

Plan for the week

- M: Classical Statistics
 - Confidence Intervals
- W: Classical Statistics (cont'd)
 - Hypothesis Testing
- F: Guest Speaker: Kate Miller
 - Attendance will be taken
 - But regardless, you definitely don't want to miss this lecture!

Confidence Intervals

Point Estimate

Given a statistical model of a population, a **point estimate** is a single value used to estimate a model parameter.

Examples:

- A **sample mean** is often used to estimate the mean (i.e., the “true” mean) of a normally distributed random variable.
- Likewise, a **sample proportion** is often used to estimate the probability of success of a binomially distributed random variable.

But even the very best, data-driven point estimate is often wrong!

Interval Estimate

- Goal is to find not just a single point estimate, but an **interval estimate**, which is a plausible range of values for the parameter of interest
- As such, an interval estimate is delimited by an **upper** and **lower bound**
- Further, it quantifies the uncertainty in the estimate, via a confidence level, α
 - α is usually small
 - $\alpha = 0.05$ implies a 95% confidence interval
 - $\alpha = 0.10$ implies a 90% confidence interval
- $\Pr[z_{lo} \leq Z \leq z_{hi}] = 1 - \alpha$

Potential Pitfall

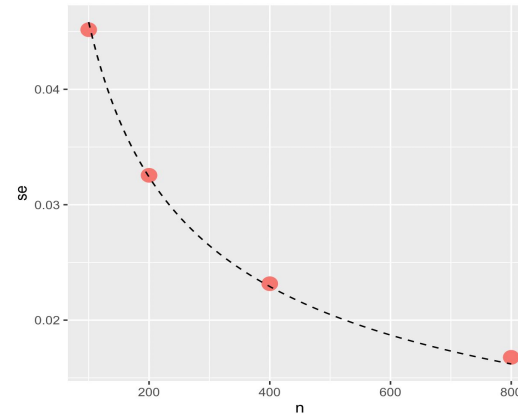
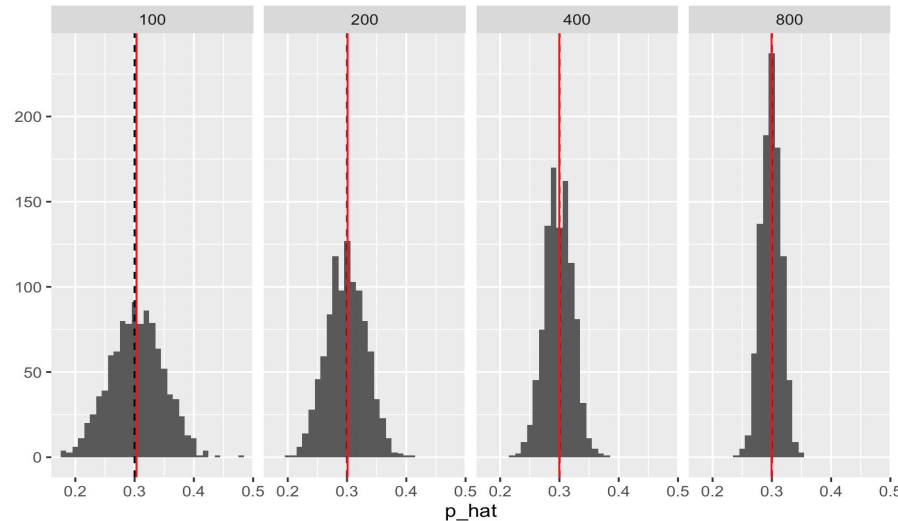
A confidence interval is a pair $\theta_{lo} \leq \theta_{hi}$ such that $\Pr[\theta_{lo} \leq \theta \leq \theta_{hi}] \geq 1 - \alpha$, where the randomness stems from the sampling process and impacts θ_{lo} and θ_{hi}

A confidence interval is a **NOT** pair $\theta_{lo} \leq \theta_{hi}$ such that $\Pr_{\theta}[\theta_{lo} \leq \theta \leq \theta_{hi}] \geq 1 - \alpha$, because θ is **NOT** a random quantity (in classical statistics, anyway!)

Aside: A **credible** interval **IS** a pair $\theta_{lo} \leq \theta_{hi}$ such that $\Pr_{\theta}[\theta_{lo} \leq \theta \leq \theta_{hi}] \geq 1 - \alpha$, because θ **IS** a random quantity in Bayesian statistics

The Mighty Central Limit Theorem

- The CLT tells us how a sample mean, say \bar{Y} , is distributed:
 - It is normally distributed with mean μ and standard error $SE = \sigma/\sqrt{n}$, where μ and σ are the population mean and standard deviation



The Mighty Central Limit Theorem

- The CLT tells us how a sample mean, say Y , is distributed:
 - It is normally distributed with mean μ and standard error $SE = \sigma/\sqrt{n}$, where μ and σ are the population mean and standard deviation
- $Z = (Y - \mu) / SE$ is thus distributed according to the standard normal:
 $\Pr[z_{lo} \leq Z \leq z_{hi}] = 1 - \alpha$

The Mighty Central Limit Theorem

- The CLT tells us how a sample mean, say \bar{Y} , is distributed:
 - It is normally distributed with mean μ and standard error $SE = \sigma/\sqrt{n}$, where μ and σ are the population mean and standard deviation

- $Z = (\bar{Y} - \mu) / SE$ is thus distributed according to the standard normal:

$$\Pr[z_{lo} \leq Z \leq z_{hi}] = 1 - \alpha$$

- $1 - \alpha =$

$$\Pr[z_{lo} \leq Z \leq z_{hi}] =$$

This is what we know: a confidence interval around

Z

$$\Pr[z_{lo} \leq \bar{Y} - \mu / SE \leq z_{hi}] =$$

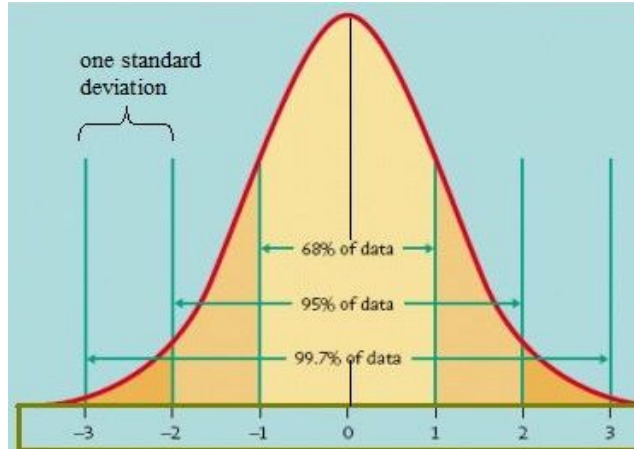
$$\Pr[\mu + z_{lo} SE \leq \bar{Y} \leq \mu + z_{hi} SE] =$$

$$\Pr[\bar{Y} + z_{lo} SE \leq \mu \leq \bar{Y} + z_{hi} SE]$$

This is what we want: a confidence interval around μ

Critical Values

- $\Pr[Y + z_{lo}SE \leq \mu \leq Y + z_{hi}SE] = 1 - \alpha$
- The choice of α dictates values for z_{lo} and z_{hi} : i.e., the **critical values**:
e.g., $\alpha = 0.5$ implies $z_{lo} = -1.96$ and $z_{hi} = 1.96$



[Image Sources](#)

Standard Normal Table

In the olden days (back when I was a student), we used the standard normal table to answer queries.

Find z_{lo} and z_{hi} s.t. $\Pr[z_{lo} \leq Z \leq z_{hi}] = 1 - \alpha$

- If $1 - \alpha = 90\%$, then $|z_{\alpha/2}| = z_{1-\alpha/2} = 1.645$
- If $1 - \alpha = 95\%$, then $|z_{\alpha/2}| = z_{1-\alpha/2} = 1.96$
- If $1 - \alpha = 98\%$, then $|z_{\alpha/2}| = z_{1-\alpha/2} = 2.33$
- If $1 - \alpha = 99\%$, then $|z_{\alpha/2}| = z_{1-\alpha/2} = 2.58$

$$z_{lo} = -z_{\alpha/2} \text{ \& } z_{hi} = z_{1-\alpha/2}$$

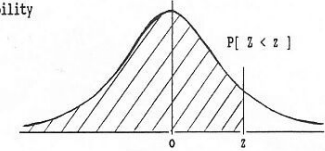
E.g., $\Pr[-1.96 \leq Z \leq 1.96] = .95$

STANDARD STATISTICAL TABLES

1. Areas under the Normal Distribution

The table gives the cumulative probability up to the standardised normal value z i.e.

$$P[Z < z] = \int_{-\infty}^z \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{1}{2}z^2\right) dz$$



z	0.00	0.01	0.02	0.03	0.04	0.05	0.06	0.07	0.08	0.09
0.0	0.5000	0.5040	0.5080	0.5120	0.5159	0.5199	0.5239	0.5279	0.5319	0.5359
0.1	0.5398	0.5438	0.5478	0.5517	0.5557	0.5596	0.5636	0.5675	0.5714	0.5753
0.2	0.5793	0.5832	0.5871	0.5910	0.5948	0.5987	0.6026	0.6064	0.6103	0.6141
0.3	0.6179	0.6217	0.6255	0.6293	0.6331	0.6368	0.6406	0.6443	0.6480	0.6517
0.4	0.6554	0.6591	0.6628	0.6664	0.6700	0.6736	0.6772	0.6808	0.6844	0.6879
0.5	0.6915	0.6950	0.6985	0.7019	0.7054	0.7088	0.7123	0.7157	0.7190	0.7224
0.6	0.7257	0.7291	0.7324	0.7357	0.7389	0.7422	0.7454	0.7486	0.7517	0.7549
0.7	0.7580	0.7611	0.7642	0.7673	0.7704	0.7734	0.7764	0.7794	0.7823	0.7854
0.8	0.7881	0.7910	0.7939	0.7967	0.7995	0.8023	0.8051	0.8078	0.8106	0.8133
0.9	0.8159	0.8186	0.8212	0.8238	0.8264	0.8289	0.8315	0.8340	0.8365	0.8389
1.0	0.8413	0.8438	0.8461	0.8485	0.8508	0.8531	0.8554	0.8577	0.8599	0.8621
1.1	0.8643	0.8665	0.8686	0.8708	0.8729	0.8749	0.8770	0.8790	0.8804	0.8830
1.2	0.8849	0.8869	0.8888	0.8907	0.8925	0.8944	0.8962	0.8980	0.8997	0.9015
1.3	0.9032	0.9049	0.9066	0.9082	0.9099	0.9115	0.9131	0.9147	0.9162	0.9177
1.4	0.9192	0.9207	0.9222	0.9236	0.9251	0.9265	0.9279	0.9292	0.9306	0.9319
1.5	0.9332	0.9345	0.9357	0.9370	0.9382	0.9394	0.9406	0.9418	0.9429	0.9441
1.6	0.9452	0.9463	0.9474	0.9484	0.9495	0.9505	0.9515	0.9525	0.9535	0.9545
1.7	0.9554	0.9564	0.9573	0.9582	0.9591	0.9599	0.9608	0.9616	0.9625	0.9633
1.8	0.9641	0.9649	0.9656	0.9664	0.9671	0.9678	0.9686	0.9693	0.9699	0.9706
1.9	0.9713	0.9719	0.9726	0.9732	0.9738	0.9744	0.9750	0.9756	0.9761	0.9767
2.0	0.9773	0.9778	0.9783	0.9788	0.9793	0.9798	0.9803	0.9808	0.9812	0.9817
2.1	0.9821	0.9826	0.9830	0.9834	0.9838	0.9842	0.9846	0.9850	0.9854	0.9857
2.2	0.9861	0.9865	0.9868	0.9871	0.9874	0.9878	0.9881	0.9884	0.9887	0.9890
2.3	0.9893	0.9896	0.9898	0.9900	0.9904	0.9906	0.9909	0.9911	0.9913	0.9916
2.4	0.9918	0.9920	0.9922	0.9924	0.9927	0.9929	0.9931	0.9932	0.9934	0.9936
2.5	0.9938	0.9940	0.9941	0.9943	0.9945	0.9946	0.9948	0.9949	0.9951	0.9952
2.6	0.9953	0.9955	0.9956	0.9957	0.9959	0.9960	0.9961	0.9962	0.9963	0.9964
2.7	0.9965	0.9966	0.9967	0.9968	0.9969	0.9970	0.9971	0.9972	0.9973	0.9974
2.8	0.9974	0.9975	0.9976	0.9977	0.9977	0.9978	0.9979	0.9980	0.9980	0.9981
2.9	0.9981	0.9982	0.9982	0.9983	0.9984	0.9984	0.9985	0.9985	0.9986	0.9986
z	3.00	3.10	3.20	3.30	3.40	3.50	3.60	3.70	3.80	3.90
P	0.9986	0.9990	0.9993	0.9995	0.9997	0.9998	0.9998	0.9999	0.9999	1.0000

Standard Normal Table

These days, we use R:

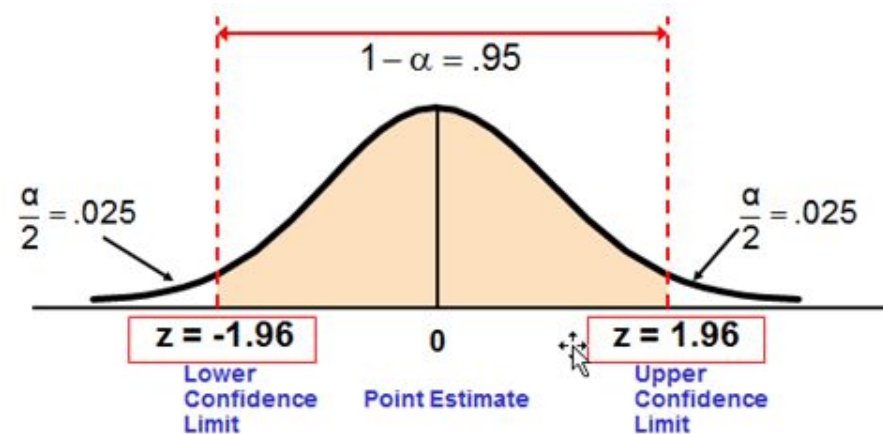
- `qnorm(0.975) = 1.959964`

Find z_{lo} and z_{hi} s.t. $\Pr[z_{lo} \leq Z \leq z_{hi}] = 1 - \alpha$

- If $1 - \alpha = 90\%$, then $|z_{\alpha/2}| = z_{1-\alpha/2} = 1.645$
- If $1 - \alpha = 95\%$, then $|z_{\alpha/2}| = z_{1-\alpha/2} = 1.96$
- If $1 - \alpha = 98\%$, then $|z_{\alpha/2}| = z_{1-\alpha/2} = 2.33$
- If $1 - \alpha = 99\%$, then $|z_{\alpha/2}| = z_{1-\alpha/2} = 2.58$

$$z_{lo} = -z_{\alpha/2} \text{ \& } z_{hi} = z_{1-\alpha/2}$$

E.g., $\Pr[-1.96 \leq Z \leq 1.96] = .95$



[Image Source](#)

Let's say you are running for mayor!*

- You hire a polling agency to determine if you are likely to win or not.
- After sampling 100 likely voters, the agency reports that 55 of those 100 support you. Woo hoo! **You are predicted to be the winner!**
- **But wait! There are more than 100 likely voters!**
- **How does 0.55, the sample mean, relate to the population mean?**
 - The polling agency also reports a **95% confidence interval** around the number 0.55: e.g., (0.45, 0.65). If this interval includes values below 0.5, your win is less certain.
 - Polling agencies tend to report margins of error (MoE = 0.1).
 - MoE and confidence intervals are related.
 - But wait...how did the polling agency come up with this confidence interval/MoE?

*Sumbul Siddiqui '10 is now in her third term as the mayor of the City of Cambridge, Massachusetts. Additionally, two Brown alumni were recently elected mayor in Greece, in Athens and in Thessaloniki.

Building a Confidence Interval

- Let X be the number of people who support you in a poll, and let $p_{\text{hat}} = X/n$.
- By the central limit theorem, for large enough n , p_{hat} is approximately normal, with mean μ and standard deviation SE.
- It follows that $\Pr[p_{\text{hat}} + z_{\text{lo}} \text{SE} \leq \mu \leq p_{\text{hat}} + z_{\text{hi}} \text{SE}] = 1 - \alpha$
- Choose $\alpha = 0.05$: $\Pr[p_{\text{hat}} - 1.96\text{SE} \leq \mu \leq p_{\text{hat}} + 1.96\text{SE}] = 0.95$
- Sounds good, but what is the standard error of p_{hat} ?

The mean and variance of X

- As $SE = \sigma/\sqrt{n}$, this begs the question: what is the standard deviation of X ?
- Well, how is it distributed? Easy: X is binomially distributed.
 - Mean = np
 - Variance = $np(1-p)$
- The mean of p_{hat} is p .
 - $\mu = E[p_{\text{hat}}] = E[X/n] = np/n = p$
- The variance of p_{hat} is $p(1 - p)/n$:
 - $\text{Var}[p_{\text{hat}}] = \text{Var}[X/n] = 1/n^2 \text{Var}[X] = 1/n^2 (np)(1 - p) = p(1 - p)/n$
 - $SE = \sqrt{p(1 - p)/n}$

Building a Confidence Interval (cont'd)

- Let X be the number of people who support you in a poll, and let $p_{\text{hat}} = X/n$.
- By the central limit theorem, for large enough n , p_{hat} is approximately normal, with mean μ and standard deviation SE.
- It follows that $\Pr[p_{\text{hat}} + z_{\text{lo}} \text{SE} \leq \mu \leq p_{\text{hat}} + z_{\text{hi}} \text{SE}] = 1 - \alpha$
- Choose $\alpha = 0.05$: $\Pr[p_{\text{hat}} - 1.96\text{SE} \leq \mu \leq p_{\text{hat}} + 1.96\text{SE}] = 0.95$
- Sounds good, but what is the standard error of p_{hat} ?
- Now we know σ : $\text{SE} = \sqrt{p(1 - p)/n}$
- Or do we?

Building a Confidence Interval (cont'd)

New problem: We don't know p !

- We estimate p by p_{hat} .
 - N.B. This is cheating, but only negligibly so.

Hence, we build our 95% confidence interval as follows:

- Let $\sigma_{\text{hat}} = \sqrt{p_{\text{hat}}(1 - p_{\text{hat}})/n} = \sqrt{(.55)(.45)/100} = 0.05$
- $\Pr[p_{\text{hat}} + z_{\text{lo}} \sigma_{\text{hat}} \leq \mu \leq p_{\text{hat}} + z_{\text{hi}} \sigma_{\text{hat}}] = .95$
 - Lower Bound: $0.55 + (-1.96)(0.05) = .45$
 - Upper Bound: $0.55 + (1.96)(0.05) = .65$
- The value $(1.96)(0.05) \approx 0.1$ is called the **margin of error**.

John Snow's Grand Experiment, Revisited

Data collected by John Snow

Supply Area	# of Houses	Cholera Deaths	Deaths/10,000 Houses
S&V	40,046	1,263	315
Lambeth	26,107	98	37
Rest of London	256,423	1,422	59

Let's compute confidence intervals

- Choose $\alpha = 0.05$: $\Pr[p_{\text{hat}} - 1.96\sigma_{\text{hat}} \leq \mu \leq p_{\text{hat}} + 1.96\sigma_{\text{hat}}] = 0.95$
- We need to calculate p_{hat} and σ_{hat} to find a confidence interval
- For S&V:
 - $P_{\text{hat}} = 1263/40046 \approx 0.0315$
 - $\text{Var}[p_{\text{hat}}] = (0.0315)(1 - 0.0315) / 40046 = 7.62 \times 10^{-7}$
 - The standard error is the square root of the variance: $\text{SE}[p_{\text{hat}}] = \sqrt{7.62 \times 10^{-7}} = 0.00087$
- So the confidence interval at the 95% level is:
 - $[0.0315 - (1.96)(0.00087), 0.0315 + (1.96)(0.00087)] = [0.03, 0.033]$

Let's compute confidence intervals

- Choose $\alpha = 0.05$: $\Pr[p_{\text{hat}} - 1.96\sigma_{\text{hat}} \leq \mu \leq p_{\text{hat}} + 1.96\sigma_{\text{hat}}] = 0.95$
- We need to calculate p_{hat} and σ_{hat} to find a confidence interval
- For Lambeth:
 - $P_{\text{hat}} = 98/26107 \approx 0.00375$
 - $\text{Var}[p_{\text{hat}}] = (0.00375)(1 - 0.00375) / 26107 = 1.43 \times 10^{-7}$
 - The standard error is the square root of the variance: $\text{SE}[p_{\text{hat}}] = \sqrt{1.43 \times 10^{-7}} = 0.000378$
- So the confidence interval at the 95% level is:
 - $[0.00375 - (1.96)(0.000378), 0.00375 + (1.96)(0.000378)] = [0.003, 0.0044]$

Let's compute confidence intervals

- Choose $\alpha = 0.05$: $\Pr[p_{\text{hat}} - 1.96\sigma_{\text{hat}} \leq \mu \leq p_{\text{hat}} + 1.96\sigma_{\text{hat}}] = 0.95$
- We need to calculate p_{hat} and σ_{hat} to find a confidence interval
- For Lambeth:
 - $P_{\text{hat}} = 1422/256423 \approx 0.0055$
 - $\text{Var}[p_{\text{hat}}] = (0.0055)(1 - 0.0055) / 256423 = 2.133 \times 10^{-8}$
 - The standard error is the square root of the variance: $\text{SE}[p_{\text{hat}}] = \sqrt{2.133 \times 10^{-8}} = 0.00014$
- So the confidence interval at the 95% level is:
 - $[0.0055 - (1.96)(0.00014), 0.0055 + (1.96)(0.00014)] = [0.005, 0.006]$

Let's code this up!

```
# data
area <- c("S&V", "Lambeth", "London")
houses <- c(40046, 26107, 256423)
deaths <- c(1263, 98, 1422)
cholera <- data.frame(area, houses, deaths)
```

```
> cholera
```

	area	houses	deaths
1	S&V	40046	1263
2	Lambeth	26107	98
3	London	256423	1422

Let's code this up!

```
# statistics
cholera$phat <- deaths / houses
cholera$variance <- (cholera$phat) * (1 - cholera$phat) / cholera$houses
cholera$se <- sqrt(cholera$variance)
```

```
> cholera
```

	area	houses	deaths	phat	variance	se
1	S&V	40046	1263	0.031538730	7.627238e-07	0.0008733406
2	Lambeth	26107	98	0.003753783	1.432448e-07	0.0003784769
3	London	256423	1422	0.005545524	2.150654e-08	0.0001466511

qnorm

```
> qnorm(0.01)
```

```
[1] -2.326348
```

```
> qnorm(0.025)
```

```
[1] -1.959964
```

```
> qnorm(0.05)
```

```
[1] -1.644854
```

```
> qnorm(0.95)
```

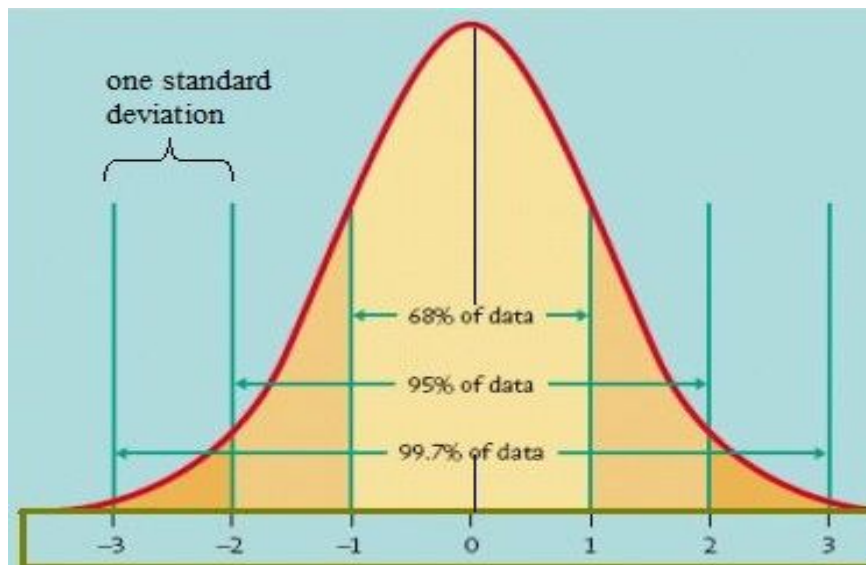
```
[1] 1.644854
```

```
> qnorm(0.975)
```

```
[1] 1.959964
```

```
> qnorm(0.99)
```

```
[1] 2.326348
```



[Image Sources](#)

Let's code this up!

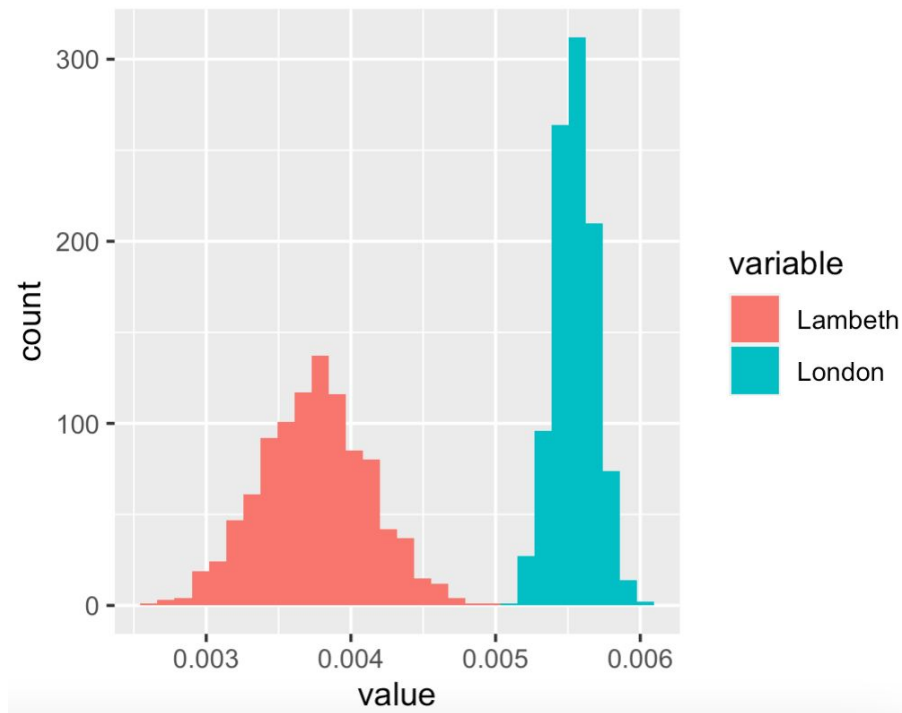
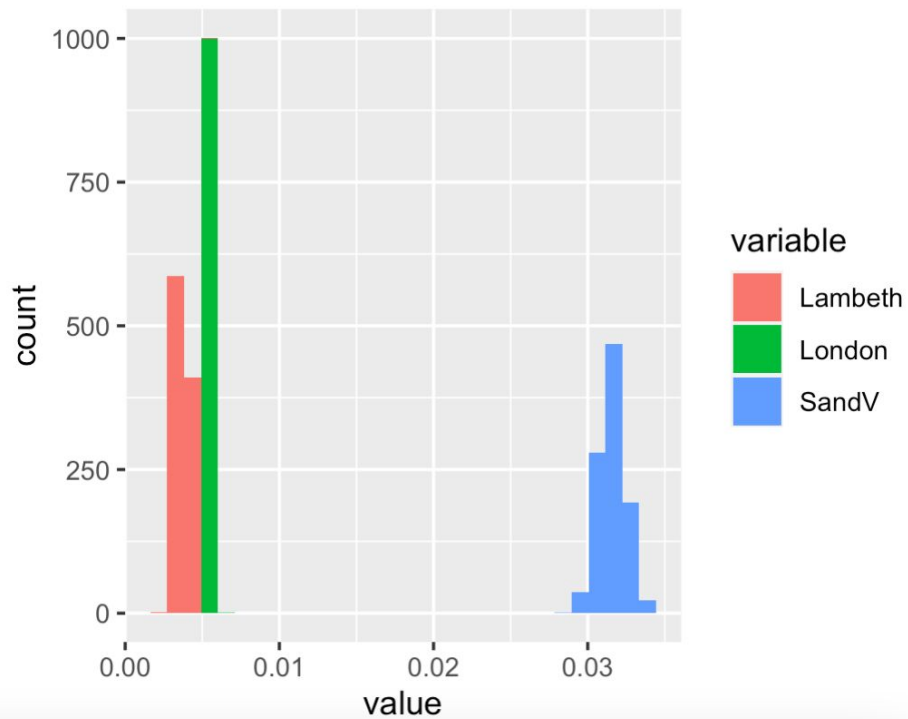
```
# confidence intervals
alpha = 0.05
quantile <- -qnorm(alpha / 2)
cholera$CI_lower <- cholera$phat - quantile * cholera$se
cholera$CI_upper <- cholera$phat + quantile * cholera$se
```

```
> quantile
[1] 1.959964
```

```
> cholera %>% select(area, phat, CI_lower:CI_upper)
```

	area	phat	CI_lower	CI_upper
1	S&V	0.031538730	0.029827014	0.033250447
2	Lambeth	0.003753783	0.003011981	0.004495584
3	London	0.005545524	0.005258094	0.005832955

Some Histograms



Deriving Standard Error

The Standard Error of the Sample Mean

- Let \bar{X}_n represent the sample mean:
 - $\bar{X}_n = (X_1 + \dots + X_n)/n$
- First, we need the variance of the sample mean:
 - X_i are normally distributed with variance σ^2
 - $\text{Var}[\bar{X}_n] = \text{Var}[(X_1 + \dots + X_n)/n] = (1/n^2) \text{Var}[X_1 + \dots + X_n] = (1/n^2) (\text{Var}[X_1] + \text{Var}[X_2] + \dots + \text{Var}[X_n]) = (1/n^2) (n) \text{Var}[X_1] = \sigma^2/n$
- We now have a formula for the standard error of the sample mean:
 - $\text{SE}[\bar{X}_n] = \sigma/\sqrt{n}$

The Standard Error of the Sample Proportion

- Let P_{hat} represent the sample proportion:
 - $P_{\text{hat}} = X/n$, where X is a binomial random variable distributed according to (n, p)
- First, we need the variance of the sample proportion:
 - $X \sim B(n, p)$
 - $\text{Var}[P_{\text{hat}}] = \text{Var}[X/n] = (1/n^2)(np)(1 - p) = p(1 - p)/n$
- We now have a formula for the standard error of the sample proportion:
 - $\text{SE}[P_{\text{hat}}] = \sqrt{p(1 - p)/n}$

Difference of Two Sample Means

- Let $X_N - Y_M$ represent the difference between two sample means.
- $X_N - Y_M = (X_1 + \dots + X_n)/n - (Y_1 + \dots + Y_m)/m$
 - X_i and Y_i are normally distributed with variance σ_X and σ_Y
 - $\text{Var}[X_N - Y_M] = \text{Var}[X_N] + \text{Var}[Y_M] = \sigma_X^2/n + \sigma_Y^2/m$
- We now have a formula for the standard error of the difference between two sample means:
 - $\text{SE}[X_N - Y_M] = \sqrt{(\sigma_X^2/n + \sigma_Y^2/m)}$

Difference of Two Sample Proportions

- Let P_1 represent one proportion, and P_2 a second proportion.
- $P_1 = X/n$ and $P_2 = Y/m$
 - $X \sim B(n, p_1)$ and $Y \sim B(m, p_2)$
 - $\text{Var}[P_1 - P_2] = \text{Var}[X/n] + \text{Var}[Y/m] = p_1(1 - p_1)/n + p_2(1 - p_2)/m$
- We now have a formula for the standard error of the difference between two sample proportions:
 - $\text{SE}[P_1 - P_2] = \sqrt{p_1(1 - p_1)/n + p_2(1 - p_2)/m}$

Student t -Distribution

Building Confidence Intervals

- The only difference between building confidence intervals using the t -distribution and building them using the normal is: the critical values for the normal distribution are familiar numbers to most statisticians. Recall:
 - Find z_{lo} and z_{hi} s.t. $P[z_{lo} \leq Z \leq z_{hi}] = 1 - \alpha\%$
 - If $1 - \alpha = 90\%$, then $|z_{\alpha/2}| = z_{1-\alpha/2} = 1.645$
 - If $1 - \alpha = 95\%$, then $|z_{\alpha/2}| = z_{1-\alpha/2} = 1.96$
 - If $1 - \alpha = 98\%$, then $|z_{\alpha/2}| = z_{1-\alpha/2} = 2.33$
 - If $1 - \alpha = 99\%$, then $|z_{\alpha/2}| = z_{1-\alpha/2} = 2.58$
 - $z_{lo} = z_{\alpha/2}$ & $z_{hi} = z_{1-\alpha/2}$ are called **critical values**.
- For the t -distribution, we have to look up the critical values.
 - There is more uncertainty because there are fewer samples.
 - So fixing α , the magnitude of these values is higher.

An Example

- Let's say we polled 30 people to see who they thought would win the 2016 presidency, and 60% of them said Clinton.
- We can build a 95% confidence interval as follows:
 - The estimate is 0.6.
 - The SE = $\sqrt{(.6)(.4)/30} = 0.09$.
 - We look up the critical t -values:
 - $qt(0.025, 29) = -2.05$ and $qt(0.975, 29) = 2.05$
 - The number 29 is the degrees of freedom, which is the number of values that are free to vary
- The 95% CI is $(0.6 - (2.05)(0.09), 0.6 + (2.05)(0.09)) = (0.42, 0.78)$

An Example

- Let's say we polled 20 people to see who they thought would win the 2016 presidency, and 60% of them said Clinton.
- We can build a 95% confidence interval as follows:
 - The estimate is 0.6.
 - The SE = $\sqrt{(.6)(.4)/20} = 0.11$.
 - We look up the critical t -values:
 - $qt(0.025, 19) = -2.09$ and $qt(0.975, 19) = 2.09$
 - The number 19 is the degrees of freedom, which is the number of values that are free to vary
- The 95% CI is $(0.6 - (2.09)(0.11), 0.6 + (2.09)(0.11)) = (0.37, 0.83)$

An Example

- If we poll 30 people, the 95% CI is $(0.6 - (2.05)(0.09), 0.6 + (2.05)(0.09)) = (0.42, 0.78)$
- If we poll 20 people, the 95% CI is $(0.6 - (2.09)(0.11), 0.6 + (2.09)(0.11)) = (0.37, 0.83)$
- If we poll 10 people, the 95% CI is $(0.6 - (2.26)(0.155), 0.6 + (2.26)(0.155)) = (0.25, 0.95)$
- Notice how the width of the confidence interval increases as the sample size decreases, holding the significance level α constant.

Back to your run for mayor, using Student t

- The polling agency polled 100 people to see if they support you for mayor, and 55 people said they did.
- We can build a 95% confidence interval as follows:
 - The estimate is 0.55.
 - The SE = $\sqrt{(.55)(.45)/100} = 0.05$.
 - We look up the critical t -values:
 - `qt(0.025, 99) = -1.98` and `qt(0.975, 99) = 1.98`
 - The number 99 is the degrees of freedom, which is the number of values that are free to vary
- The 95% CI is $(0.55 - (1.98)(0.05), 0.55 + (1.98)(0.05)) = (0.451, 0.649)$
- The 95% CI **was** $(0.55 - (1.96)(0.05), 0.55 + (1.96)(0.05)) = (0.45, 0.65)$

Extras

Some words of warning!

- Our methods make use of the central limit theorem.
 - In the examples, we do not assume anything about how preferences are distributed.
 - However, by the CLT, we know the sampling distribution is approximately normal.
 - The surveys were large enough for the CLT to apply. However, the CLT may not have applied if the sample size were smaller. (Rule of thumb: CLT applies whenever $n \geq 30$.)
- **Very common mistake:** a 95% confidence interval between 1 and 2 is often interpreted as a 95% chance the parameter lies between 1 and 2.
 - This is a misconception, because the true parameter isn't random! It is a fixed value.
 - A 95% interval means that if we repeated the experiment 100 times, 95 of the resulting 100 intervals would contain the true parameter.
 - These two interpretations are not the same. Be careful to avoid this potential pitfall!