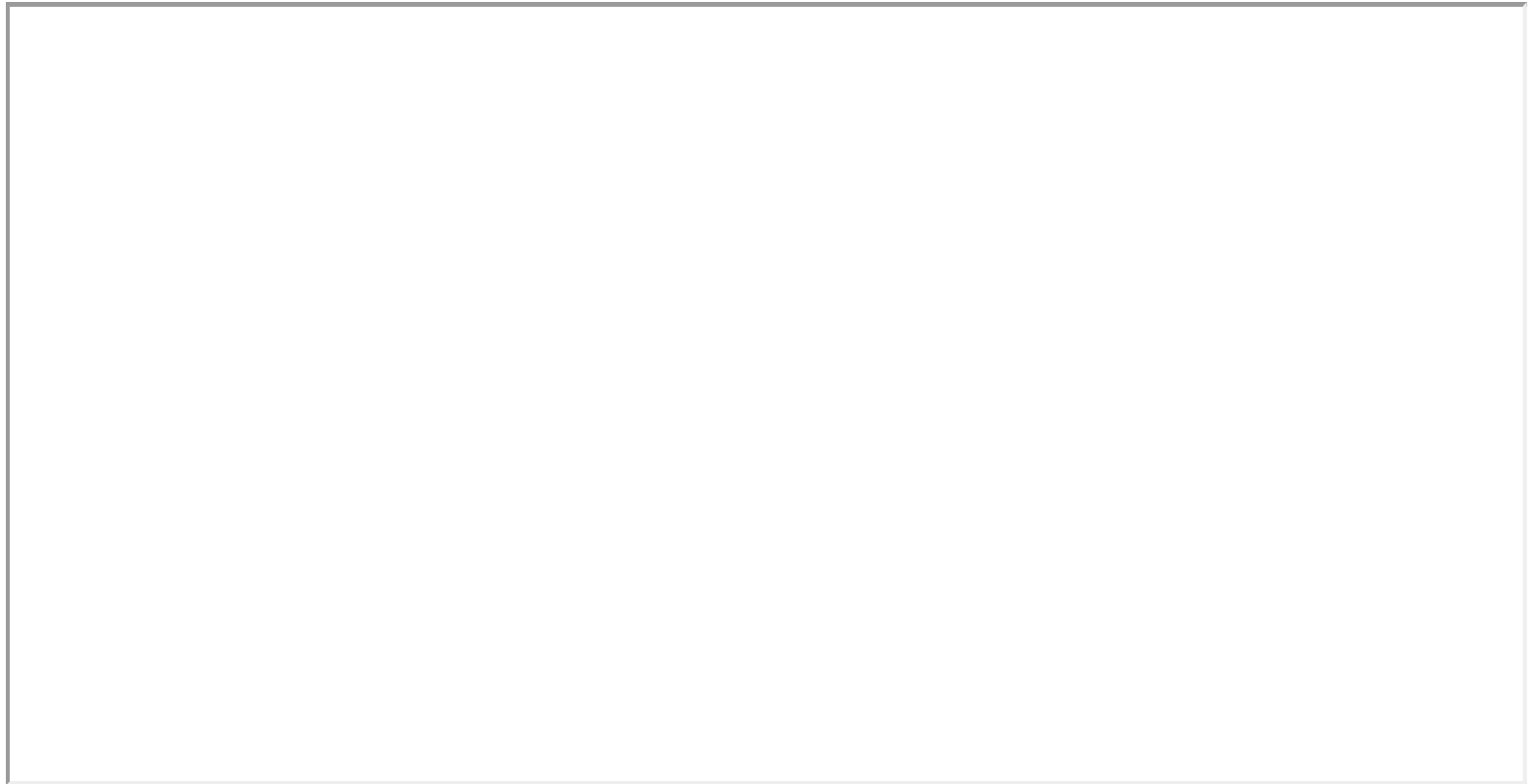
- Hence 95% of the draws of sample means will be within distance of $1.96 \frac{\sigma_X}{\sqrt{n}}$ to the true parameter
- There is only 5% chance that we have draw sample weird enough that \bar{X} is further from μ_X by more than $1.96 \frac{\sigma_X}{\sqrt{n}}$
- Confidence interval of \bar{X} will cover μ_X as long as $|\mu_X - \bar{X}| < 1.96 \frac{\sigma_X}{\sqrt{n}}$

Sampling Distribution of Sample Mean



Sampling distribution

- Suppose we draw many samples from the same distribution
- For each sample we compute the sample mean and we construct the interval
- 95% of them will cover the true population mean!



Source: [<https://seeing-theory.brown.edu/frequentist-inference/index.html#section2>]

Calculation Procedure

Use this procedure if $n > 30$

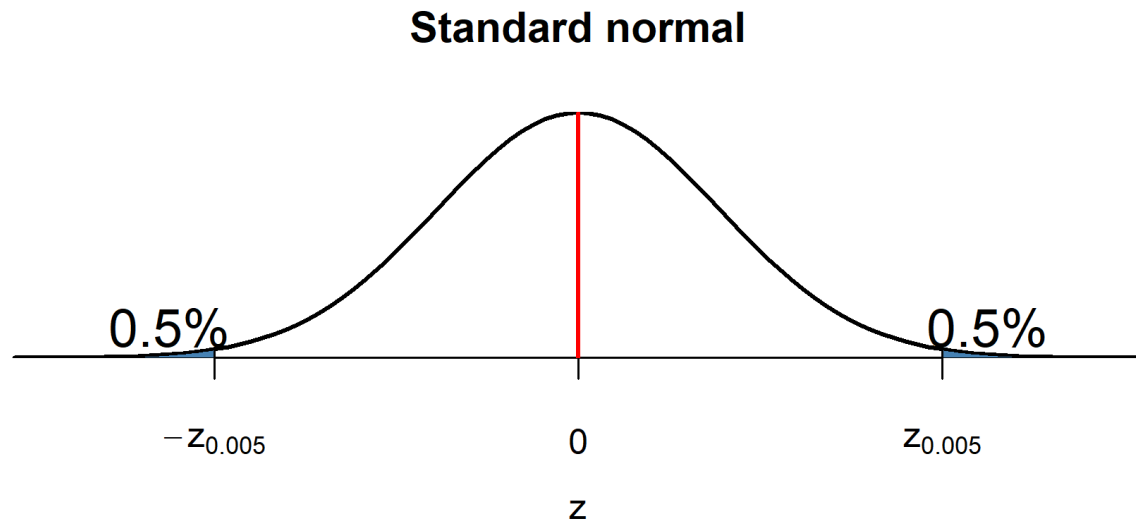
1. Take an IID sample
2. Calculate mean \bar{x} and standard deviation s in your sample
 - Standard Error is standard deviation of the estimator $SE = \frac{s}{\sqrt{n}}$
3. Pick confidence level (usually 90,95,99%)
 - We typically denote the confidence level $1 - \alpha$
 - α is probability of making a Type 1 error (more about it later)
 - **Example**: if confidence level is 95%, $\alpha = 0.05$
4. Find the corresponding critical values $z_{\frac{\alpha}{2}}$
 - Critical values are such that $P(-z_{\frac{\alpha}{2}} < Z < z_{\frac{\alpha}{2}}) = 1 - \alpha$
 - **Example**: if confidence level is 95%, $z_{\frac{\alpha}{2}} = z_{0.025} = 1.96$
5. Construct the confidence interval as:

$$\left\{ \bar{x} - \underbrace{z_{\frac{\alpha}{2}} * \frac{s}{\sqrt{n}}}_{SE}, \bar{x} + z_{\frac{\alpha}{2}} * \frac{s}{\sqrt{n}} \right\}$$

Finding Critical Values

- Suppose confidence interval is 99%.
- Then $\alpha = 0.01$
- We are looking for $z_{\frac{\alpha}{2}}$ such that:

$$P(-z_{\frac{\alpha}{2}} < Z < z_{\frac{\alpha}{2}}) = 0.99$$



$$P(Z > z_{0.005}) = 0.005 \quad \text{or} \quad P(Z < z_{0.005}) = 0.995$$

Standard Normal Distribution Table

We may use [Cookies](#)

OK

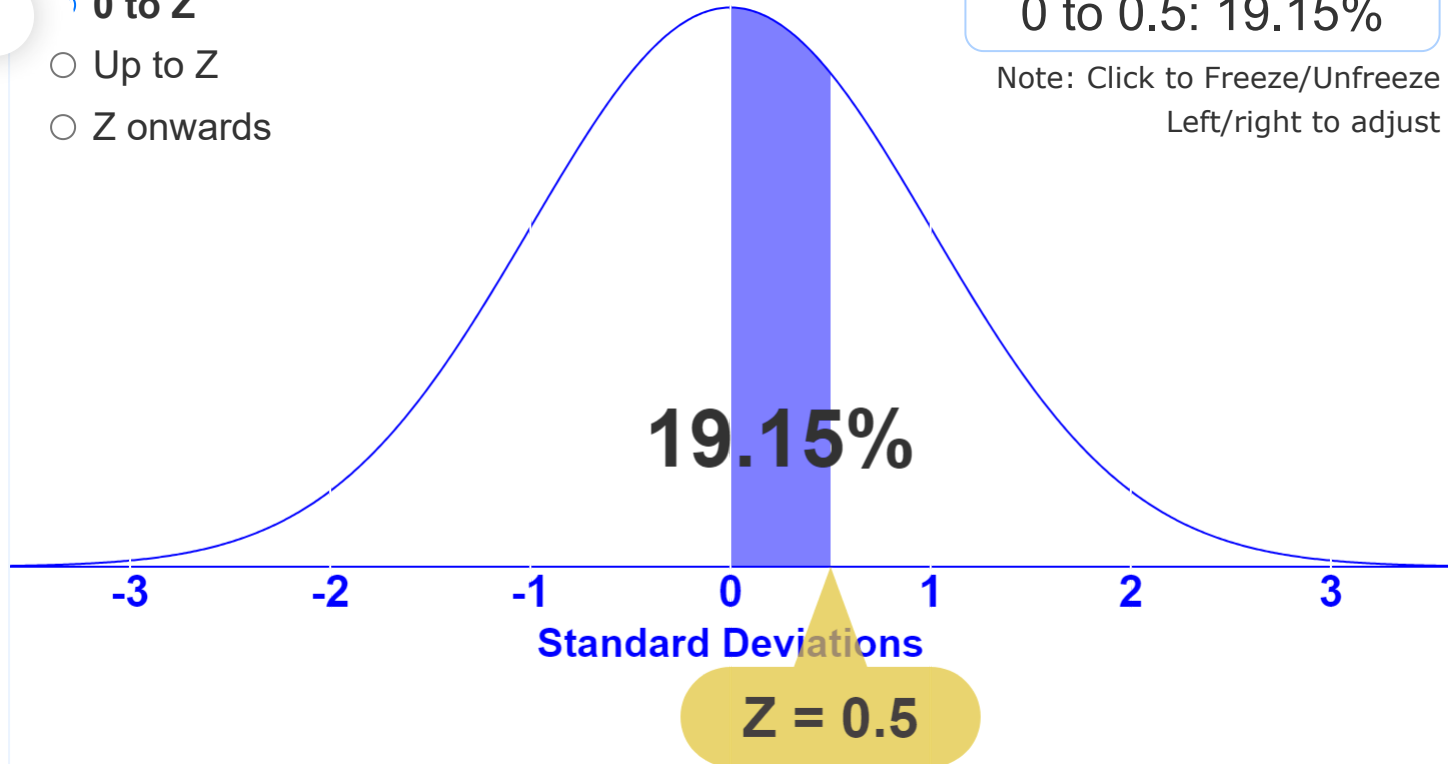
☒ 0 to Z

☐ Up to Z

☐ Z onwards

0 to 0.5: 19.15%

Note: Click to Freeze/Unfreeze
Left/right to adjust

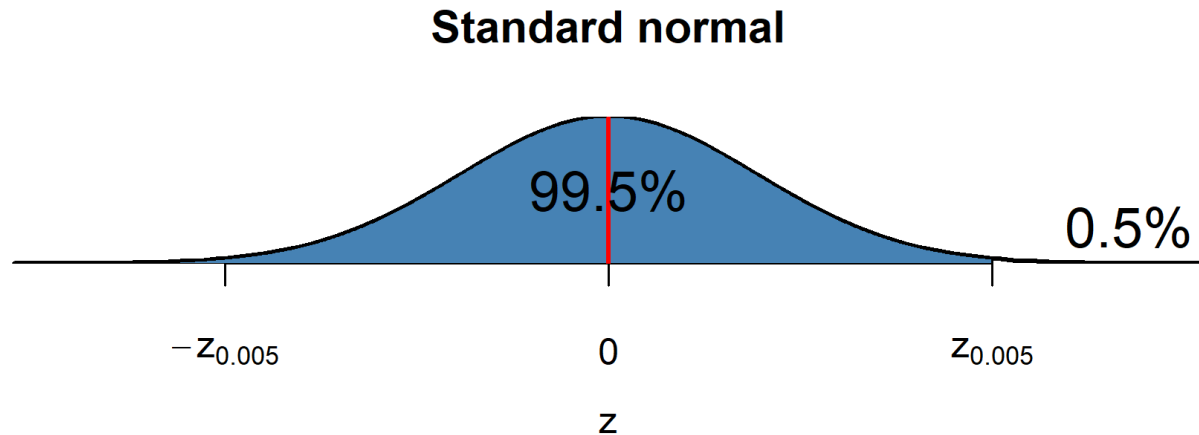


© 2021 MathsIsFun.com v0.78

Source: [<https://www.mathsisfun.com/data/standard-normal-distribution-table.html>]

Finding Critical Values

$P(Z < z_{\frac{\alpha}{2}}) = 0.995$ $z_{\frac{\alpha}{2}}$, is 99.5% quantile of standard normal $\rightarrow z_{\frac{\alpha}{2}} = 2.58$



2.1	0	0.01	0.02	0.03	0.04	0.05	0.06	0.07	0.08	0.09
2.2	0.9860965525	0.9864474189	0.9867906162	0.9871262786	0.9874545386	0.9877755273	0.9880893746	0.9883962085	0.9886961558	0.9889893417
2.3	0.98927589	0.9895559229	0.9898295613	0.9900969244	0.9903581301	0.9906132945	0.9908625325	0.9911059574	0.991343681	0.9915758136
2.4	0.9918024641	0.9920237397	0.9922397464	0.9924505886	0.992656369	0.9928571893	0.9930531492	0.9932443474	0.9934308809	0.9936128452
2.5	0.9937903347	0.9939634419	0.9941322583	0.9942968737	0.9944573766	0.994613854	0.9947663918	0.9949150743	0.9950599842	0.9952012034
2.6	0.995338812	0.9954728889	0.9956035117	0.9957307566	0.9958546986	0.9959754115	0.9960929674	0.9962074377	0.996318892	0.996427399
2.7	0.9965330262	0.9966358396	0.9967359042	0.9968332837	0.9969280408	0.9970202368	0.9971099319	0.9971971854	0.9972820551	0.9973645979
2.8	0.9974448697	0.997522925	0.9975988175	0.9976725998	0.9977443233	0.9978140385	0.997881795	0.997947641	0.9980116241	0.9980737909
2.9	0.9981341867	0.9981928562	0.9982498431	0.99830519	0.9983589388	0.9984111304	0.9984618048	0.9985110013	0.9985587581	0.9986051128
3	0.998650102	0.9986937616	0.9987361266	0.9987772313	0.9988171093	0.9988557932	0.998893315	0.9989297061	0.998964997	0.9989992175

Constructing CI: example

Let's calculate 90% CI for average price of listing with grade>4.5

1. Take an IID sample
 - $n = 100$ ✓
2. Calculate mean \bar{x} and standard deviation s
 - $\bar{x} = 1245.43$ and $s = 961.9$
3. Pick confidence level
 - We pick 90%, so $\alpha = 0.1$
4. Find the corresponding critical values $z_{\frac{\alpha}{2}}$
 - Find $z_{\frac{\alpha}{2}}$ such that $P(Z > z_{\frac{\alpha}{2}}) = 0.05$ (or $P(Z < z_{\frac{\alpha}{2}}) = 0.95$)
 - $z_{0.05} = 1.65$
5. Construct the confidence interval as:

$$\left\{ \bar{x} - z_{\frac{\alpha}{2}} * \frac{s}{\sqrt{n}}, \bar{x} + z_{\frac{\alpha}{2}} * \frac{s}{\sqrt{n}} \right\}$$

$$\left\{ 1245.43 - 1.65 \frac{961.9}{\sqrt{100}}, 1245.43 + 1.65 \frac{961.9}{\sqrt{100}} \right\}$$

Interpreting confidence intervals

$$CI_{90} = \{1086.64, 1404.22\}$$

How do we interpret a 90% confidence interval we computed?

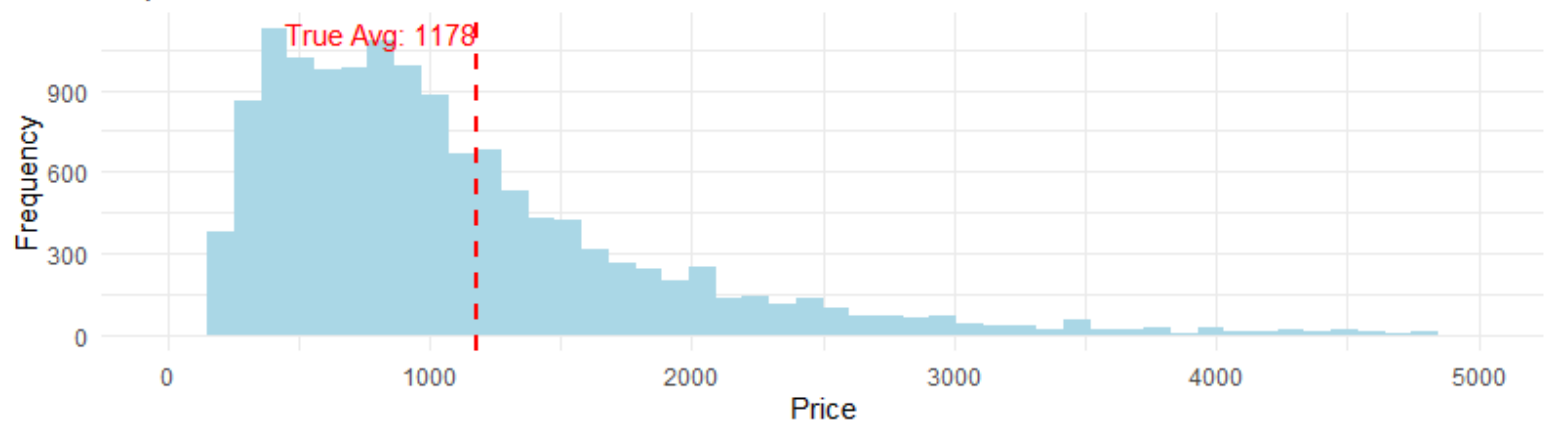
- **Correct Interpretation**

- We are 90% confident that the interval captures the true mean
- We are 90% confident that the true mean price of listings with grade>4.5 is between 1086.64 and 1404.22

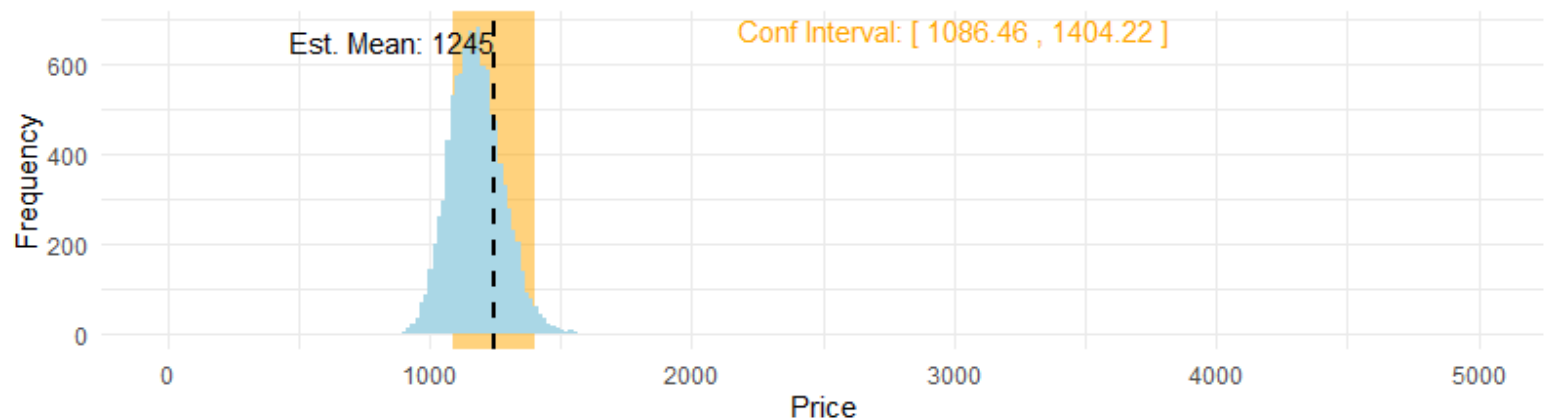
- **Incorrect**

- With 90% probability the true mean is between 1086.64 and 1404.22
 - Computed interval is not-random and true mean is not random, so can't make probabilistic statements.
 - Interval is a function of random variables only **before** we draw a sample and make any computation.
 - After we have a sample, nothing is random. The true mean is either between 1086.64 and 1404.22 or not.

Population Distribution



Sampling Distribution of the Mean and Our Estimates



Shape of confidence intervals

Confidence intervals $\{\bar{x} - z_{\frac{\alpha}{2}} * \frac{s}{\sqrt{n}}, \bar{x} + z_{\frac{\alpha}{2}} * \frac{s}{\sqrt{n}}\}$ are wider when:

- Confidence level is higher (99% is wider than 90%)
- When n is small
- When σ is large

Practice

Suppose we want to know what is the average commute time for ITAM students. We take a sample of 60 students. We calculate sample mean to be $\bar{x} = 23$ and sample standard deviation to be $\sum x^2 = 35516$. Calculate 99% confidence interval and say whether interpretation is correct:



Practice

Suppose we want to know what is the average commute time for ITAM students. We take a sample of 60 students. We calculate sample mean to be $\bar{x} = 23$ and sample standard deviation to be $\sum x^2 = 35516$. Calculate 99% confidence interval and say whether interpretation is correct:



Join by Web PollEv.com/krzysztofzaremba186



What is the standard deviation of commute time?

Nobody has responded yet.

Hang tight! Responses are coming in.

Practice

Suppose we want to know what is the average commute time for ITAM students. We take a sample of 60 students. We calculate sample mean to be $\bar{x} = 23$ and sample standard deviation to be $\sum x^2 = 35516$. Calculate 99% confidence interval and say whether interpretation is correct:



Practice

Suppose we want to know what is the average commute time for ITAM students. We take a sample of 60 students. We calculate sample mean to be $\bar{x} = 23$ and sample standard deviation to be $\sum x^2 = 35516$. Calculate 99% confidence interval and say whether interpretation is correct:

We are 99% confident that the average commute of these 60 students is between (20.3364, 25.6636)

0

We are 99% confident that the average commute of all ITAM students is between (20.3364, 25.6636) ✓

0

A 95% confidence interval would be wider

0

99% of random samples would have mean between (20.3364, 25.6636)

0

99% of random samples would capture the true mean ✓

0

With 99% probability true mean is between this and this (20.3364, 25.6636)

0

What critical values?

When should we use critical values from Normal Distribution?

1. Original distribution (of \bar{X}) is not normal:

- If $n > 30$ - use critical values from normal distribution (40 if σ unknown)
- If $n < 30$ - you are screwed

2. Original distribution (of \bar{X}) is normal:

- If you know σ , you can use critical values from normal (n doesn't matter)
 - If \bar{X} is normal, then use σ instead of s and $\frac{\bar{X}-\mu}{\frac{\sigma}{\sqrt{n}}} \sim N(0, 1)$
- If you don't know σ but $n > 40$, you can use critical values from normal
 - CLT kicks in
- If you don't know σ and $n < 40$, you use critical values from [student's t](#).
 - $\frac{\bar{X}-\mu}{\frac{s}{\sqrt{n}}}$ is not normal. s is not a good approx. of σ when n is low

What's Student's t?

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

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

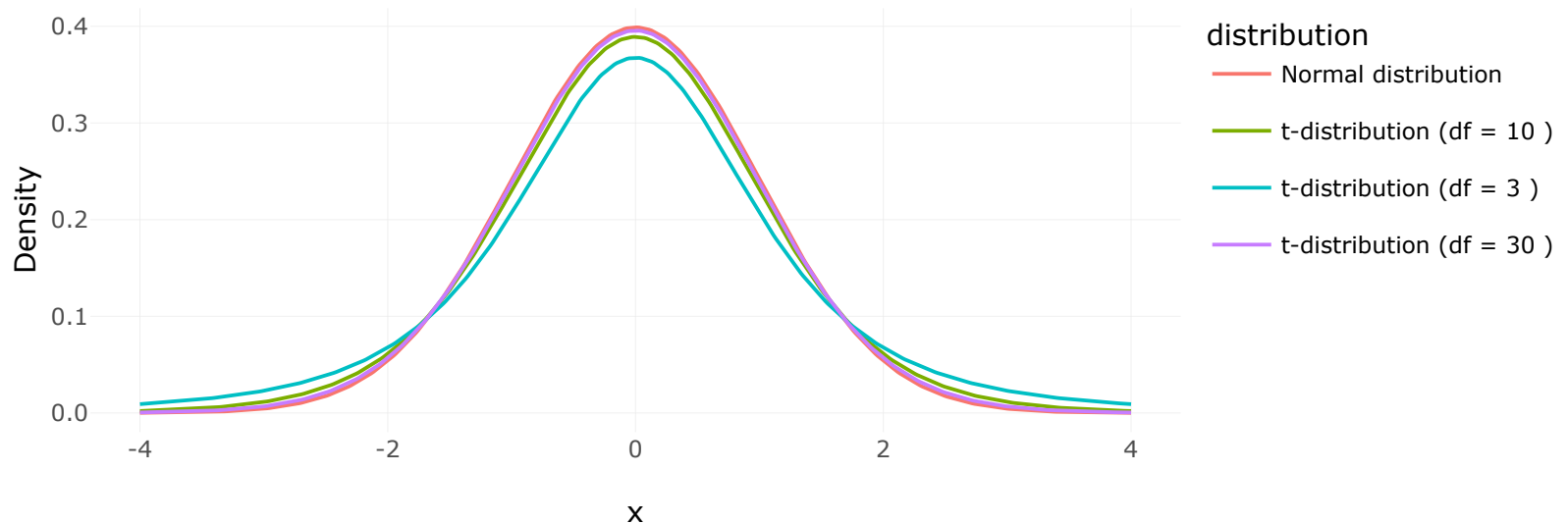
Where s is sample standard deviation.

T has a student's t distribution with $n-1$ degrees of freedom

$$T \sim t_{n-1}$$

What's Student's t?

- Bell shaped and symmetric around 0
- More spread out - heavier tails, more uncertainty (because we don't know standard deviation)
- Shape determined by the degrees of freedom.
 - As n increases (and hence degrees of freedom), it tends to standard normal (as it should by CLT!)
 - Less uncertainty because we are better at estimating standard deviation

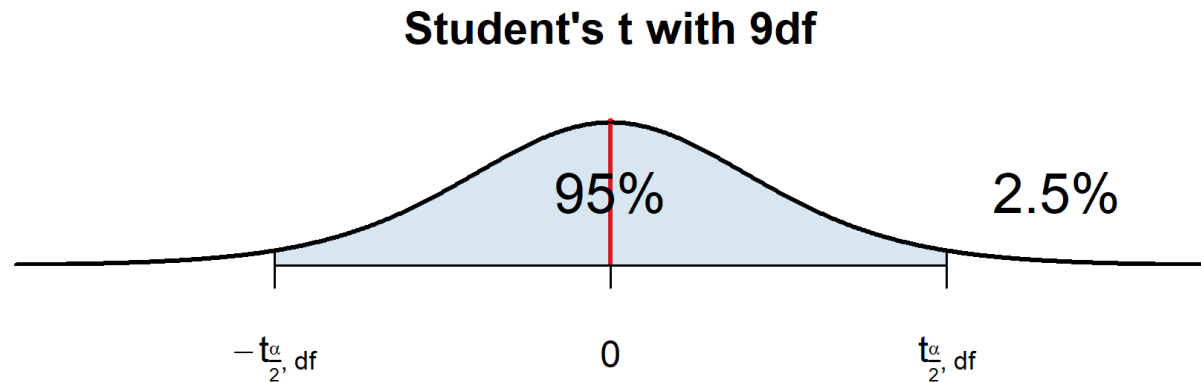


Student's t critical values

Finding critical values for student's t distribution:

1. Determine what is the right number of degrees of freedom ($n - 1$)!
2. Determine what's your confidence level and your ($1 - \alpha$)
 - From this figure out $\alpha/2$
3. Find the percentile such that

$$P(T > t_{\frac{\alpha}{2}, \underbrace{n-1}_{d.f.}}) = \frac{\alpha}{2} \quad \text{or} \quad P(T < t_{\frac{\alpha}{2}, n-1}) = 1 - \frac{\alpha}{2}$$



Example

- $n = 10 \rightarrow df = 9$
- Confidence level is 95% $\rightarrow 1 - \alpha = 0.95$ and $\frac{\alpha}{2} = 0.025$
- What's $t_{0.025,9}$ such that $P(T < t_{0.025,9}) = 0.975$

5. Distribución t de Student

$$T \sim t_\nu$$

siendo ν los grados de libertad.

$$p = P(T \leq t) = \int_{-\infty}^t \phi_T(u) du = 1 - \alpha$$

donde

$$\phi_T(u) = \frac{1}{\sqrt{\nu\pi}} \frac{\Gamma\left(\frac{\nu+1}{2}\right)}{\Gamma\left(\frac{\nu}{2}\right)} \left(1 + \frac{u^2}{\nu}\right)^{-\frac{\nu+1}{2}}$$

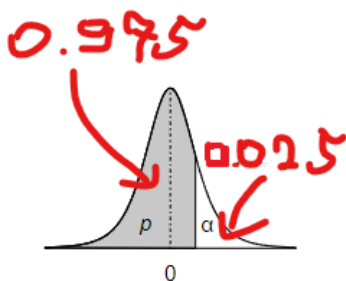
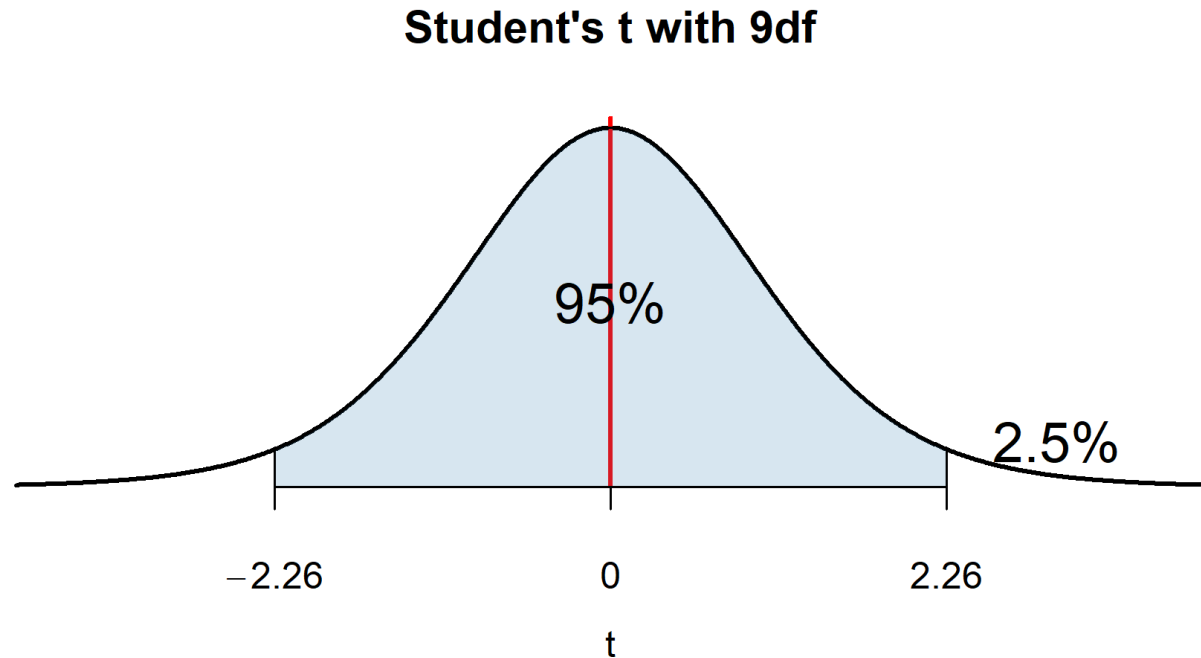


Tabla 5. Valores críticos $t_{(\alpha;\nu)}$ de la distribución t de Student.

ν	0.75	0.80	0.90	0.95	0.975	p 0.99	0.995	0.999	0.9995	0.9999
	0.25	0.20	0.10	0.05	0.025	α 0.01	0.005	0.001	0.0005	0.0001
1	1.000	1.376	3.078	6.314	12.706	31.821	63.657	318.309	636.619	3183.099
2	0.816	1.061	1.886	2.920	4.303	6.965	9.925	22.327	31.599	70.700
3	0.765	0.978	1.638	2.353	3.182	4.541	5.841	10.215	12.924	22.204
4	0.741	0.941	1.533	2.132	2.776	3.747	4.604	7.173	8.610	13.034
5	0.727	0.920	1.476	2.015	2.571	3.365	4.032	5.893	6.869	9.678
6	0.718	0.906	1.440	1.943	2.447	3.143	3.707	5.208	5.959	8.025
7	0.711	0.896	1.415	1.895	2.365	2.998	3.499	4.785	5.408	7.063
8	0.706	0.889	1.397	1.860	2.306	2.896	3.355	4.501	5.041	6.442
9	0.703	0.883	1.383	1.833	2.262	2.821	3.250	4.297	4.781	6.010
10	0.700	0.879	1.372	1.812	2.228	2.764	3.169	4.144	4.587	5.694
11	0.697	0.876	1.363	1.796	2.201	2.718	3.106	4.025	4.437	5.453
12	0.695	0.873	1.356	1.782	2.179	2.681	3.055	3.930	4.318	5.263
13	0.694	0.870	1.350	1.771	2.160	2.650	3.012	3.852	4.221	5.111
14	0.692	0.868	1.345	1.761	2.145	2.624	2.977	3.787	4.140	4.985
15	0.691	0.866	1.341	1.753	2.131	2.602	2.947	3.733	4.073	4.880
16	0.690	0.865	1.337	1.746	2.120	2.583	2.921	3.686	4.015	4.791
17	0.689	0.863	1.333	1.740	2.110	2.567	2.898	3.646	3.965	4.714
18	0.688	0.862	1.330	1.734	2.101	2.552	2.878	3.610	3.922	4.648
19	0.688	0.861	1.328	1.729	2.093	2.539	2.861	3.579	3.883	4.590
20	0.687	0.860	1.325	1.725	2.086	2.528	2.845	3.559	3.859	4.539

df

- $n = 10 \rightarrow df = 9$
- Confidence level is 95% $\rightarrow 1 - \alpha = 0.95$ and $\frac{\alpha}{2} = 0.025$
- What's $t_{0.025,9}$ such that $P(T < t_{0.025,9}) = 0.975$



- Once we have critical value, we construct the CI as before:

$$\left\{ \bar{x} - t_{\frac{\alpha}{2}, n-1} * \frac{s}{\sqrt{n}}, \bar{x} + t_{\frac{\alpha}{2}, n-1} * \frac{s}{\sqrt{n}} \right\}$$

Practice:

Your company implemented free shipping for a random group of customers. They want to know whether it increased spending. Here is your data:

\$157.80, \$192.45, \$210.20, \$175.60, \$198.30, \$180.90, \$205.75, \$185.20, \$177.40, \$195.60

a) Calculate 90% confidence interval. What assumptions you need?

- Hint 1: $\sum_i x_i = 1869.80$
- Hint 2: $\sum_i x_i^2 = 361752.55$

b) Average spending without free shipping is \$182, can say anything about whether free shipping increased spending?

Confidence Intervals for Variance

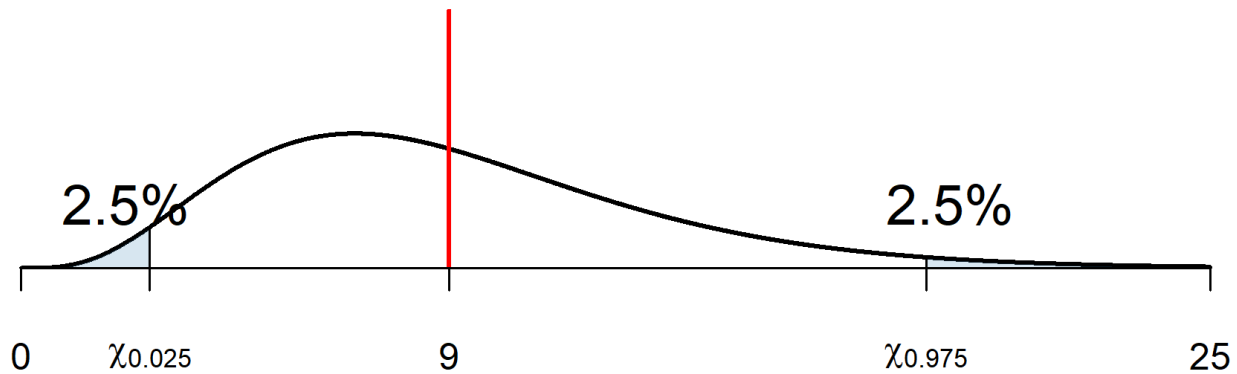
- Suppose X_1, X_2, \dots, X_n come from normal distribution
- The sampling distribution of the sample variance $S^2 = \frac{\sum_i (X_i - \bar{X})^2}{n-1}$ is:

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

- We will use the fact that:

$$P(\chi_{0.025, n-1} < \frac{(n-1)S^2}{\sigma^2} < \chi_{0.975, n-1}) = 0.95$$

Chi-Square with 9df



Confidence Intervals for Variance

How we use it to construct the confidence interval?

$$\begin{aligned} 0.95 &= P(\chi_{0.025,n-1} < \frac{(n-1)S^2}{\sigma^2} < \chi_{0.975,n-1}) \\ &= P(\frac{1}{\chi_{0.975,n-1}} < \frac{\sigma^2}{(n-1)S^2} < \frac{1}{\chi_{0.025,n-1}}) \\ &= P(\frac{(n-1)S^2}{\chi_{0.975,n-1}} < \sigma^2 < \frac{(n-1)S^2}{\chi_{0.025,n-1}}) \end{aligned}$$

So more generally, the confidence interval for the sample variance is

$$CI_{1-\alpha} = \left\{ \frac{(n-1)S^2}{\chi_{1-\frac{\alpha}{2},n-1}}, \frac{(n-1)S^2}{\chi_{\frac{\alpha}{2},n-1}} \right\}$$

- Where $\chi_{1-\frac{\alpha}{2},n-1}$ and $\chi_{\frac{\alpha}{2},n-1}$ are quantiles of χ_{n-1} distribution, such that $P(X < \chi_{1-\frac{\alpha}{2},n-1}) = 1 - \frac{\alpha}{2}$ and $P(X < \chi_{\frac{\alpha}{2},n-1}) = \frac{\alpha}{2}$
- You can read them off the tables

Practice

Suppose you produce sausages. As a quality control, you measure the level of fat in your sausages. You take a random sample of 12 sausages and you find the variance of 20 (*grams*²). Find 99% confidence interval for the variance. What assumptions you need?

Exercises:

- Review Exercises:
 - PDF 4: 1,2,3,4,5,6,7,8,9

