

Reference: Master statistics & machine learning: intuition, math, code

Table of Content

Math prerequisite
What are (is_) data_
Visualising Data
Descriptive Statistics
Data Normalisation and Outliers
Probability Theory
Hypothesis Testing
The t-test Family
Confidence Intervals on Parameters
Correlation
ANOVA
Regression
Statistical Power and Sample Sizes

1. Math prerequisite

Table of Content

Scientific Notation
What is Logistic Function
Rank and tied-Rank
How to deal with ties in rank transformation ?

Scientific Notation

Example

Scientific notation

$$200 = 2 \times 10^2 = 2 \times 10 \times 10$$

$$240 = 2.4 \times 10^2$$

$$240 = 2.4e2$$

Master stats and ML — MX Cohen — sincxpress.com



Few other examples:

$$200 = 2 \times 10^2 = 2e2$$

$$0.02 = 2 \times 10^{-2} = 2e-2$$

What is Logistic Function

Logistic function $\rightarrow \log(p / 1 - p) = \beta$

And when you solve for p , you will get the solution like the sigmoid formula.

Why the logistic function in statistics?

$$\ln \frac{p}{1-p} = \beta \quad p = e^\beta - e^\beta p$$

$$\exp \left(\ln \frac{p}{1-p} \right) = e^\beta \quad e^\beta = p + e^\beta p$$

$$\frac{p}{1-p} = e^\beta \quad p = \frac{e^\beta}{1+e^\beta} \frac{e^{-\beta}}{e^{-\beta}}$$

$$p = e^\beta(1-p) \quad p = \frac{1}{1+e^{-\beta}}$$

Master stats and ML — MX Cohen — sincxpress.com



Rank and tied-Rank

- Say, $a = \{4.001, 1, 4, -10, 987\}$
- Rank is basically sort these numbers and their index will be rank.
- Therefore rank of a is $ar = \{4, 2, 3, 1, 5\}$
- This is called rank-transformation.
- This is non-linear transformation (directly map the index according to their position on number line)
- But the disadvantage is, consider 4 and 4.001, they are only one index apart, even though their distance is small. Also, consider 4.001 and 987, even though their distance is large enough but still one rank apart.
- This is lossy transformation and also it is non-invertible transformation. You cannot convert this rank transformation into the original set.

How to deal with ties in rank transformation ?

- say $a = \{10, -10, 1, 10\}$
- Now what should be rank transformation ?
- We know it like this = $ar = \{?, 1, 2, ?\}$
- Now there is a tie for which 10 should take the rank 3 and which should take 4.
- The solution is, take the average of their possible ranks (i.e 3 and 4) and assign it to both.
- Therefore, $ar = \{3.5, 1, 2, 3.5\}$
- This is called tied-rank.

Some Python Basics Code: [Github](#)

2. What are (is_) data_

Table of Content

- | |
|-------------------------------------------------|
| Is the word Data singular or plural ? |
| Type (most common) of data: |
| Why do these types matter ? |
| Sample vs Population: |
| Why should we care about sample vs population ? |
| When should you create fake data ? |

Is the word *Data* singular or plural ?

- Singular: Datun OR data point
- Plural: Data
- Therefore the sentence like: "Data is", is wrong because this means the word Data is singular, instead we should say, "these data are...".

Type (most common) of data:

- Numeric → Interval → Numeric Scale with meaningful intervals e.g Temp in Celcius. Also, zero degree celsius does not mean absence of temperature.
- Numeric → Ratio → Just like interval but with a meaningful value for zero. e.g height, money. Height of zero means absence of height.
- Numeric → Discrete → No arbitrary precision (integers) , e.g population of country.
- Categorical → Ordinal → We can sort the data e.g education degree (high school, secondary school, bachelors..)
- Categorical → Nominal → Non-sortable e.g movie genres (we cannot say action is greater than romance, ...)

Why do these types matter ?

→ Some statistical analyses can be made only to certain data types and also some visualisation methods are valid only for some data types.

Sample vs Population:

- **Population:** All the Data e.g Salaries of all employees in a department, all facebook groups (we can gather this data)
- **Sample:** Some amount of data (hopefully randomly sampled from the population, so that it would be sufficient representative of population)
e.g salaries of all employees in a country (we cannot gather this data.)

Why should we care about sample vs population ?

- Many statistical methods are designed either for sample OR population.
- Therefore, applying a procedure to the wrong data type leads to incorrect results and incorrect interpretations!!
- But, in practice, if you have a large sample size, then your results will not get affected much.
- Most of the time, you will work with sample data and you will try to generalise your sample results for the large population.

When should you create fake data ?

- Validate analysis methods.
- Learn the pros and cons of analysis methods.
- Controls signals, noise and effect sizes. (you cannot control the real data but you can control fake data)
- Think more carefully and critically about data.
- It allows one to gain great intuition about the statistical methods.

3. Visualising Data

Table of Content

[Bar plots:](#)

[Box and Whisker Plots:](#)

[Histograms:](#)

[PIE CHARTS:](#)

[LINEAR SCALING VS LOGARITHMIC SCALING](#)

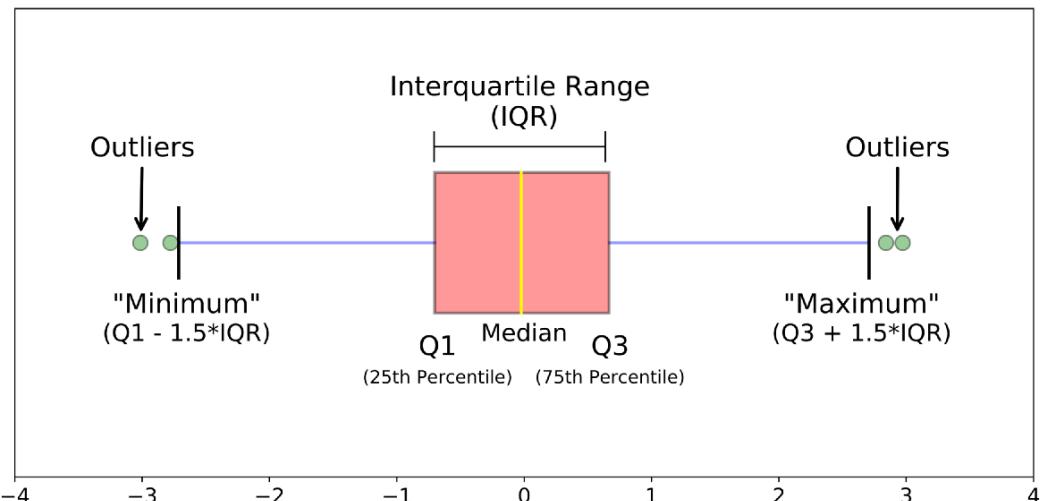
Bar plots:

- 1) What kind of data can be used in the bar plot ?
 - x-axis: Nominal OR Ordinal ||| y-axis -> continuous
 - OR if you want to plot continuous data on x-axis, then you can discretize the data (like create bins in histograms)
 - Also you must take care of the number of bins, and also the width of each bin.(Technically this is histogram)
- 2) What are error bars in bar plots ?
 - Error bars typically show std-dev, std. error OR confidence intervals around some parameter.
 - We will understand these terms more in-depth later. But whenever you plot error bars in bar plots, you should clearly explain what they represent.

[Code:](#)

Box and Whisker Plots:

[Source](#)



Now what you can do is, say you had a colour attribute which contains black and white. You can create one box plot of another variable for colour=black and one for colour=white. And then compare two box plots for more insights.

[Code:](#)

Histograms:

- 1)
 - Difference between bar plots and histogram is: Bar plots have categories on x-axis and histograms have binned continuous values on x-axis.
 - Another Important question might be, can you swap the locations of the bins on x-axis ?
 - If yes, then histogram OR bar plot, anything can be used.
 - If not, then because the shape of distribution in histogram is important and the shape of the distribution in bar plot is not.
- 2) You can have two types of histograms:
 - Histogram of counts (where you count on the y-axis)
 - Histogram of proportions (where you normalise the count and use the proportions on y-axis)
- 3)

Histogram of counts	Histogram of proportions
1) Can be meaningful to interpret (raw numbers) (we can say how many people have height between 80 and 90.)	1) Can take extra effort to relate to the raw data. (It will take extra effort to get the count of people between 80 and 90 as we have proportions on y-axis)
2) Can be difficult to compare across datasets. (Say person A went out and stored the height of 125 people and plot histogram and person B thought 125 is not enough, so B stored heights of 400 people and plot histogram. Now it doesn't make sense to compare the heights of two histograms, as they have different sample sizes.)	3) Easy to compare across datasets. (But if both person A and B convert their data into proportions (normalised) and plot the histograms, then it becomes easy to compare across datasets)
3) Does not need to sum to 1 OR 100%.	3) Sums to 1 OR 100%.
4) Usually better for qualitative inspection.(i.e usually better for one dataset)	4) Better for quantitative analysis.(better for multiple datasets.)

4) Converting counts to proportions:

Say we had heights of 125 people, and in one bin we had count=17. Now we want to convert this count to proportion

$$\text{so, proportion} = (\text{count} / N) * 100 \\ = (17 / 125) * 100 = 13.6\%$$

(Now instead of using count on y-axis we will be using proportion)

How many bins should we have in histogram ?

→ There are general guidelines, but they are based on descriptive statistics and therefore we will answer this question later.

When to use line plots vs histograms ?

- For the bar plots, the line doesn't make any sense.
- For the histogram, where we have a very large number of bins, the line chart can be used.
- You can have line charts, when you want to compare continuous data for two different groups.
e.g comparing heights of male and female using line charts, because histograms become hard to interpret.

[Code:](#)

PIE CHARTS:

What kinds of data can be used for pie charts ?

- Nominal, Ordinal and discrete.
- If you want to use continuous data for pie charts, make them discrete first use and then use pie charts.

[Code](#)

LINEAR SCALING VS LOGARITHMIC SCALING

- 1) You can scale your axis either linearly or logarithmically.
- 2) If you do log scale on one axis, then it is called log plot, if both of the axes are log scaled then it is called log log plot.

Linear scaling	Log scaling
1) Often easier to interpret.	1) Might need an explanation for the general audience.
2) Easily scaled for big, small or negative numbers.	2) Might not work with negative numbers.
3) Can obscure trends or comparisons across variables	3) Often appropriate for physics, finance, etc. for large differences between the numbers.
4) Use unless you have good reason to use log plots.	4) Use only when required

[Code](#)

4. Descriptive Statistics

Table of Content

Difference Between Descriptive Statistics Vs Inferential Statistics
Difference Between Accuracy, Precision And Resolution:
Data Distributions:
Central Tendency
Measures of Dispersion
QQ PLOTS (Quartile-Quartile)
STATISTICAL MOMENTS:
How many bins ?
Violin Plots
Entropy

Difference Between Descriptive Statistics Vs Inferential Statistics

DESCRIPTIVE STATISTICS	INFERRENTIAL STATISTICS
1) Describe the characteristics of the dataset like mean, median, mode, skewness, kurtosis, etc.	1) Use the features of the dataset to make the claims about the population like calculating p-values, T/F/chi-square values, Hypothesis testing, confidence intervals, etc.
2) Note: There is no relation to the population OR generalisation to the other datasets OR comparing multiple groups.	2) Note: The entire purpose is to relate with the population OR generalisation with the other groups.

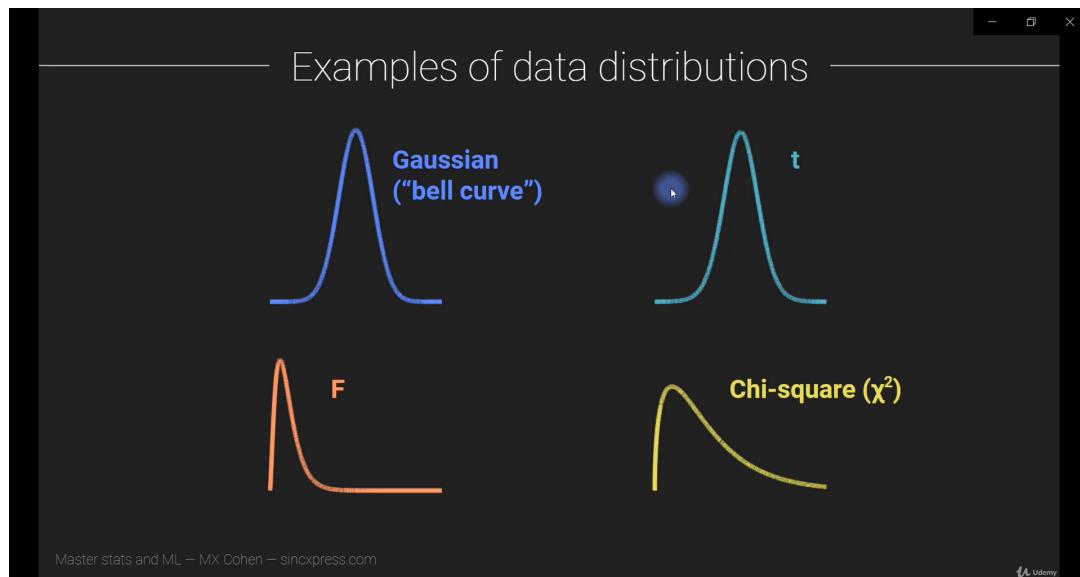
Difference Between Accuracy, Precision And Resolution:

- **Accuracy:** The relationship between the measurement and the actual truth
- **Precision:** The certainty of each measurement (inversely related to variance)
e.g say we are measuring the height of a person multiple times, each time we are getting the value which are almost similar, then precision is high.
If we get the values which are very different, then the precision is low.
- **Resolution:** The number of data points per unit measurement (time, space, ...)
Say we are measuring some data every one second and each second we are collecting 100 data points. Therefore our resolution will be 100 hz.

Data Distributions:

- The t-distributions are for evaluating the statistical significance.
- The f-distributions are generally used in general linear models like anova, regression.
- Chi-square distribution is highly used for inferential statistics.
- Note: f-distribution and chi-square distribution has only positive values including zero.

Code



Who cares about the distributions ?

- Most statistical procedures are based on the assumptions about the distributions.
- Data distributions provide insights about the nature of the data.

Beauty and simplicity of gaussian

- The formula for pdf of gaussian, in its basic simple form is e^{-x^2} , and when you plot this, this will be a nice bell curve.
- And this simple expression forms the basis for many of the things in the statistics like Central Limit Distribution.

Central Tendency

Central tendency is different from the "Expected Value".

Expected Value is the data times its probability of occurrence.

Mean:

Mean as valid descriptor in first diagram and mean as inappropriate descriptor in second diagram:

— Mean (a.k.a. arithmetic mean a.k.a. average) —

Formula: $\bar{x} = n^{-1} \sum_{i=1}^n x_i$

Suitable for:
roughly normally distributed
data.

Suitable data types:
Interval, ratio

"Failure" scenarios:

Master stats and ML – MX Cohen – sincxpress.com

Even though in the second case we can calculate the mean mathematically but it does not make sense

Is the mean suitable for discrete data types ?

e.g Average US family has 1.9 children.

Here 1.9 children doesn't really make sense, you have to interpret it correctly when using discrete data.

Is the mean suitable for Ordinal data ?

e.g Average course rating has 4.3/5

The mean for ordinal data is not suitable because the difference between 3 stars and 4 stars is not equal to the difference between 4 and 5 stars, but still 4.3 rating is informative to make a decision.

Is the mean suitable for nominal data?

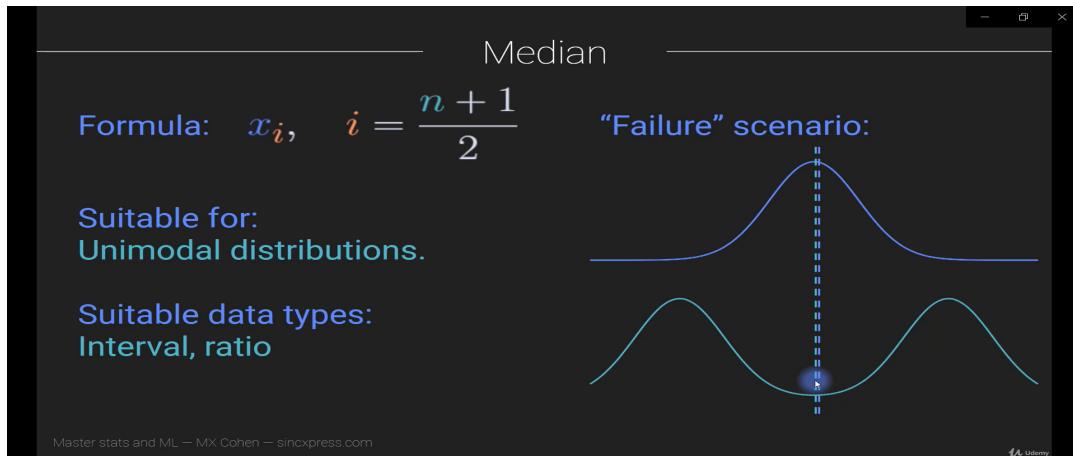
NO. e.g take the average of red, blue and green colour.

IT DOESN'T MAKES SENSE.

Therefore to conclude:

- 1) Mean is best applied for interval and ratio based data.
- 2) Meaning of discrete and ordinal data can be useful, but carefully interpret it.
- 3) The mean is not appropriate for nominal data.
- 4) Appropriate for roughly gaussian data.

Median:



- 1) Median is suitable for unimodal distributions and interval and ratio based types.
- 2) Robust to outliers.

Mode

- 1) Suitable for any distribution.
- 2) Suitable for any datatype, but continuous data should be converted to discrete first.
- 3) Mostly used for nominal data.
- 4) Also, it is possible to have multiple modes

Summary of Central Tendency

- 1) Use a mean with symmetric distribution. (unimodal)
- 2) Use median with skewed distribution (unimodal)
- 3) Use mode with multiple modes.

Code

Measures of Dispersion

Variance

[Image Source](#)

$$S^2 = \frac{\sum(x_i - \bar{x})^2}{n - 1}$$

- 1) Suitable for any distribution but they are easily interpretable when distribution is unimodal and roughly gaussian.
- 2) Suitable for numerical and ordinal types (but requires mean)
- 3) In the variance formula, we have $(x_i - \bar{x})^2$, this $(x_i - \bar{x})$ is called mean-centering.
- 4) **Why mean-centering in variance ?**
 - Because we want the dispersion around the mean.
e.g $d1 = [1,2,2,2,1]$
 $d2 = [101, 102, 102, 102, 101]$ (basically $d1 + 100$)

The variance for both $d1$ and $d2$ is the same and that's what we want.

- 5) **What if we don't square the term ?**
 - The variance will ALWAYS BE ZERO
- 6) **Why not take an absolute difference ?**
 - Squaring emphasises large values, which is good for better optimization (continuous and differentiable), it is also the second "moment" of the distribution and also squaring has some nice properties.
- 7) Mean absolute difference is robust to outliers but less commonly used.
- 8) Mean absolute difference is also closely related to Median Absolute Difference which is used for outliers and cleaning data.
- 9) Dividing by $n-1$ is for sample variance and when the denominator is n , then it is population variance.
- 10) Population mean is theoretical quantity (because it does not change).
The sample mean is empirical quantity (as it depends on how you sample)
(therefore sample mean will not remain constant)
- 11) Suppose I roll the die 4 times, and then the mean calculated is 3.
Now my question is, how many data points of dice do you need to know to get the mean of 3 ?
 - You only need to know 3 values, and then you can calculate the 4th observation.
Say, you get the data from die as 1,2,4,x
$$\frac{1+2+4+x}{4} = 3$$

Therefore, $x=5$.

Hence, you need to know only $N-1$ values OR there are $N-1$ free values OR degrees of freedom.

As the 4th output is dependent on the previous output, therefore it is not free.

Standard Deviation

$\text{std} = \sqrt{\text{variance}}$

Fano Factor

Fano factor = variance / mean

Coefficient of variation

Coefficient of variation = std / mean

Note: whenever you use

`np.std(x)`, this basically calculates the population standard deviation and the same for `np.var(x)`

Therefore always pass the parameter `ddof` (denominator degree of freedom) therefore the correct way for the sample standard deviation is
`np.std(x, ddof=1)`

[Code](#)

Interquartile Range

$IQR = \text{Quartile3} - \text{Quartile1}$
= Q3 - Q1

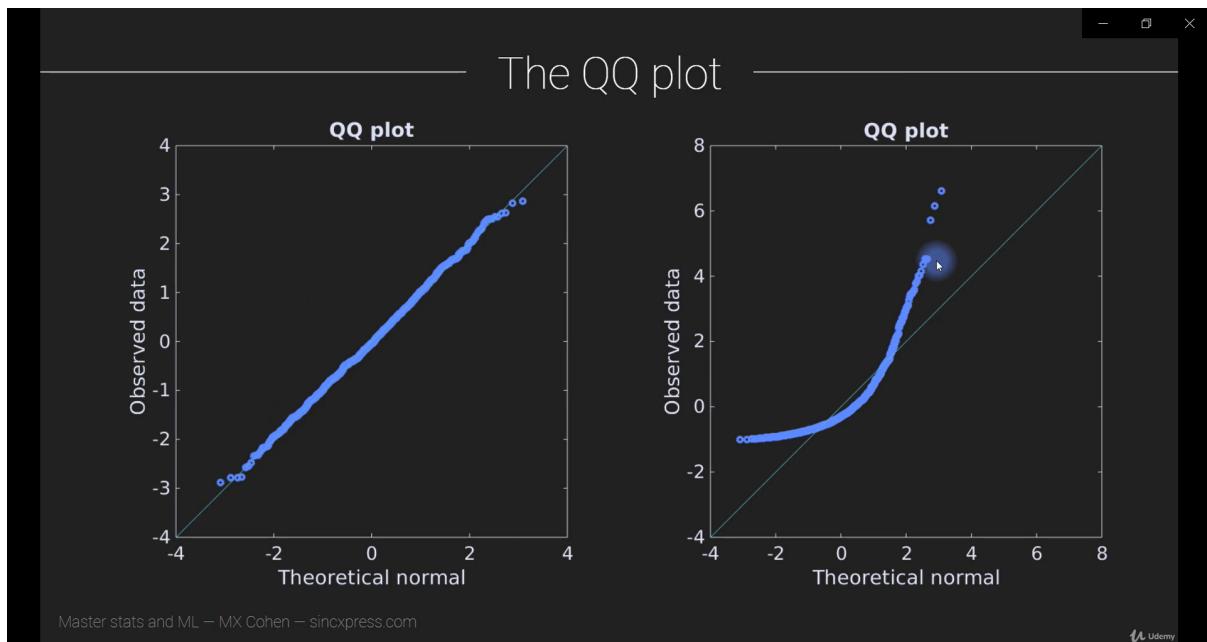
For detailed IQR: check this out

 [Range | Interquartile Range \(IQR\) | Box and whisker plot](#)

[Code](#)

QQ PLOTS (Quartile-Quartile)

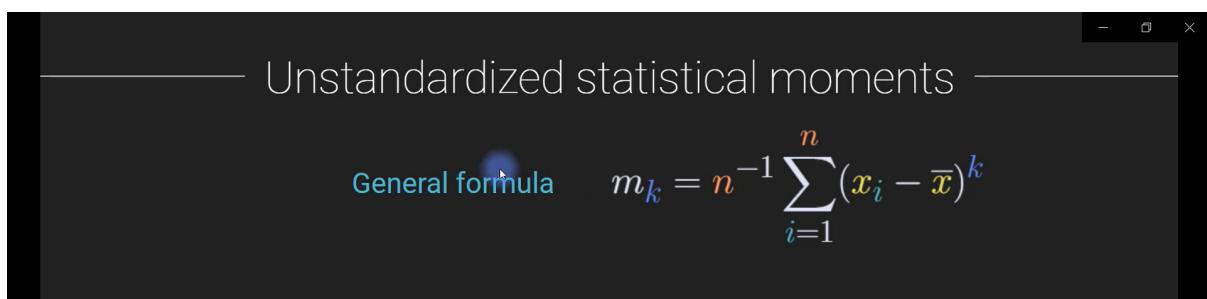
- Suppose you are given some data and I ask you to tell me whether the data came from a gaussian process?
- One way would be to plot the histogram along with the theoretical gaussian to compare.
- Another way is, you plot the observed data on y-axis and theoretical (generated) normal on x-axis and if your data is actually gaussian, then more or less the data points should lie on line ($y=x$) This is called QQ plot.
- This is useful for qualitative judgement.



Code

STATISTICAL MOMENTS:

- 1) Statistical moments bring different kinds of descriptive statistics under a generalised framework.



- 2) In the first moment, if $k=1$, then the m_1 will always be zero as we are subtracting x_i with \bar{x} . Hence, in the first moment we ignore the \bar{x} .

Unstandardized statistical moments

General formula $m_k = n^{-1} \sum_{i=1}^n (x_i - \bar{x})^k$

First moment: mean $m_1 = n^{-1} \sum_{i=1}^n x_i$

Second moment: variance $m_2 = n^{-1} \sum_{i=1}^n (x_i - \bar{x})^2$

Master stats and ML – MX Cohen – sincxpress.com

Udemy

- 3) The second moment is variance.
4) Skewness is the third moment and fourth moment is called kurtosis.

Standardized statistical moments

Mean: Average value $m_1 = n^{-1} \sum_{i=1}^n x_i$

Variance: Dispersion $m_2 = n^{-1} \sum_{i=1}^n (x_i - \bar{x})^2$

Skewness: Dispersion asymmetry $m_3 = (n\sigma^3)^{-1} \sum_{i=1}^n (x_i - \bar{x})^3$

Kurtosis: Tail “fatness” $m_4 = (n\sigma^4)^{-1} \sum_{i=1}^n (x_i - \bar{x})^4$

Master stats and ML – MX Cohen – sincxpress.com

Udemy

How many bins ?

Remember, in the visualisation section, in the Histograms, we have to decide the number of bins. So here are the few general guidelines for deciding the number of bins.

- First way is $k = \text{number of bins} = \text{ceil}((\max(x) - \min(x)) / h)$
 $h = \text{width of the bin you want}$
- The other ways are mentioned in the image. The Freedman-Diaconis guideline is generally used as it considers data count and the dispersion as well.

Guideline	Formula	Key advantage
Sturges	$k = \lceil \log 2(n) \rceil + 1$	Depends on data count.
Freedman-Diaconis	$h = 2 \frac{\text{IQR}}{\sqrt[3]{n}}$	Depends on data count and on data spread.
Arbitrary	$k = 40$	Easy to use.

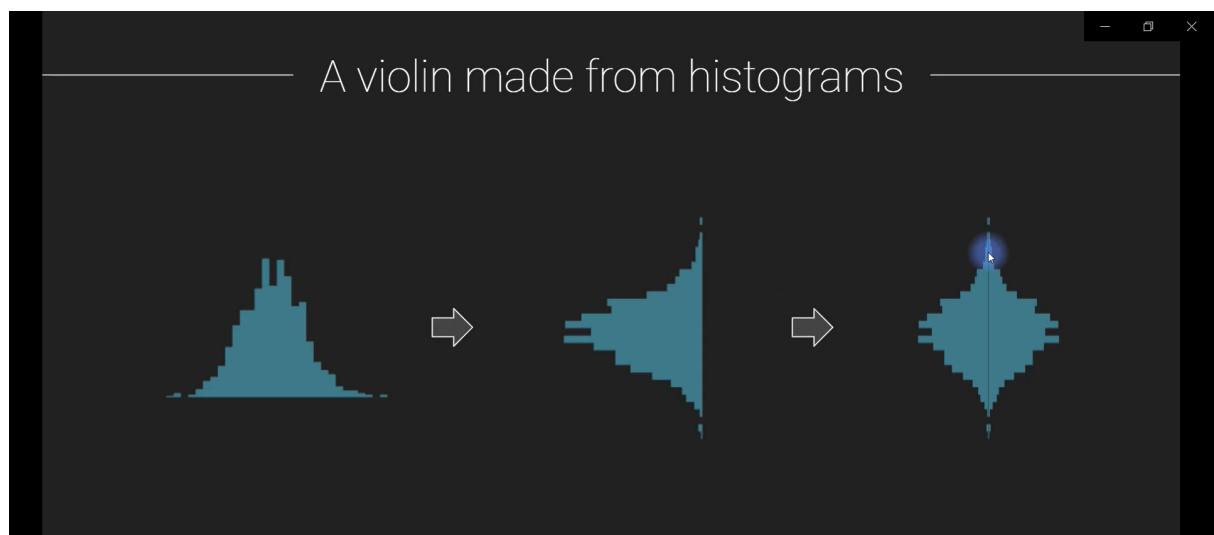
Master stats and ML — MX Cohen — sinxpress.com

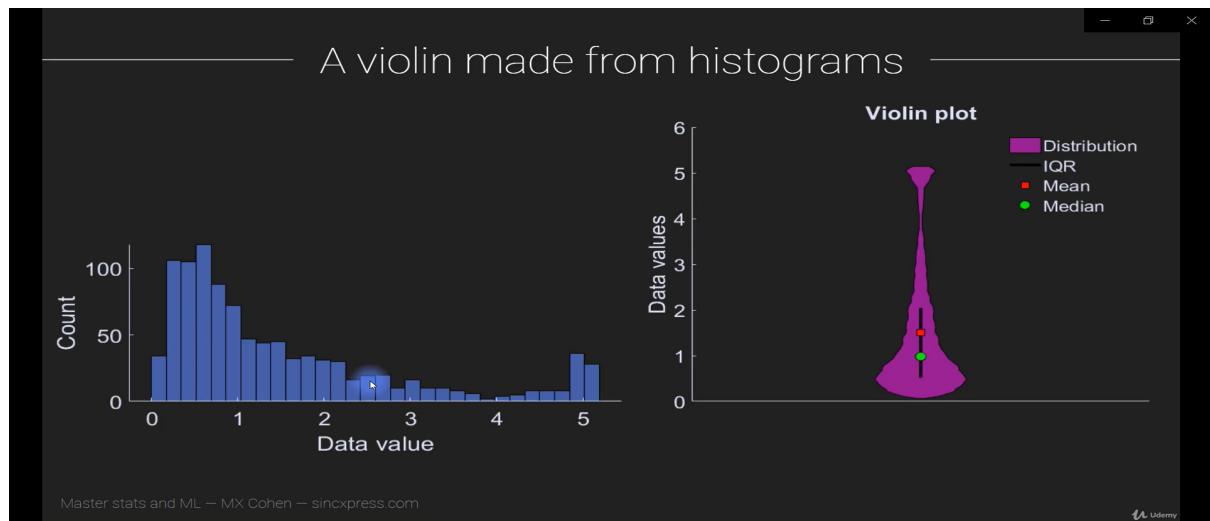
Udemy

Code

Violin Plots

- It is one of the beautiful way of illustrating the distribution of data
- Basically you start with the histogram, rotate it and then combine it with the mirror of it.
- And you can show the IQR, mean, median on the plot.





- What other thing you can do is, plot the distribution of one variable on the left side and distribution of another variable on the right and check the violin now. (Make sure that range of the two variables should be similar otherwise you will get weird shape.)

[Code](#)

Entropy

Check this video for GREAT explanation:

[Deep Learning\(CS7015\): Lec 4.10 Information content, Entropy & cross entropy](#)

- Basically, lower the probability of an event, greater the amount of information will be stored in that event. Therefore, from this we can say that information gain is inversely proportional to the probability.
- Formula for Information gain = $\log_2(1 / p) = -\log_2(p)$
Where p = probability of an event
- Entropy is the **Expected** amount of information gained.
entropy = 0
for i=1 to n:
 entropy += (-p(x) * log₂(p(x)))
p(x) = probability of event x occurring.
- What data types can be used with this formula ?
→ Nominal, Ordinal, Discrete

Why not interval or ratio ?

1. The probability of any continuous variable at some specific point is zero, hence entropy is not for these types.
 2. The other option would be to convert continuous variable into discrete (convert it into histogram), then use entropy.
 3. **Note:** It is important that entropy would be affected by bin size and number of bins.
- Interpreting the entropy
 - 1) Higher the entropy, greater the amount of information in the data (i.e data has high variance)
 - 2) Lower the entropy, lower the information in data (i.e data has lower variance)
 - How is entropy different from variance ?
 - 1) Entropy is non-linear and makes no assumptions about the distribution.
 - 2) Variance depends on the validity of the means and therefore appropriate for roughly normal data.

[Code](#)

5. Data Normalisation and Outliers

Table of Content

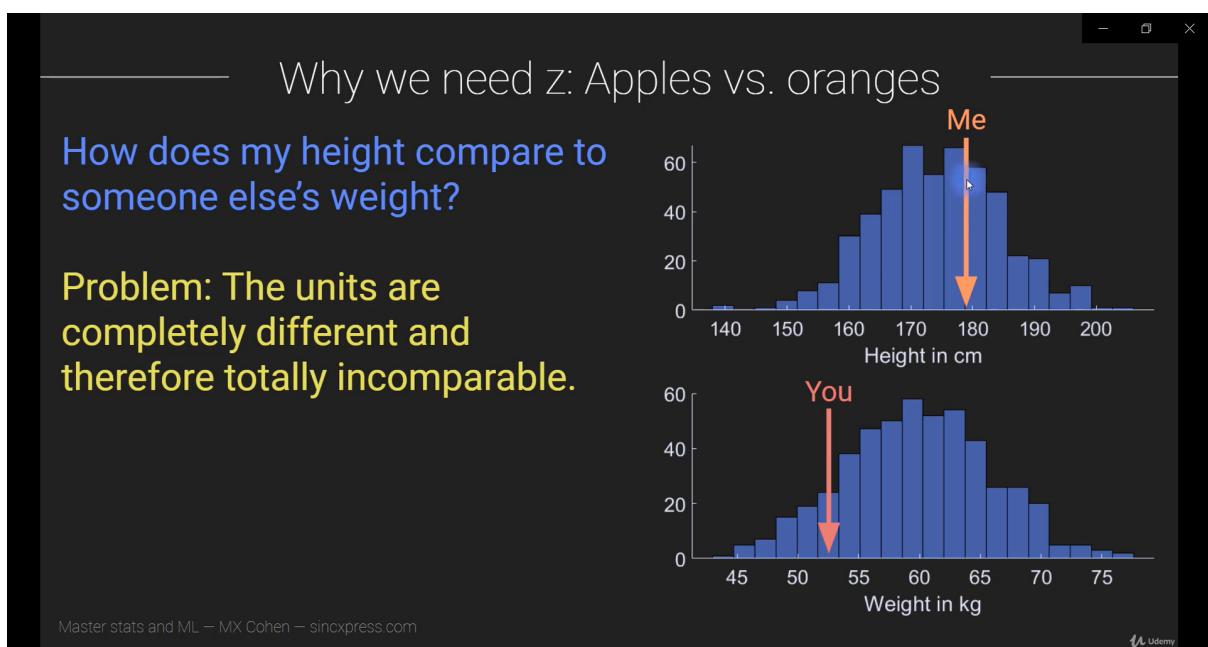
Garbage-In Garbage-out
Z-Score Normalisation
Min-Max Scaling
Outliers
Removing Outliers using Z-Scores
Modified z-transformation for non-normal distributions
Multivariate Outlier Detection using z-scores
Removing Outliers based on Data Trimming
Non Parametric Solutions to the Outliers

Garbage-In Garbage-out

- This means if you start with crappy data then you get crappy results.
- Also, awesome data doesn't guarantee awesome results. (maybe the hypothesis is not correct)
- So, when possible, always clean data.
- **Note:** While cleaning the data, be patient, critical, use domain knowledge and be unbiased.

Z-Score Normalisation

- Why do we need a z-score ?
 1. Say we want to compare my height (in cms) with your weight (in kg)
 2. The problem is, the units are completely different and therefore totally incomparable.
 3. Therefore we need a way to compare things independent of data scale.
 4. So, maybe we don't need to care about my height in cms, instead how my height compares with the population of other people's heights.
 5. So now we collect heights of 200 people and find out my height compares with the mean. (say my height is close to mean)
 6. Also we take the weight of the same 200 people and then compare my weight with the mean. (say my weight is one standard deviation less than mean)
 7. There now we can compare these two things based on how far I am from the mean.



- *This is the key insight:* The value on its own is difficult to interpret, but a value relative to its distribution is easy to interpret. Therefore we want to normalise the things into unitless distribution.
- That;s why we need z-normalisation

- Formula

$$z = \frac{x - \mu}{\sigma}$$

μ = Mean

σ = Standard Deviation.

$x - \mu$ = mean_centering

division by std: normalisation

- See, the units get cancelled out (say u divide cm / cm) and therefore z is the number of std. away from the mean of the distribution.
- Z-transformation shifts and stretches, but it keeps the overall shape of the distribution the same
- **What is the key assumption of z-transformation ?**
Mean and standard deviation are APPROPRIATE descriptors of central tendency and dispersion. Basically, it means the distribution is ROUGHLY normal.
- **Does this mean we cannot use z-transformation when you have a distribution which is not normal ?**
Well it depends on the specific application and if you are doing z-transform then you must be VERY CAREFUL WHILE INTERPRETING THEM.
- Remember that even if you have data with all +ve values, your z-scores might be -ve, therefore whether or not interpretation of -ve values is awkward, it depends on specific application.

Code

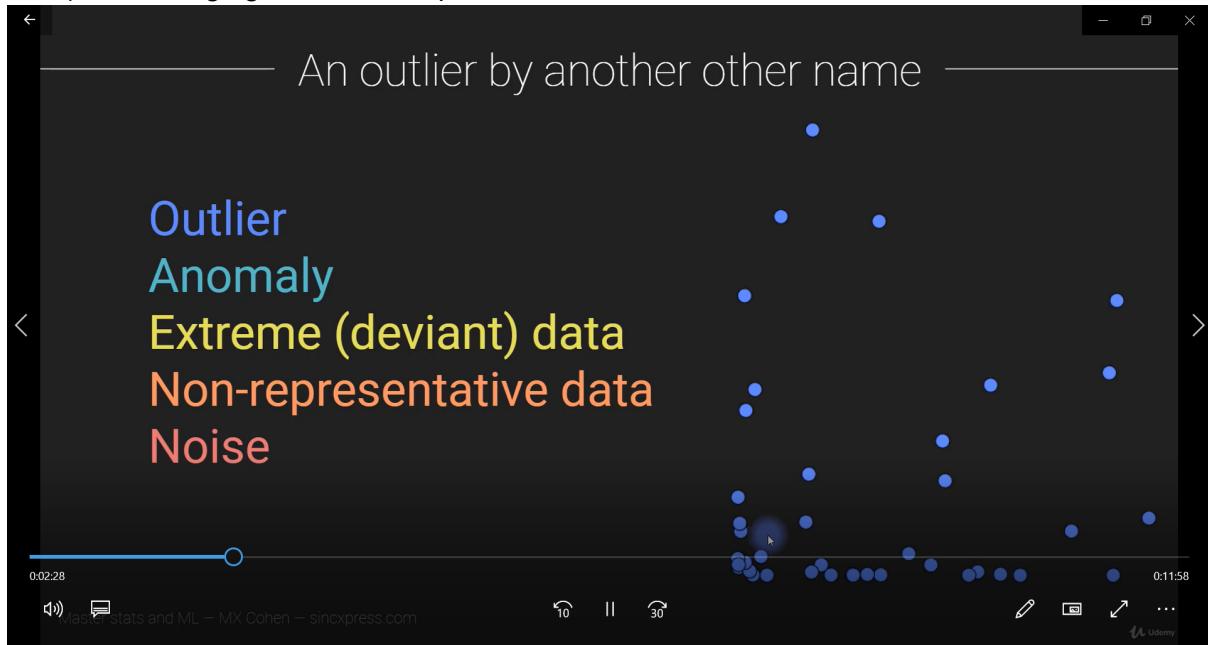
Min-Max Scaling

- Scale the data into any arbitrary range, most commonly used range [0,1] and scaling the data in [0,1] range is called unity-normed data scale.
- The shape of the distribution is preserved.
- This is not lossy transformation, this means from the transformed data we can get back to the original value.
- Formula: $x_{\text{transformed}} = (x_i - x_{\text{min}}) / (x_{\text{max}} - x_{\text{min}})$, this transforms the data into [0,1] range
- Now suppose we want to scale between any arbitrary range [a,b]
then, this would be 2 step
First use the above formula to scale between [0,1]
Then, scale to range between a and b →
 $x_{\text{betweenAandB}} = a + x_{\text{transformed}} * (b-a)$
- This is a useful transformation, especially when you are working on some analysis which requires the data into a certain range.
- You can also combine the two step formula in one like this
 $x_{\text{betweenAandB}} = a + ((x_i - x_{\text{min}}) / (x_{\text{max}} - x_{\text{min}})) * (b - a)$

Code

Outliers

- 1) The other names for the outlier are: Anomaly, Extreme (deviant) data, Non-representative data, Noise
- 2) There are cases when outliers are easy to identify and there are cases when outliers are not easy to identify. (check this image)
- 3) The things gets more complicated when we have multi-dimensional data

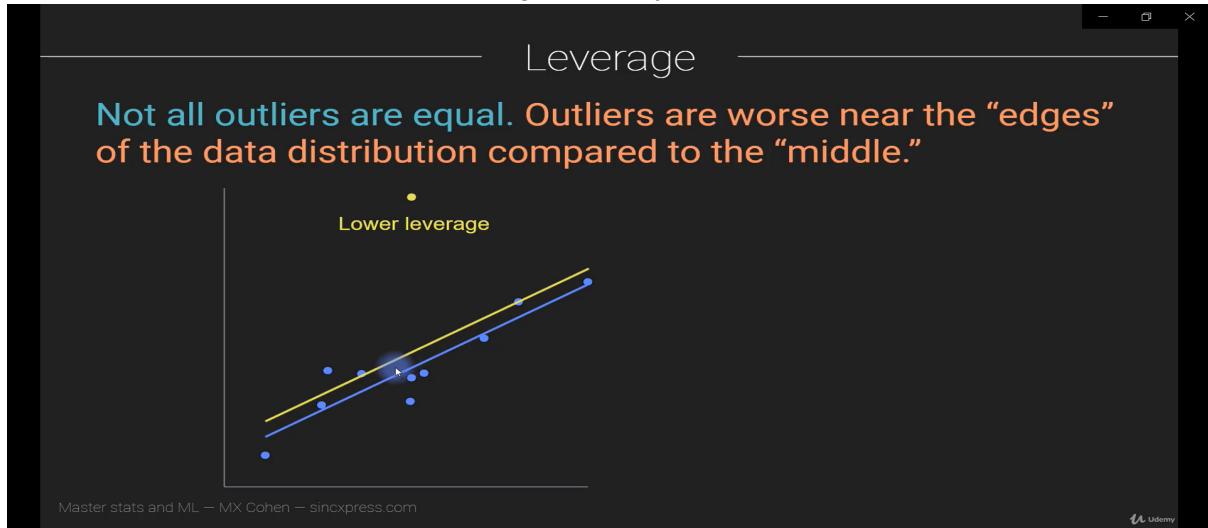


- 4) Where do the outliers come from ?
→ Noisy data, faulty equipment, human error, people filling the forms in a very weird way :D, OR it is natural variation in data.
- 5) Why are outliers bad ?
 - Most statistical analysis uses squared terms like (variance, anova, polynomials, GLM, correlation,etc.), so when you square them, the outliers become HUGE.
 - Therefore, the results of your analysis may become incorrect, OR become skewed just because of outliers.
 - Outliers can have more impact on analysis when the sample size is small.
 - Therefore, you should always worry about outliers even if you have large data and if you have small n, then you should be more worried.

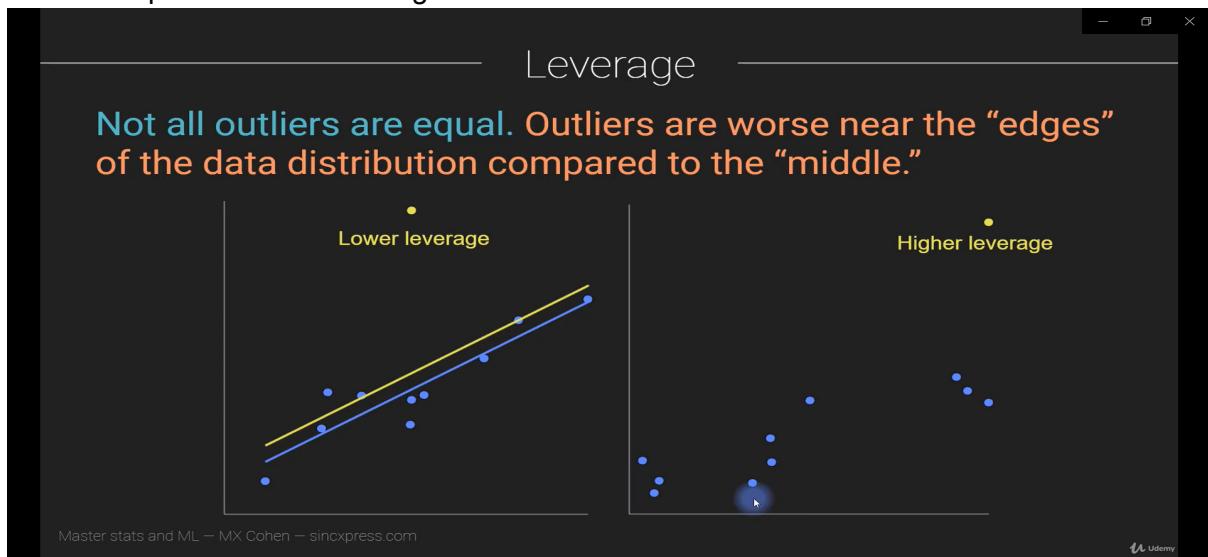
6) What is leverage ?

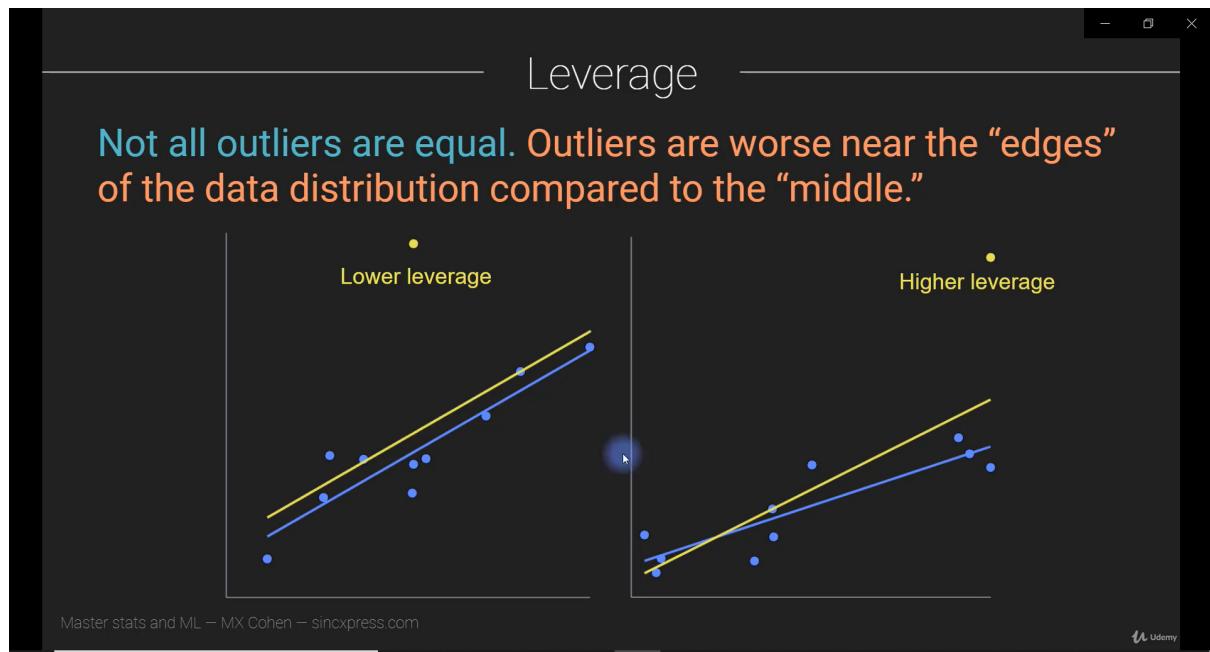
- Not all outliers are equal. Outliers are worse "near the edges (or ends)" of the data distribution and outliers near the middle of distribution are less worse.

The outlier in the image does not really have much impact on the slope of best fit, just it gets shifted. Therefore the slope of the line remains almost same with or without that outlier, therefore that outlier has "Lower Leverage (can say effect)"

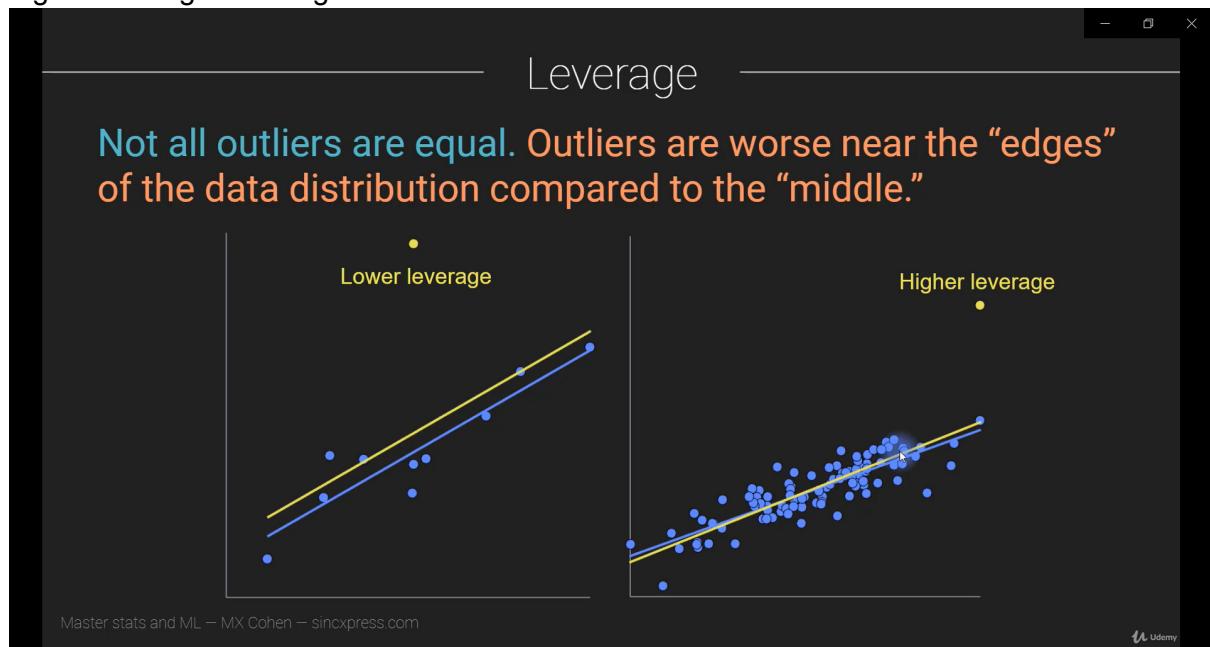


Now check the effect of the outlier with higher leverage. Also see, when the n was small, the impact on the result due to the outlier was quite high, and when we have large amount of data the impact decreases to a great extent.





Higher leverage with large amounts of data.



- 7) How to deal with outliers ?
 - There are two general strategies:
 - 1) Identify the outliers and remove them before running any analysis.
But you are making the assumption that your outliers are just noise or otherwise invalid.
 - 2) You keep the outliers in the data (even though you know that can negatively impact your results) and use the robust methods for the analysis that minimises the effect of outliers on the results.

Assumption: Outliers are unusual but valid data.

Robust methods like: non-parametric t-tests, permutation testing, spearman correlation, robust regression and iteratively weighted regression, etc

- 3) Therefore, you should always investigate the outliers and evaluate them.

8) Never remove outliers without any thought, think about why you are removing outliers and whether it is justified removing those outliers. **Outliers can be informative based on the experiment that you are running and they might be a crucial part of your research.**

Removing Outliers using Z-Scores

- 1) As said before, the z-score is appropriate if your distribution is roughly gaussian, if your data is strictly non-normal then you can use a modified z-score (which we will learn later).
- 2) Algorithm:
 - Convert your data to z-scores.
 - z-scores with "x" standard deviation greater are considered as outliers (typically $x = 3$, but it can change based on use-case)
 - We can use z-scores because we are not changing the relationship between data points, we are not changing distribution shape, all we are doing is z-transformation.
- 3) Also, why typically $x=3$, and not 2.17 or 2.9 or 4, what is done in many cases, you plot the z-scores and inspect what could be the appropriate threshold for your data.
- 4) This algorithm can be extended with iterative z-score based outlier removal.
- 5) Algorithm:
 - Convert your data to z-scores
 - z-scores with "x" standard deviation greater are considered as outliers.
 - Repeat steps 1 and 2 until there are no more outliers.
- 6) But I am not a huge fan of this iterative method, because there can be few data points which are at the edge of distributions which may not be considered as outliers but after several iterations, they might get removed, even though it was a potential inlier.

Code

Modified z-transformation for non-normal distributions

- 1) Formula:

$$\text{modifiedZScore}_i = \frac{(0.6745 * (x_i - \text{median}(x)))}{\text{MedianAbsoluteDifference}}$$

MedianAbsoluteDifference

MedianAbsoluteDifference(MAD) is basically, you take each x_i , subtract it with the median(x), take the absolute difference, and you do this for all x_i 's then,
 $\text{MAD} = \text{median}(|x_i - \text{median}(x)|)$

It is like the mean absolute difference, just instead of mean, we are taking the median.

- 2) This 0.6745 acts as a normalizer and it also helps us to treat this modifiedZScores as the normal zScores, and we can use this method for outlier removal, also an iterative method.
- 3) 0.6745 is the value on the x-axis for normal distribution when the area under the distribution is 0.75

Multivariate Outlier Detection using z-scores

Z-score method for multivariate data

Multivariate algorithm:

1. Compute the data mean.
2. Compute the distance from each data point to the mean.
3. Convert distances to Z-score.
4. Remove outliers based on threshold, as shown previously.

Master stats and ML – MX Cohen – sincxpress.com

Udemy

Code

Removing Outliers based on Data Trimming

Algorithm:

1. Mean centre the data
2. Sort the mean centred data.
3. Remove the extreme k data points OR remove the extreme $k\%$ of data points. (where k is the threshold you decide)

Data trimming: How it works

Algorithm:

1. Sort the mean-centered data.
2. Remove the most extreme k values, or the most extreme $k\%$.

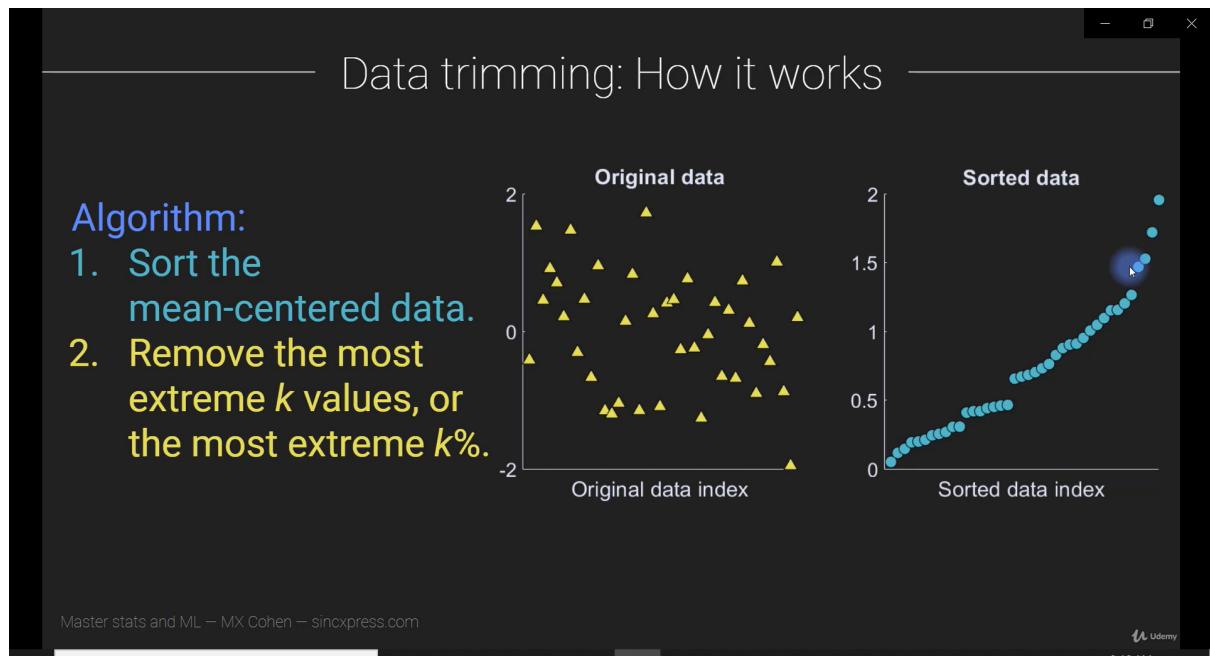
Original data

Original data index

Sorted data

Sorted data index

Check this image where it will also remove some of the non-outliers as well.



Code

Non Parametric Solutions to the Outliers

1. Remember that we discussed two general strategies of dealing with outliers, in the first we assumed that outliers are invalid data and therefore we removed them prior to any analysis. So till now whatever solutions that we discussed z-scores, trimming, modified z-scores follow strategy 1.
2. But there are some cases where we have to follow strategy 2, which assumes that outliers are unusual data but valid data. Therefore we have to use some robust methods to decrease the effect of unusual data on the analysis.
3. Why are these non-parametric methods less sensitive to outliers ?
→ Because they are based on median OR ranks, which are insensitive to outliers.

6. Probability Theory

From now on, we will cover the topics from inferential statistics, which is based on probabilities.

Table of Content

Basics
Proportion vs Probability
Valid Data types for probability
Probability and Odds
Probability mass and density functions:
CDF
Creating Sample Estimate Distributions (V.V. Imp)
Monte Carlo
Sampling Variability, Noise And Annoyances:
Expected Value
Conditional Probability
Law of Large Numbers
Central Limit Theorem

Basics

What is probability ?

- It is a numerical description of how likely an event is to occur. It is a number between 0 and 1, where 0 means impossibility and 1 means certainty.

Now I will give you a statement and its three interpretations, tell which ones are right or wrong

Statement: There is a 20% probability of rain today.

Interpretation 1: It will rain for 20% of the day i.e $0.2 \times 24 = 4.8$ hours

Interpretation 2: There is a 1 in 5 chance that it will rain today.

Interpretation 3: We can be 20% confident that it will rain today.

- Now interpretation 1 is incorrect, interpretation 2 is correct, tell me why interpretation 3 is incorrect ?
- Because confidence is different from probability. We could be 99% confident of a 20% chance of rain. In this case, 20% is our parameter estimate, and the confidence interval might be [19%, 21%]
- Why is interpretation 2 correct ?
If we could repeat the exact same conditions today a very large number of times (approaching infinity), then it would rain on 20% of those days.

When do we need probability ?

- We need probability when there is uncertainty about the outcome of the event.
- Say, for example I say, tomorrow the sun will rise. Now in this statement, we know that an event WITH 100% CERTAINTY is going to occur. So, we don't need probability here.
- Say, some medical test stated that you have cancer, then you want to know what are the chances that you have cancer GIVEN THAT the test is positive. Here we will need probability.

Code

Proportion vs Probability

- 1) Probability: Likelihood of an event to occur.
Proportion: The fraction of a whole.
- 2) I will give one statement and three phrases, tell which statements are correct or incorrect.
Statement: I spend a total of 5.1 minutes each day brushing my teeth out of a total of 17 hours (1020 minutes)(when I am not sleeping).
Phrase 1: The proportion of my waking day, spent brushing my teeth is $5.1/1020 = 5\%$
Phrase 2: The probability that a randomly selected minute of my day involves teeth brushing = 5%
Phrase 3: The probability that I will brush my teeth during the day is 1.

All the three phrases are CORRECT.

Therefore the probability and proportion might be different depending on how the question OR how the statement is phrased.

- 3) **Statement:** Say I toss the fair coin 10 times, and we get 6 heads and 4 tails.
What is the probability of getting heads ?

→ The probability = 0.5, because the probability of getting heads will always remain the same irrespective of the number of heads we get.

What is the proportion of getting heads?

→ Proportion = 0.6, because remember the proportion is basically the fraction as a whole.

Say I used the same statement, then what is probability, out of 10 flips, I select any random flip and its outcome is heads ?

→ 0.6

Valid Data types for probability

- Remember that we have 5 types of data: interval, ratio, discrete, nominal, ordinal
- Now the probability is valid for discrete, nominal, ordinal and not for interval and ratio based data.
- What you have to do is convert them to discrete then calculate probability.
- e.g, what is probability of height of randomly selected men = 122.64982913874323
You cannot calculate the height of men with such precision and chances = 0
But we can calculate $p(\text{height} > 121 \text{ and height} < 123)$ and you can vary these intervals.

Also the probability is valid if the data have mutually exclusive outcomes.

For e.g when you toss a coin, you will get either heads or tails but not both at the same time.

When you roll a die, you will get a number from 1 to 6, and not 2 different numbers at the same time.

Probability and Odds

Image source

Odds is basically the ratio of events not occurring divided by event occurring.

It is commonly called the odds ratio.

$$\text{Odds} = \frac{\text{probability of the event}}{1 - \text{probability of event}} = \frac{P}{1 - P}$$

OR

$$P = \frac{\text{Odds}}{1 + \text{Odds}}$$

e.g what are the odds of drawing the king from the randomly shuffled deck of cards.

$$r = \frac{4}{52} = \frac{1:12}{48:52}$$

i.e the odds of drawing a king from a deck of cards is 1:12

Probability mass and density functions:

- **PMF:** For discrete events, generally bar plots or histograms are used for visualisation.
- **PDF:** for continuous events, line plot, with values on x-axis and their probabilities will be area under the curve.

- In statistics, applied statistics, ML, medicine, etc. we generally try to discretize things and hence we generally use pmf.
- Because, in computers you really cannot represent the analog signal with high precision, hence we try to estimate the true probability density function using pmf.

Code

CDF

————— Cumulative density function (cdf) ————

A cdf is the cumulative sum (or integral) of the probability distribution (or density).

The y-axis value at each x-value is the sum of all probabilities to the left of that x-value.

A cdf starts at 0 and increases monotonically to 1. The sum of the cdf is more than 1.

$$C(x_a) = p(X \leq x_a)$$

$$C(x_a) = \int_{-\infty}^{x_a} p(x_t) dt$$

$$C(x_a) = \sum_{i=1}^{a} p(x_i)$$

Master stats and ML — MX Cohen — sincxpress.com



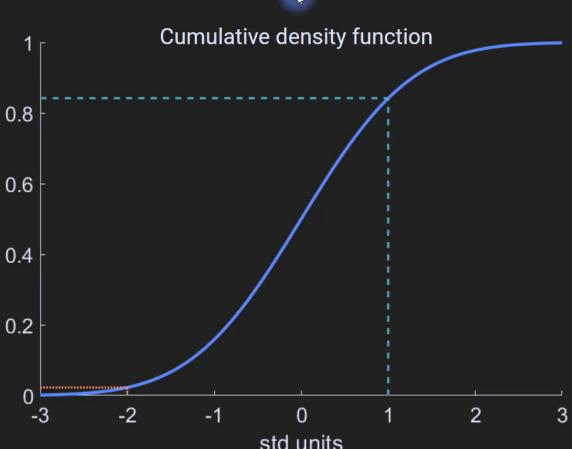
Cumulative Distribution Function : used for pmf
Cumulative Density Function : used for pdf

————— Cumulative density function (cdf) ————

cdf's are used to evaluate the probability of obtaining a value up to X or at least X.

Example: *What is the probability of getting at least 1 std higher on the SATs than average?*

Example: *What is the probability of an elephant weighing less than 2 std below the average?*



Master stats and ML — MX Cohen — sincxpress.com



The answer for first question:
 $0.15: p(X \geq 1) = 1 - p(X < 1) = 1 - 0.85$

The answer for second question:
 $p(x \leq -2) = 0.03$

Code

Creating Sample Estimate Distributions (V.V. Imp)

- It is very important for inferential statistics, for generalising from sample to population, which we often want to do.

- Let's start by asking one question

How tall are the giraffes ?

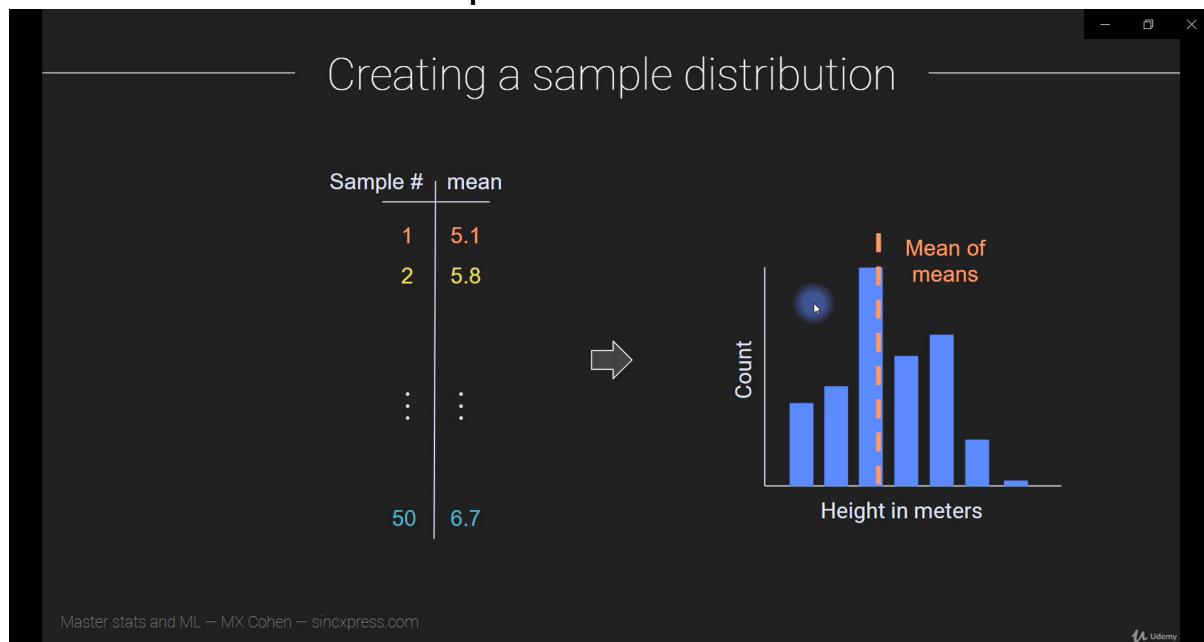
The statistical way of asking this question is
What is the population parameter of giraffe height ?
IT IS IMPOSSIBLE TO ANSWER THIS QUESTION!!!

- WHY?

→ Because we cannot measure the height of ALL the giraffes.

- Instead what we can do is, we can measure the heights of samples of giraffes.

1. So we start by randomly sampling giraffes from the whole population of giraffes and start measuring their heights.
2. Say we collected the heights of 100 giraffes and plotted their histogram along with the vertical line at mean.
3. IMPORTANT: This data distribution (HISTOGRAM) produces one sample estimate (mean)
4. Now we again sample the giraffes height (sampling with replacement i.e giraffe in the first sampling might be present in the second sampling also)
5. Then again plot the histogram and calculate the mean. (This mean might be different from the first one because of the sampling variability and so on)
6. We repeat this experiment N times, and calculate N means.
7. Now plot these N means and it is called **Sample Distribution** and the mean of **Sample Distribution** is called **Mean of the Means**.



NOTE: "mean" is just one example for parameter estimate. (it could be var, std, ..)

Now the basics are clear, let's go further and ask a new question.

Are giraffes taller than cats ?

- Remember that when we are quite certain about the event, then there is no need for probability. Here we don't need probability, but let's say we need probability.
- Now, we sample giraffes and sample cats (say sample size is 40), say the sample mean of giraffe is gi and of cats is ci , we calculate their difference and store it.

```
meanDiffs = []
for i = 1 to N:
    gi = sample(giraffe).mean()
    ci = sample(cats).mean()
    meanDiffs.append(gi - ci)
```

Note: you can do $ci - gi$, but this will generate -ve difference, which can lead to harder interpretation OR you can simply use absolute value, since you are interested only in magnitude and answering that question.

- Even though it is possible that in the real world there exists some difference in heights of giraffes and cats, due to sampling variability ,and other stuff, the difference might be zero.

In the first question, we were trying to estimate the population parameter but in the second question, we are trying to ask a question involving two groups of population. We will continue this discussion later, but the basics are clear.

Monte Carlo

- **Monte Carlo methods:** Solve really hard problems by sampling the solution space instead of actually solving the complete problem.
i.e we are determining what the possible solution could be based on the sampling of solution space.
- e.g suppose there is a very very hard integral to calculate, so instead of solving it, we use monte carlo methods.
- It is often used in physics, statistics, deep learning, even in pure maths.
- There are a whole bunch of monte carlo methods.
- "**Monte Carlo Sampling**" is basically the same thing as randomly sampling from the population space to estimate an unknown population parameter.
- "**Markov Chain Monte Carlo**" It is variant of Monte Carlo Sampling, basically the way you pick sample depends on the previous sample
Check this more in-depth article:
<https://machinelearningmastery.com/markov-chain-monte-carlo-for-probability/>

Sampling Variability, Noise And Annoyances:

- How tall is the Indian (the person who lives in India) ?
Don't worry about the actual answer.
The actual question is: How do we go about finding the answer?
- What we do is, the first person we meet, we measure their height and claim it as the average height of an Indian (160 cm), but we met another person whose height is 192cm, that's way above average we thought, and we keep measuring heights.
- But the internet says 183.2cm (population average). The answer on the internet must be true, but how can we have people with different heights ? You know the answer, it is due to variability. Not everyone has the same height.

- Therefore "**Sampling Variability**" means Different samples from the same population can have the different values of the same measurement (parameter like mean).
- **Implication of the sampling variability:** A single measurement may be an unreliable estimate of a population parameter.

Where does this variability come from ?

- "**Natural variation**": Often seen in biology (e.g weight, height) and physics (e.g earthquake magnitude, no. of stars, etc..)
- Also, sometimes the variability might be less (there is very less variation when you calculate the size of a proton multiple times)
- Therefore, sometimes variation might be large or small depending on the use case.
- "**Measurement noise**": Imperfect instruments, less precise instruments (say you are trying to measure weight in micrograms, but your system only measures in grams)
- "**Complex System**": Say I am measuring heights ignoring the age of a person, there might be very high variability, but if we keep the age as fixed, then your variability will be less.
- "**Stochasticity**"(Randomness): The universe is wild and unpredictable place (e.g photons hitting the lens of the camera)
- Some of the sources you can control like the quality of instruments, precision, etc.

What to do about sampling variability ?

- Taking many samples !!! Large number of samples will lead to better approximation of population parameters.

**The main outcome is, the mean of the sample means approaches closer to the population mean and it leads to two important concepts in statistics:
Law of Large Numbers, Central Limit Theorem.**

Code

Expected Value

- Average = $\text{sum}(x) / n$
- ExpectedValue $E[X] = \text{sum}(x_i * p_i)$, where p_i = probability of value x_i occurring.
- When the average and Expected value are equal ?
 - When all the p_i 's = $1/n$, i.e each and every data point is equally likely to occur.
 - ALSO, when drawing large and repeated REPRESENTATIVE random samples from the population.
- Average is an empirical sample estimate based on the finite data.
- Expectation value is the expected average in the population OR from a very very large number of samples(approaching infinity).

Expected value vs. average: example

"Illegal" die



Side 1: $p=\frac{1}{4}$

Side 2: $p=\frac{1}{4}$

Side 3: $p=\frac{1}{8}$

Side 4: $p=\frac{1}{8}$

Side 5: $p=\frac{1}{8}$

Side 6: $p=\frac{1}{8}$

$$E[X] = \sum_{i=1}^n x_i p_i$$
$$= 1p_1 + 2p_2 + 3p_3 + 4p_4 + 5p_5 + 6p_6$$
$$= \frac{1}{4} + \frac{2}{4} + \frac{3}{8} + \frac{4}{8} + \frac{5}{8} + \frac{6}{8}$$
$$= 3$$

Master stats and ML — MX Cohen — sincxpress.com

Udemy

Expected value vs. average: example

"Illegal" die



Side 1: $p=\frac{1}{4}$

Side 2: $p=\frac{1}{4}$

Side 3: $p=\frac{1}{8}$

Side 4: $p=\frac{1}{8}$

Side 5: $p=\frac{1}{8}$

Side 6: $p=\frac{1}{8}$

8 random rolls of the loaded die:

1 3 4 4 4 3 2 5

Average: 3.25

Expected value: 3

Master stats and ML — MX Cohen — sincxpress.com

Udemy

Conditional Probability

For quick explanation:

► Deep Learning Part - II (CS7015): Lec 16.0 Recap of Probability Theory Example

Conditional probability: example

You look through their sales data (hopefully anonymized...) and see that out of 100 random receipts, 42 people bought toilet paper, 60 people bought canned soup, 55 bought candy bars, and 24 people bought dried fruit.

Further inspection reveals that 11 people bought toilet paper and canned soup, 32 people bought toilet paper and candy bars, and 21 people bought toilet paper and dried fruit.

Which product should t.p. be marketed with?

$$P(A|B) = \frac{P(A \cap B)}{P(B)}$$
$$A = \text{toilet paper}$$
$$B_1 = \text{canned soup}$$
$$B_2 = \text{candy bars}$$
$$B_3 = \text{dried fruit}$$

Master stats and ML — MX Cohen — sincxpress.com

Udemy

Conditional probability: example

You look through their sales data (hopefully anonymized...) and see that out of 100 random receipts, 42 people bought toilet paper, 60 people bought canned soup, 55 bought candy bars, and 24 people bought dried fruit.

Further inspection reveals that 11 people bought toilet paper and canned soup, 32 people bought toilet paper and candy bars, and 21 people bought toilet paper and dried fruit.

Which product should t.p. be marketed with?

$$P(A|B) = \frac{P(A \cap B)}{P(B)}$$
$$P(A \cap B_1) = 11/100 = .11$$
$$P(A \cap B_2) = 32/100 = .32$$
$$P(A \cap B_3) = 21/100 = .21$$

Master stats and ML — MX Cohen — sincxpress.com

Udemy

Conditional probability: example

You look through their sales data (hopefully anonymized...) and see that out of 100 random receipts, 42 people bought toilet paper, 60 people bought canned soup, 55 bought candy bars, and 24 people bought dried fruit.

Further inspection reveals that 11 people bought toilet paper and canned soup, 32 people bought toilet paper and candy bars, and 21 people bought toilet paper and dried fruit.

Which product should t.p. be marketed with?

$$P(A|B) = \frac{P(A \cap B)}{P(B)}$$
$$P(A|B_1) = \frac{.11}{.6} = .18$$
$$P(A|B_2) = \frac{.32}{.55} = .58$$
$$P(A|B_3) = \frac{.21}{.24} = .88$$

Master stats and ML — MX Cohen — sincxpress.com

Udemy

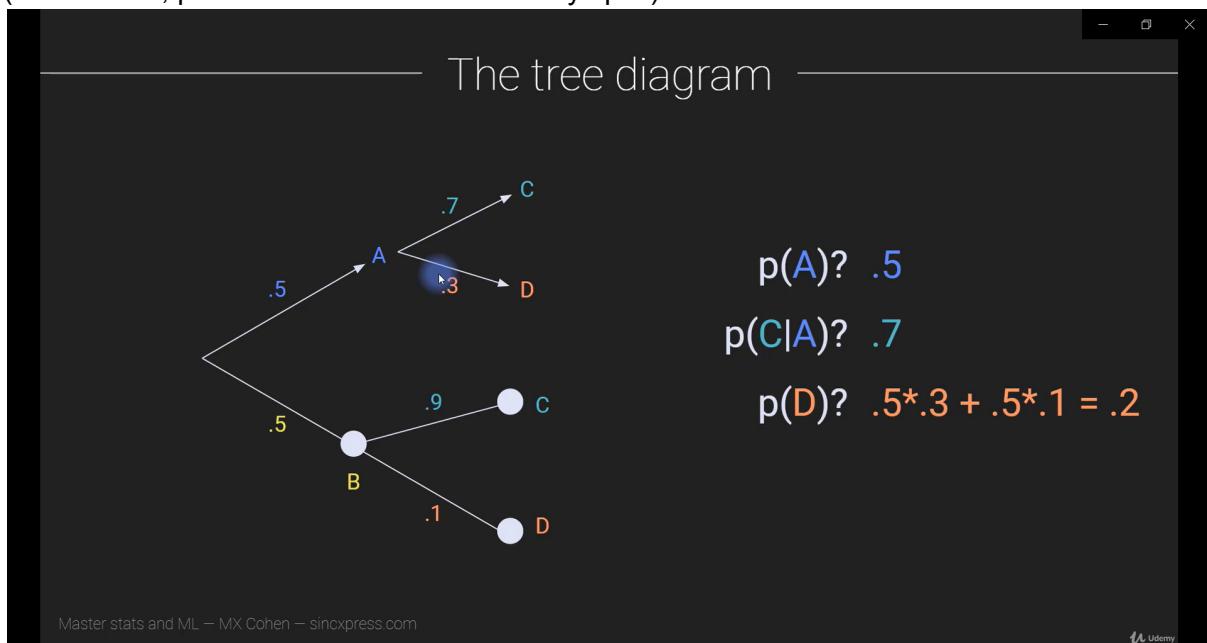
- $P(A|B1) = 0.18$ means probability of people buying toilet paper given that they have bought canned soup is 18%
- $P(A|B2) = 0.58$ means probability of people buying toilet paper given that they have bought candy bars is 58%
- $P(A|B3) = 0.88$ means probability of people buying toilet paper given that they have bought dried fruits is 88%

Be cautious, not to confuse the correlation with causation, therefore, it should NOT be interpreted like buying dried fruit causes people to buy toilet paper, what does $p(A|B3)$ means is, if someone buys dried fruits then the probability of that person buying toilet paper is 88%

- These types of questions can get a little bit tricky and it depends on how you phrase the question and how you set up the problem.
- For example, the supermarket should be advertising the candy bars with toilet paper because it should be clearly seen, out of 100, 32 people bought candy bars and toilet paper and only 21 people buy toilet paper and dried fruits.
- But if the question is, which product will maximise return on investment, then we can clearly say that people who buy dried fruits are 88% likely to also buy toilet paper.

Tree Diagram for Visualising Conditional Probabilities

(Notice that, probabilities sums to 1 at every split)



Code

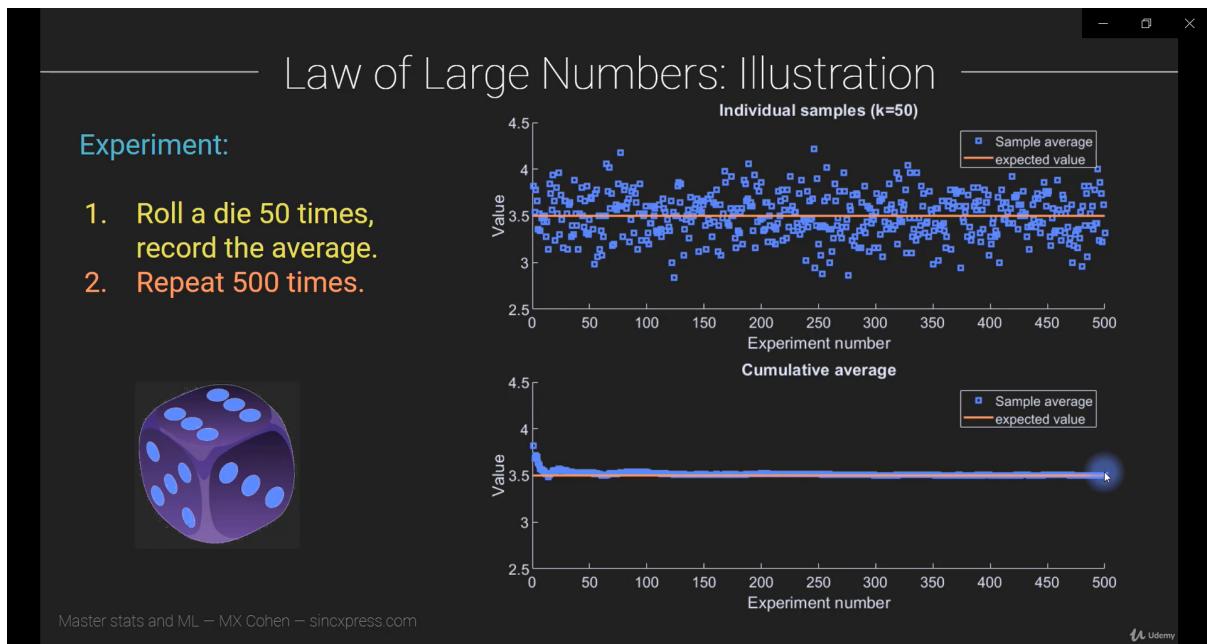
Law of Large Numbers

- Formal definition: As the number of experiment repetitions increases, the average of the sample means better approximates the population mean.

$\lim (n \rightarrow \infty) p(|\bar{x} - \mu| > \epsilon) = 0$
 n = number of times the experiment is repeated.
 \bar{x} = average of the n samples means.
 μ = population mean
 ϵ = very small number

We are making an absolute difference, as we are interested in how far we are and not whether we are underestimating or overestimating.

- As any one sample (or one experiment) is sensitive to sampling variability, noise and other sources of variation.
- But if we repeat the experiment many many times, our estimate will be really better. (independent replications of experiment)
- **Important question:** How large is large ? what does n (number of times we perform experiment) need to be to get a reasonable estimate.
→ Unfortunately, it depends on many factors like effect size, noise, system complexity, etc
- We will talk more about this in the statistical power and sample sizes section.
- If you also plot the sample means which you collected after n experiments, then you will find that it will be roughly gaussian. (Continued in Central Limit Theorem)



Code

Central Limit Theorem

- The DISTRIBUTION of the sample means (you can also take sample variance, entropy, or any other statistic) approaches the gaussian distribution, irrespective of the shape of the population distribution.
- Note: In law of large numbers, we calculate the average of sample means, but here we are talking about the distribution of the sample statistic.
- Why is the CLT so important ?
 - 1) Because many statistical analyses rely on the assumption of normality (or roughly gaussian).
 - 2) Gaussian distributions are easy to parametrize (mean, variance, skew, kurtosis) and compute confidence intervals.
 - 3) **Second Interpretation of CLT:** If we take Random samples from independent variables, their mixture will tend towards a normal distribution, even if the variables are non-normally distributed.
 - 4) **Independent Component Analysis (ICA)** is based on the assumption that the multivariate signals which are non-gaussian distributed, while random mixtures of signals are gaussian. (This is based on second interpretation of CLT)

Code

7. Hypothesis Testing

Table of Content

Stats Lingo
Model-Fitting In Statistics
Hypothesis and Models
Hypothesis proving ?
Sample Distributions under Null and Alternative hypothesis:
p-values, tails and misinterpretation
Common Misinterpretations of p-values
Degrees of Freedom
Type I and Type II Errors
Statistical Trade-offs
How to reduce Statistical Errors
Parametric and Non-Parametric Statistics
Multiple Comparisons and Bonferroni Correction
Type of Significances
Cross Validation

Stats Lingo

- DV: Dependent variable, the variable which you are trying to explain.
IV: Independent Variables, the explanatory variables, the variables that you HOPE will explain the variance in Dependent Variable.
- You generally manipulate the IVs and have control over them in the experiment and sometimes you don't have control over IV, you observe the IV from the real world and you HOPE that it will be able to explain the DV.
- Examples (Identify IV and DV):
 - 1) Effects of soil moisture on plant growth: Quite easy, IV = soil moisture, DV = plant growth
 - 2) Effects of time spent on facebook on irritability: This makes us think, IV = time spent on facebook, DV = irritability
Because we can take the other way as well, like the more you are irritated, the more time you spend on facebook. In this case, the IV = amount of irritation, DV = time spent on facebook.

- Whenever you are thinking about which is IV and which is DV, you should think about, "What is the likely flow of causality in the world".
- In statistics, you have to be very careful when discussing causality, just because we find an effect, it doesn't mean there is a causal relationship between them like in the second example.

3) Relationship between people who spend more money on clothes and go to bars frequently.

The *first case* can be: the people who buy new clothes tend to go to bars more frequently to show off their clothes.

i.e IV = money spent on new clothes

DV = frequency of going in bars.

The *other case* can be people who frequently go to bars, want to look more presentable so they end up spending more money on clothes.

i.e IV = frequency of going into bars

DV = money spent on clothes.

Therefore, sometimes it can go either way, your thoughts, use case, assumptions and how you want to interpret the data.

But sometimes it is clear, which is IV and DV.

Model-Fitting In Statistics

1) What is modelling ?

- A model is an equation(or set of equations) that explains some features in the dataset.
- For example, we want to know why there is a large variety in the height of the adults ?
- This is a really complicated question, there are many many factors that determine the height of the adult.
- Now the simplified version of this complicated question can be
$$h = a_1*x_1 + a_2*x_2 + a_3*x_3 + \text{epsilon}$$

$$h = \text{height}$$

$$x_1 = \text{gender}$$

$$x_2 = \text{parents height}$$

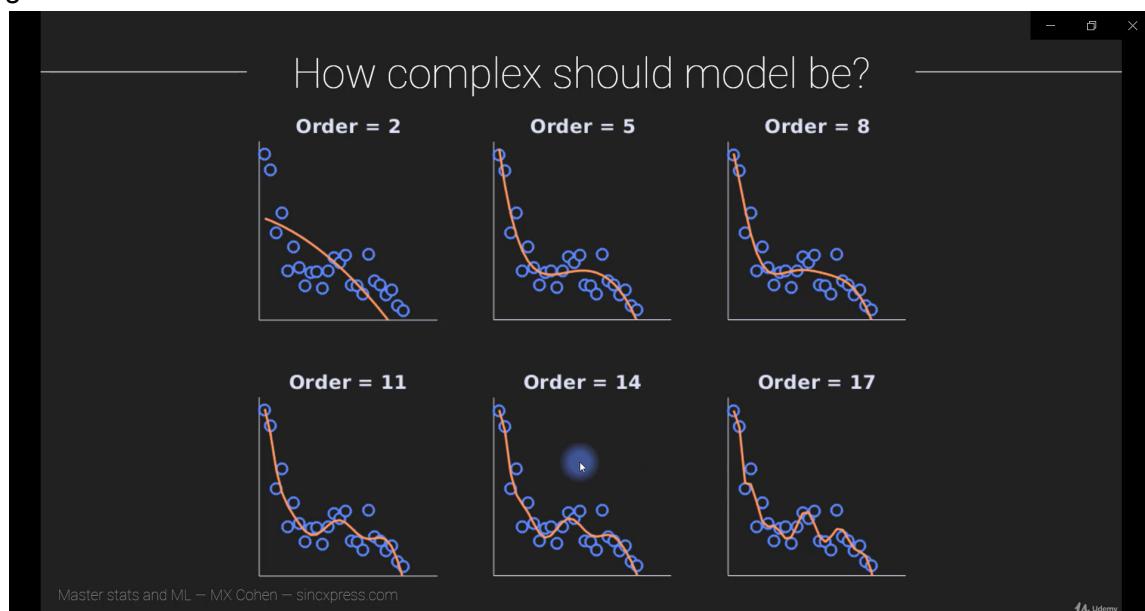
$$x_3 = \text{childhood nutrition}$$

Now we know there are many more variables that explain someone's height, but we are not trying to explain 100% of every feature about height, we want to capture the most important features in our simplified equation. Of course, these three variables are not going to explain every single pico-metre of height, therefore we add the term residual, all the differences which are not explained by our equation will go into residual.

IMPORTANT POINT IS:

- The residual should be small, but the model should be simple. (Easy to interpret)
- As the complicated models are too hard to interpret and also they can be unstable and also less generalizable.
- Therefore how you find balance between residual and simplicity is a bit hard.

Example: Order 2 is way too simple and order 17 is quite complicated and less generalizable.



Hypothesis and Models

- All the hypotheses are model comparisons.
- Hypothesis is about which model is better.
- When we say, model1 is better fit to the data than model2, that is an hypothesis and models are basically some equations that describe some real data.

What are hypotheses and how do you specify one

1) What is an hypothesis ?

- A falsifiable claim that requires verification, typically from experimental or observational data, and that allows us to predict future observations.

2) Why are hypotheses important ?

- Hypotheses improve the experiment design, critical thinking and very very important, data analysis.
- It allows you to transform your loose thoughts or ideas or assumptions into concrete and specific claims.

3) A strong hypothesis is

- Clear
- specific
- Falsifiable (can be proved wrong)
- Based on the prior data OR theory
- leads to a statistical test (can be tested with some statistical test)
- provides the prediction about the direction of the effect.
- relevant for the unobserved data OR phenomenon
- relevant for understanding the nature of data OR phenomenon.
- Statement and not a question.

4) Some examples of Not an hypothesis:

- Medical research is important for curing disease. (This is true and we 100% agree by this, but this is NOT an hypothesis)
- Will students pass this course (This is question and not the statement, we can think of generating hypothesis from this, but this thing on its own is not an hypothesis)

5) Some examples of Weak hypothesis:

- The medication has an effect.
(It is not a good hypothesis, but technically it is a hypothesis and we can try to falsify this statement, but it's not very clear and specific. And it's not relevant to understand the medication or disease. Say for example, you take medication for sleep at night, but it turns your fingers black. That would be an awful thing. The medication clearly has an effect, but the effect has nothing to do what the medication is for and it is awful effect)
- Studying improves grades (Obviously it is true, but it is not specific, and it is hard to contextualise this.)

6) Some examples of Strong hypothesis:

- The medication reduces the symptom X in the dose-response fashion.
(dose-response means the more the dose you take, the more symptom reduces)
(So this is clear hypothesis, it is specific, it is falsifiable, it leads to particular statistical test, and we can make prediction about unobserved data)
- A combination of self study and group study will improve grades by at least 10%.
(it is clear, specific, falsifiable, ..)

7) Now I will be giving a few examples and think about which is not, weak or strong hypothesis.

- 1) There are other universes with different physical rules.
- 2) Wearing purple underwear improves mood.
- 3) Plant grows differently in sugar-water
- 4) Mike's courses are awesome.
- 5) Washing hands for 20 seconds reduces the disease spread.
- 6) An apple a day keeps the doctor away.
- 7) Are people more creative after watching stand-up comedy ?

→

1. Not an hypothesis, because it is not falsifiable, we can't do experiments, we cannot collect the data or provide support or go against this statement.
2. Strong hypothesis, because it is clear, falsifiable, can be used for predicting future observations, a bit of a silly hypothesis but still it is strong.
3. Weak hypothesis: because it can be falsifiable, we can collect data and get the evidence, but it is not specific, because what does differently means, better or worse.
4. Not hypothesis
5. Strong: it is clear, specific, we can test various things, like is it really 20 seconds or 15 seconds and so on.
6. Weak hypothesis
7. Not an hypothesis: because it is a question and not the hypothesis, you can come up with a strong hypothesis using this question, but this question on its own is not an hypothesis.

Example of strong hypothesis from this question: people score 15% more on creative problem solving after watching stand up comedy.

Null Hypothesis

- Null hypothesis is the hypothesis that nothing interesting happened in the data (i.e no effect on process you did on data)
- In the research, you specify the alternative hypothesis.
- ***In statistical analysis, you always test null hypothesis and not alternate hypothesis.***
- Then the question raises, that do your data support the null hypothesis, the hope is, it do not support null hypothesis and we reject it in favour of alternate hypothesis
- Example:
Alternate hypothesis H1: People will buy more product X after seeing advertisement 1 as compared to the advertisement2.

Null Hypothesis H0: The advertisement type has no effect on purchases.

- Research is all about the alternative hypothesis, but in statistical analysis we try to quantify the evidence of the weak hypothesis using the data, if evidence is weak, then we reject it.

Hypothesis proving ?

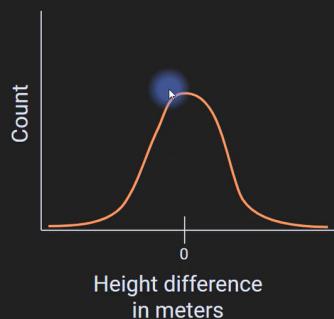
- 1) It is NOT possible to PROVE the hypothesis.
- 2) Hypothesis can be rejected OR failed to be rejected (interpreted as tentatively supported until a better hypothesis comes.)

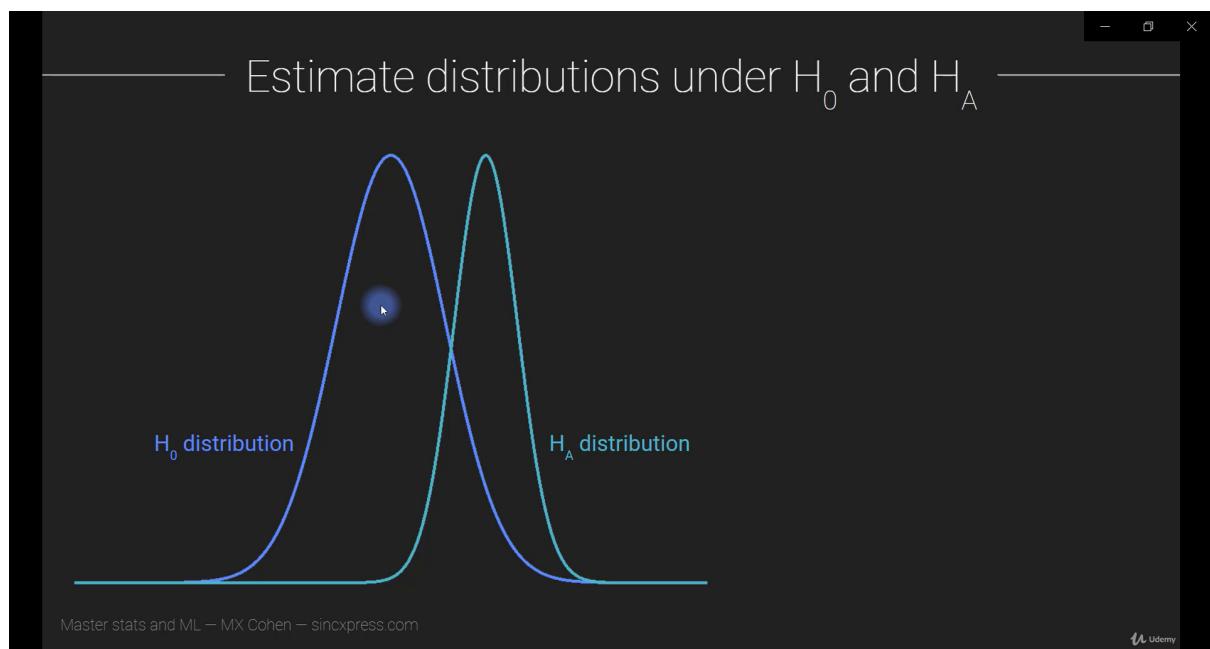
Sample Distributions under Null and Alternative hypothesis:

- Reminder: The distribution of sample means is gaussian (CLT)
- Lets create null hypothesis from the scientific question:
Are giraffes taller than giraffes ? Obviously the answer is no.
This is the null hypothesis question that we are asking here.
- The reason why we are setting up this question is to illustrate how to create the distribution (of sample parameter estimates) under the null hypothesis.
- Say we sampled the mean of heights for n_1 times and we get the distribution. The mean of the sample means was = 6 metres.
- Again we sampled mean of heights for n_2 times ($n_1 == n_2$ and different groups of giraffes), then the mean of sample means = 6.2 metres.
- Then, is it true that giraffes are 0.2 metres taller than giraffes?
NO. It means that on an average, giraffes(of group 2) are 0.2 taller than giraffes(of group 1).
The difference is because of sampling variability, noise, ...
- The point is, we expect the difference to be zero, but actually the difference is non-zero because of all causes of sampling variability.
- **Therefore, we are willing to tolerate some difference around zero.**

Creating a sample distribution of differences

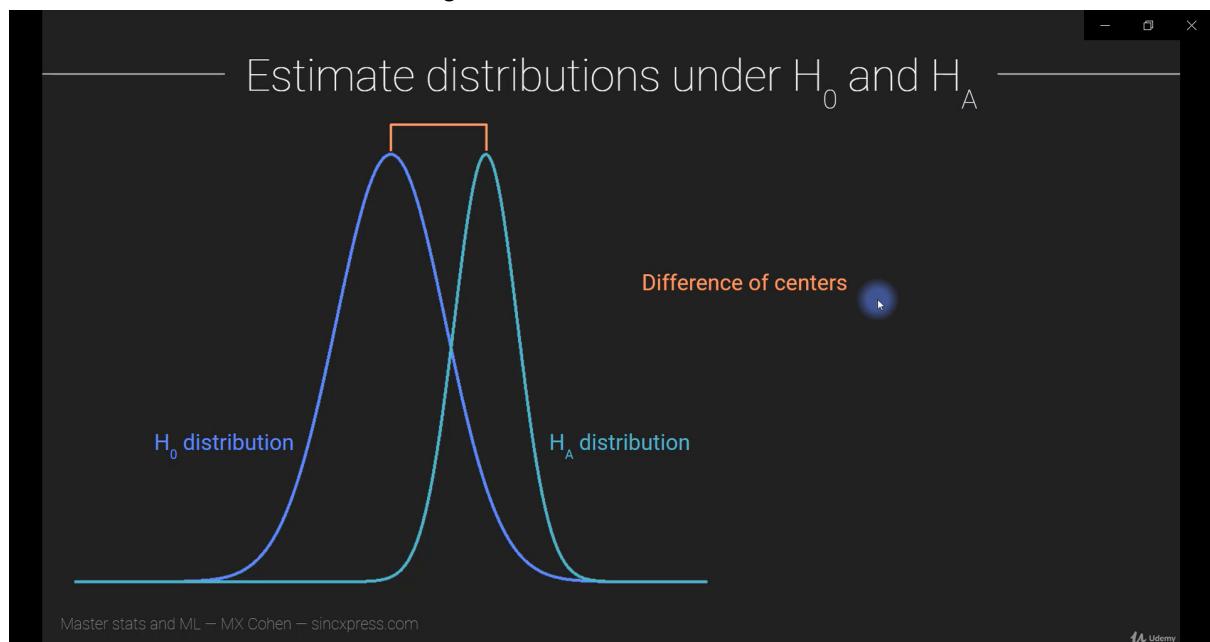
Scientific question: Are giraffes taller than giraffes?



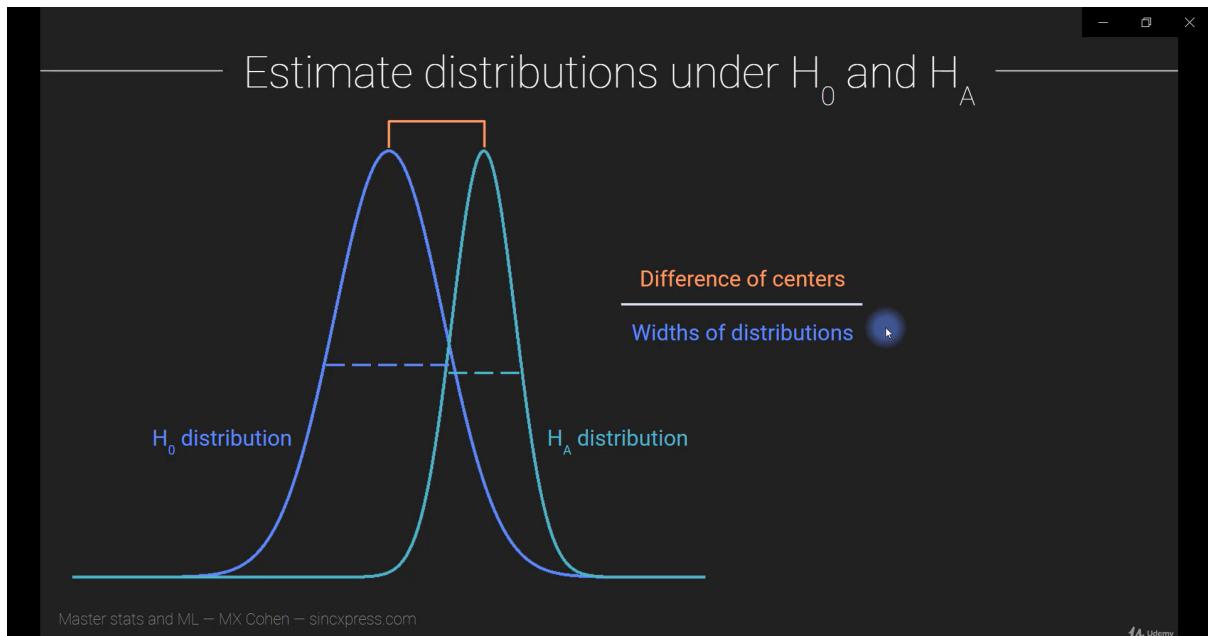


We want some probability that null hypothesis distribution and alternate hypothesis distribution are different from each other.

So one way we can say they are different based on comparing the difference in centres of distribution (means), but it will be scale dependent (metres), and we want a normalised way and also we will not be considering the width of the distribution.

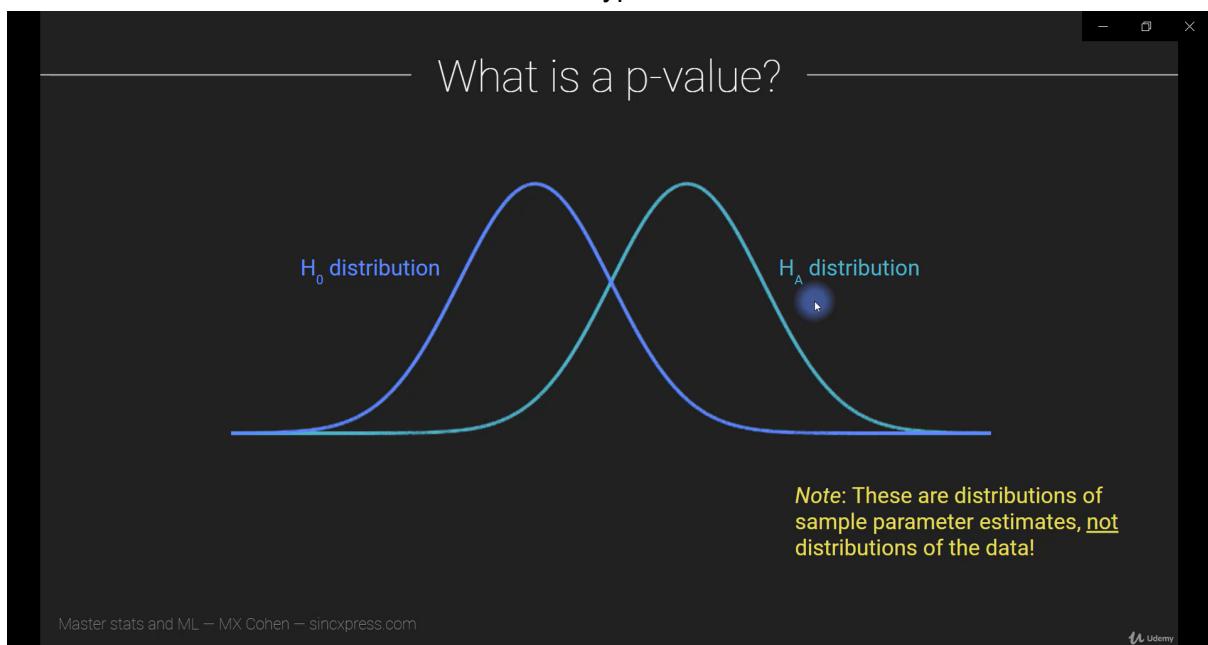


The other way could be, **difference_of_centers / width_of_distribution OR signal / noise**
Many statistical tests and basically the entire inferential statistics are based on this formula, these tests basically adapt this framework and add things, have different definitions of widths, based on some assumptions, but in general follow this framework.

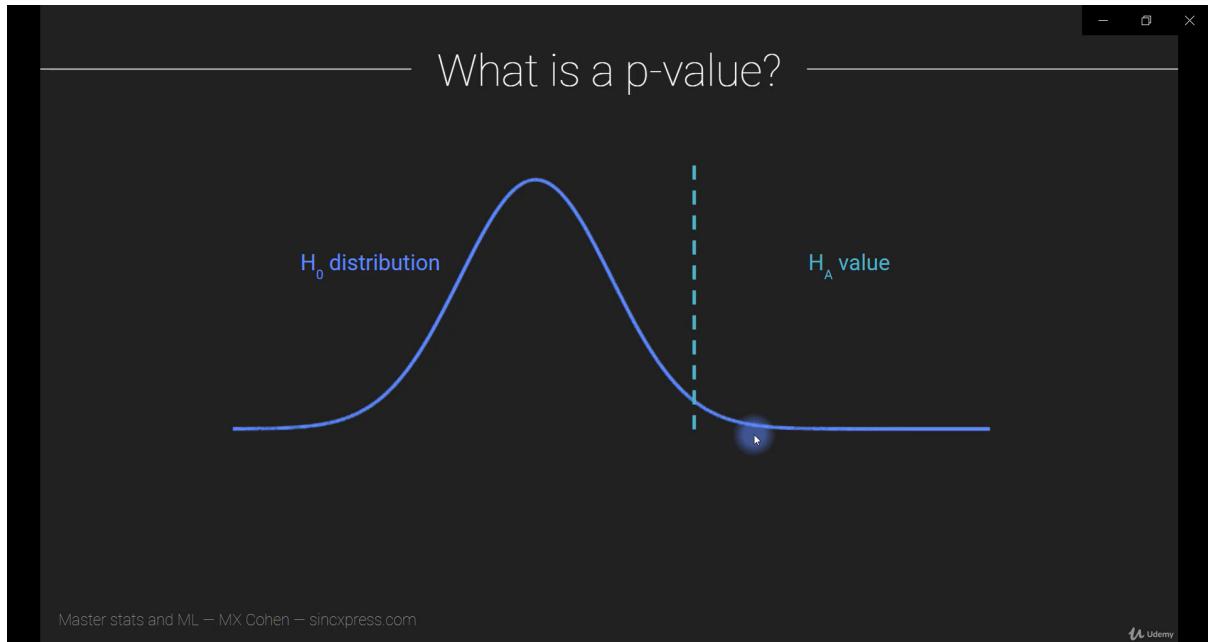


p-values, tails and misinterpretation

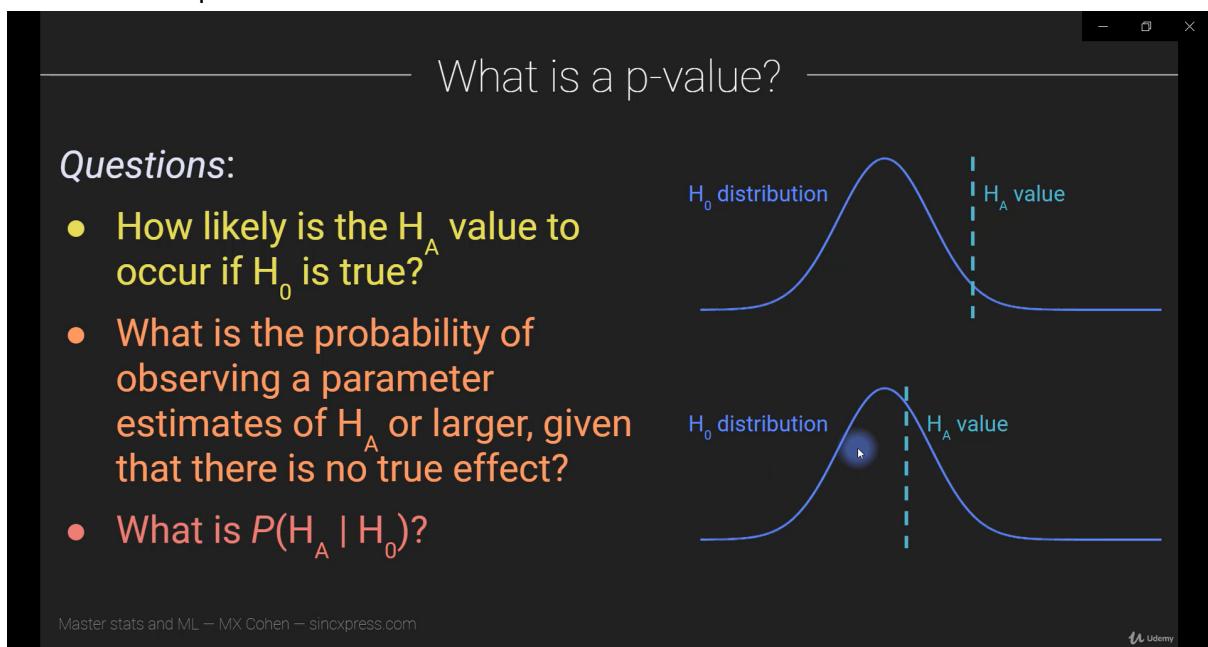
- 1) Null hypothesis distributions are based on some assumption, some formula, but we really don't know about the distribution of alternate hypotheses.



Since we don't know the distribution under alternative hypothesis



Answer these questions:



All three questions are the same.

p-value is the probability of observing the effect GIVEN the null hypothesis is true.

What is a p-value?

H_0 distribution

H_A value

Important concept:

We cannot prove that H_A is true. We can only compute the probability that the test statistic associated with H_A could be observed given that there is no true effect.

Master stats and ML — MX Cohen — sincxpress.com

Udemy

What is a p-value?

P-values are probabilities. They range from 0 to 1.

Values closer to zero indicate low probability of $H_A | H_0$, and values closer to one indicate high probability of $H_A | H_0$.

Master stats and ML — MX Cohen — sincxpress.com

Udemy

When is a finding “statistically significant”?

H_0 distribution

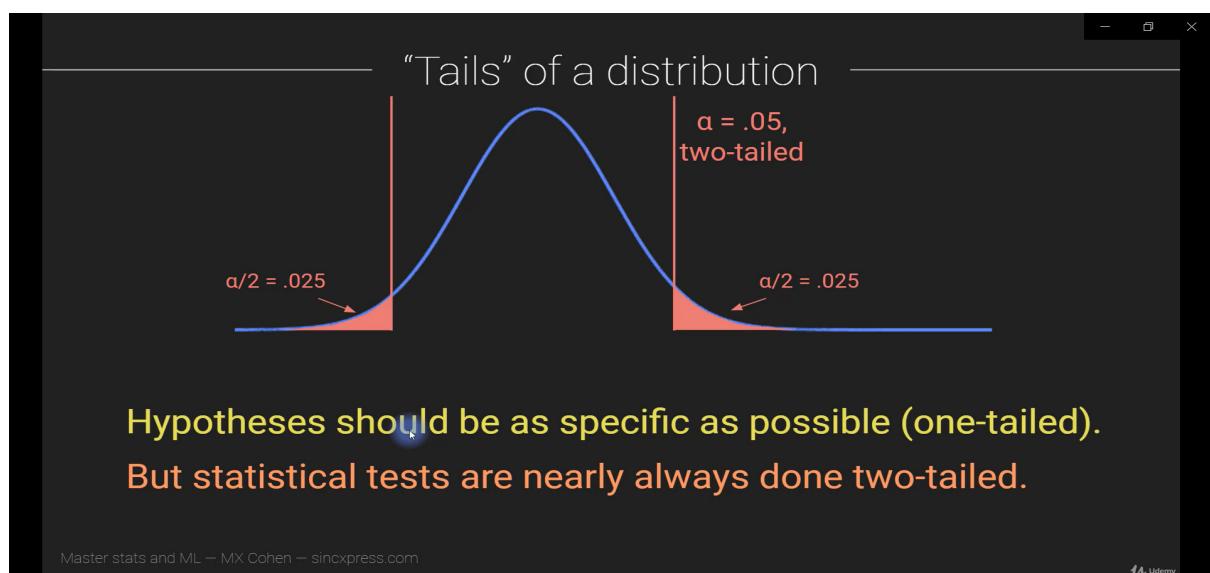
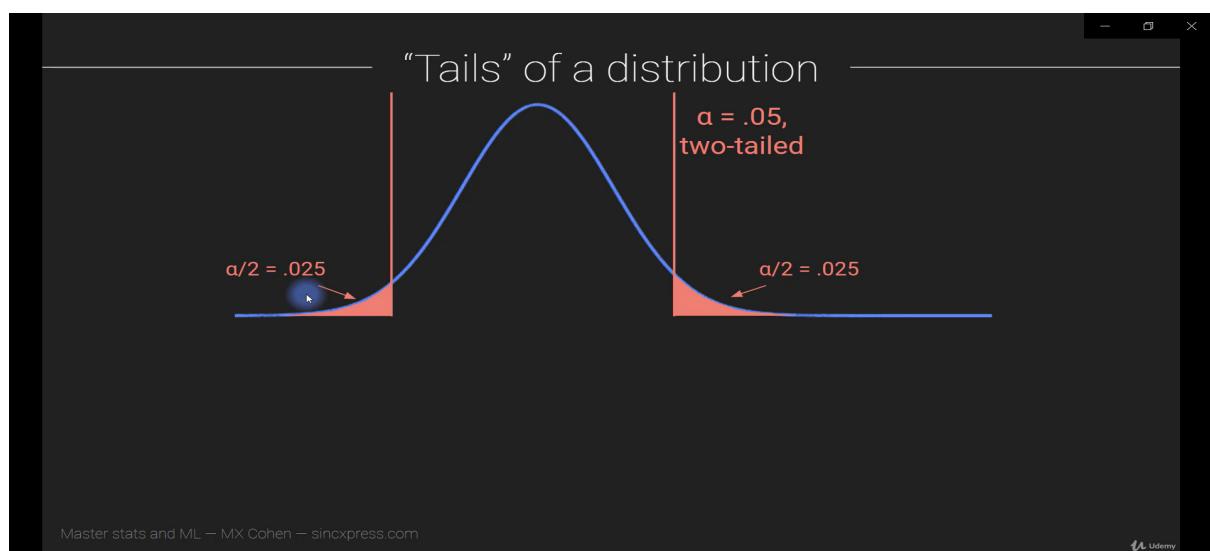
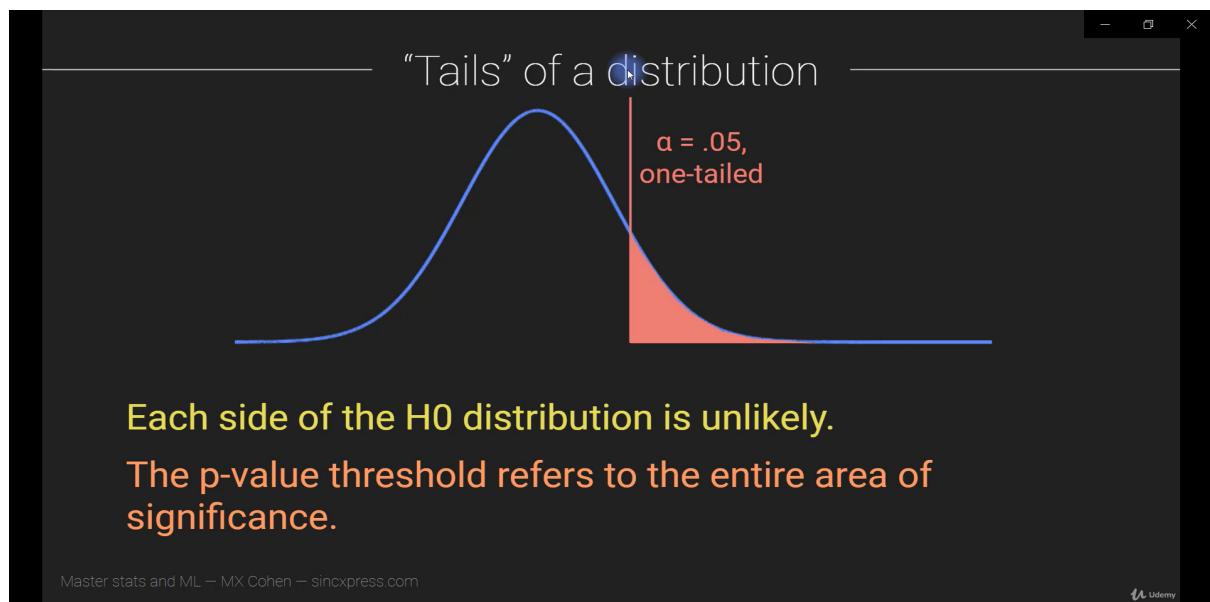
Significance threshold (α)

A finding is called “statistically significant” if the test statistic is greater than a threshold. That is, if $p(H_A) < p(\alpha)$.

Threshold is arbitrary; common values are $p < .05$ or $p < .01$.

Master stats and ML — MX Cohen — sincxpress.com

Udemy



Common Misinterpretations of p-values

Common misinterpretations of p-values

Incorrect:
“My p-value is smaller than the threshold, so therefore the effect is real.”

Correct:
“My p-value is smaller than the threshold, so it is unlikely that the effect in the sample would have been observed given the null hypothesis, assuming that the sample is representative of the population.”

Master stats and ML — MX Cohen — sincxpress.com

Udemy

Common misinterpretations of p-values

Incorrect:
“My p-value is .02, so the effect is present for 2% of the population.”

Correct:
“My p-value is .02, so there is a 2% chance that there is no effect and my large sample statistic was due to sampling variability, noise, small sample size, or systematic bias.”

Master stats and ML — MX Cohen — sincxpress.com

Udemy

Common misinterpretations of p-values

Incorrect:
“My p-value is larger than the threshold, so therefore the null hypothesis is true.”

Correct:
“My p-value is larger than the threshold, so it is *likely* that the effect in the sample would have been observed given the null hypothesis, assuming that the sample is representative of the population. There could be many other explanations for the p-value other than H_0 .”

Master stats and ML — MX Cohen — sincxpress.com

Udemy

— How to compute a p-value? —

There are several ways to compute a p-value, and they depend on the specific statistical test and assumptions made.

You will learn specific methods later in the course.

Importantly, the interpretation of a p-value is always the same, regardless of how the p-value was computed.

Degrees of Freedom

- Say you had four numbers and their mean, you know only the first three values and the fourth value is unknown, then you can calculate it easily.
- That is, the first three values can take on any values and the fourth value depends on the first three values, i.e it has 3 free values to take i.e it has 3 degrees of freedom.

— Degrees of freedom: numeric explanation —

$$x = \{a, b, c, d\}$$

$$\bar{x} = 5$$

$$5 = \frac{a + b + c + d}{4}$$

$$x = \{3, 4, 9, d\}$$

$$d = 4$$

Note: In this case, we were given the sample mean and that's why the degrees of freedom is 3

— Degrees of freedom: numeric explanation —

$$x = \{a, b, c, d\}$$
$$\bar{x} = 5$$
$$5 = \frac{a + b + c + d}{4}$$
$$x = \{3, 4, 9, d\}$$
$$d = \bar{x}n - a - b - c$$

This analysis has 3 degrees of freedom.

We know one outcome variable and there are four samples.

Any three samples are unknown; the fourth is necessarily known.

Master stats and ML – MX Cohen – sincxpress.com

Udemy

Now, say we were given the **Population Mean** and the four numbers, then the degree of freedom is 4.

— Degrees of freedom: numeric explanation —

$$x = \{a, b, c, d\}$$
$$\mu = 5$$

This analysis has 4 degrees of freedom.

The sample mean is not the same thing as the population mean.

The population parameter does not constrain the sample data values.

Master stats and ML – MX Cohen – sincxpress.com

Udemy

Why are degrees of freedom important?

Degrees of freedom (df) determine the shape of H₀ distributions.

Higher df generally indicates more power to reject the null hypothesis (related to statistical power).

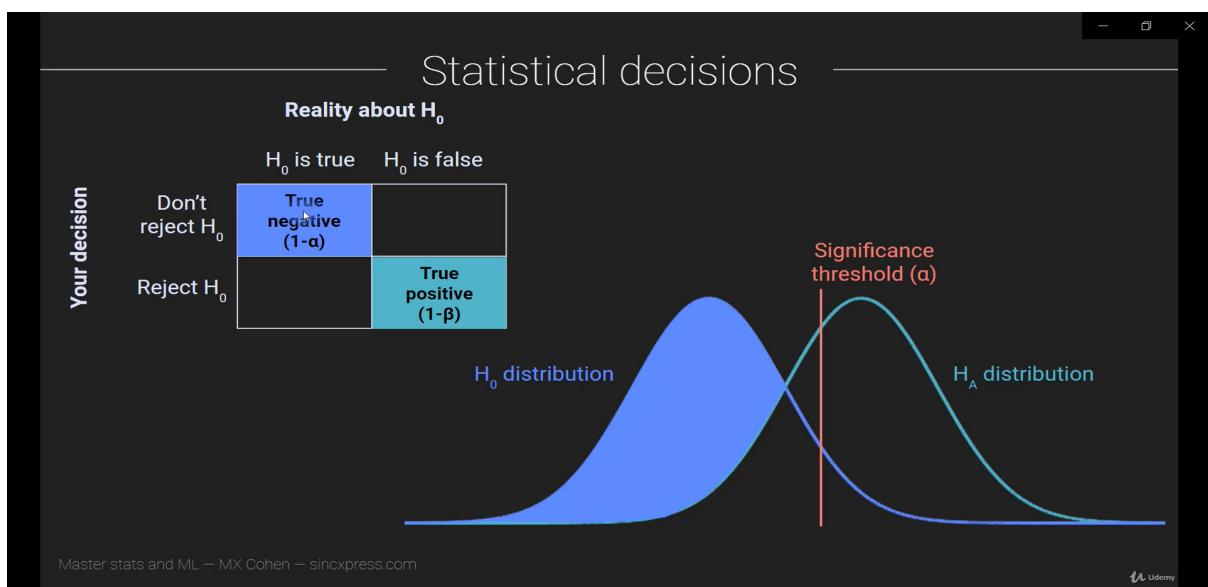
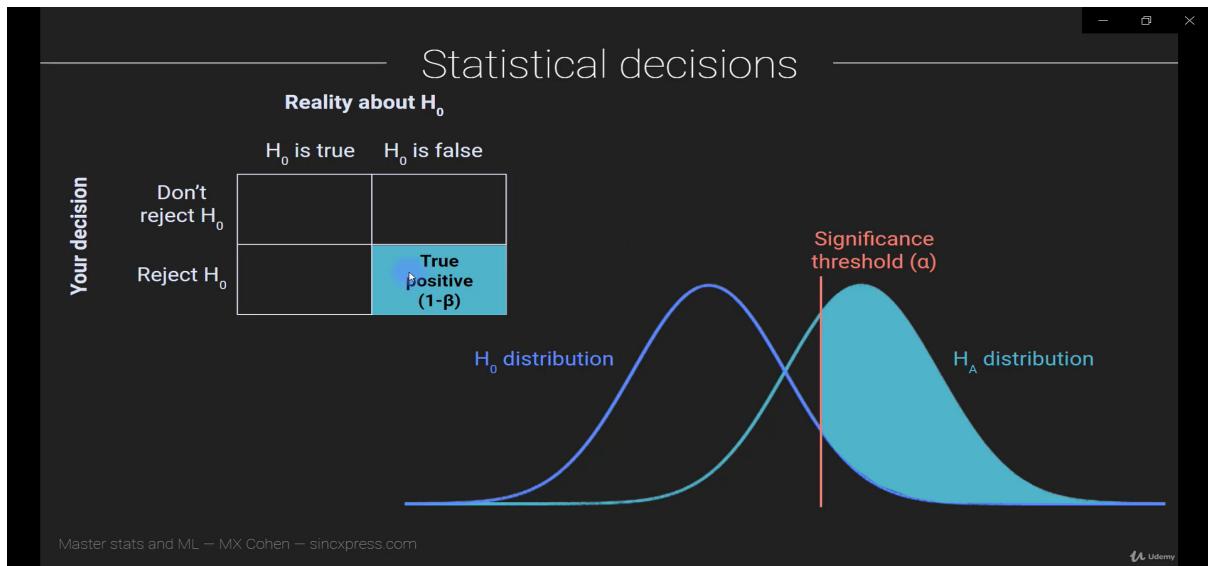
Degrees of freedom definition

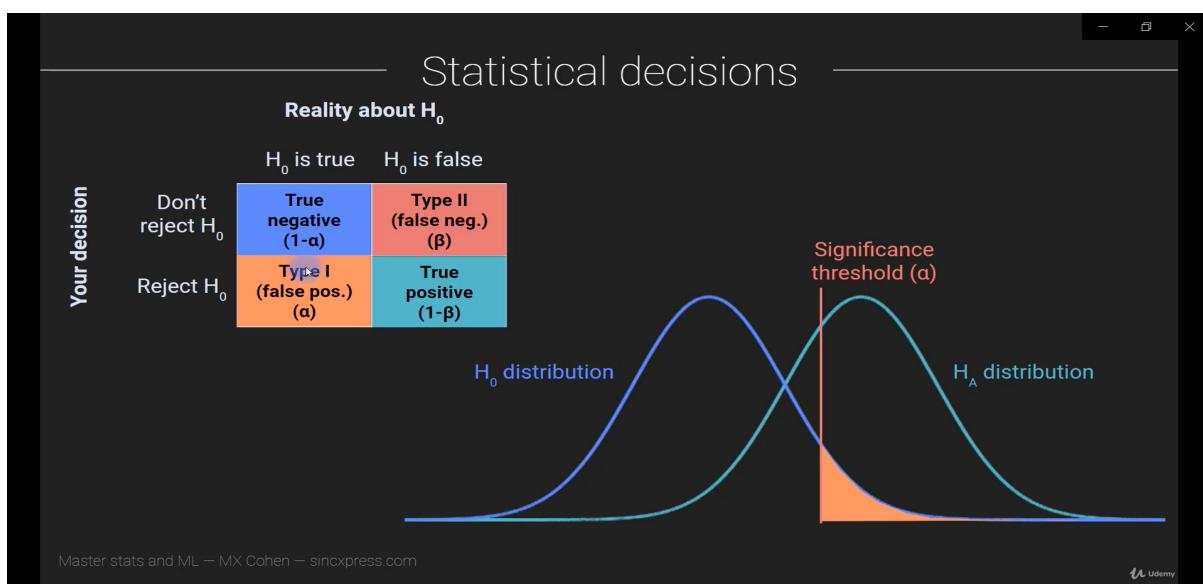
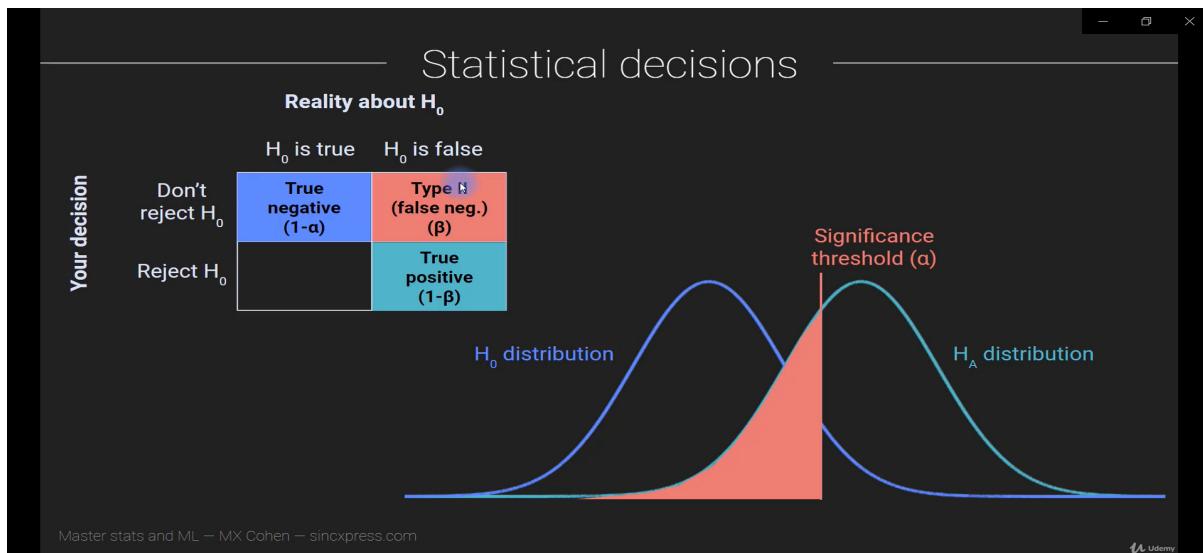
Degrees of freedom:

- The number of independent sample values.
- The number of sample data points that can vary.
- The number of data values that are unconstrained by the rest of the data values.

Generally: $df = N - k$ (N data points, k parameters)

Type I and Type II Errors





Statistical Trade-offs

Statistical trade-offs

		Reality about H_0	
		H_0 is true	H_0 is false
Your decision	Don't reject H_0	True negative ($1-\alpha$)	Type II (false neg.) (β)
	Reject H_0	Type I (false pos.) (α)	True positive ($1-\beta$)

Reducing Type I errors will also reduce True Positives

Significance threshold (α)

H_0 distribution

H_A distribution

Master stats and ML – MX Cohen – sincxpress.com

Udemy

Statistical trade-offs

		Reality about H_0	
		H_0 is true	H_0 is false
Your decision	Don't reject H_0	True negative ($1-\alpha$)	Type II (false neg.) (β)
	Reject H_0	Type I (false pos.) (α)	True positive ($1-\beta$)

Increasing True Positives will also increase Type I errors

Significance threshold (α)

H_0 distribution

H_A distribution

Master stats and ML – MX Cohen – sincxpress.com

Udemy

Statistical errors in the real world (sortof)

The image contains two side-by-side photographs. The left photograph shows a doctor in a white coat and stethoscope around his neck, looking at a patient who is wearing a blue hospital gown. A yellow speech bubble from the doctor says 'You're pregnant'. The right photograph shows a doctor in a white coat examining a pregnant woman's belly. A yellow speech bubble from the doctor says 'You're not pregnant'.

Type I error
(false positive)

Type II error
(false negative)

http://www.sethspielman.org/courses/geog5023/r_examples/Type_I_and_II_Errors.html

Master stats and ML — MX Cohen — sincxpress.com

Udemy

How to reduce Statistical Errors

How to minimize statistical errors

The image shows three vertically stacked bell-shaped curves on a dark background. In the top row, the two curves are wider and closer together. In the middle row, the two curves are narrower and further apart. In the bottom row, the two curves are very narrow and far apart, with a cursor pointing to the gap between them.

To minimize statistical errors, you need to:

Increase the distance between the distributions (bigger effects).

Decrease the width of the distributions (less variable data).

Master stats and ML — MX Cohen — sincxpress.com

Udemy

Parametric and Non-Parametric Statistics

Nonparametric statistics: definitions

What “non-parametric” doesn’t mean:
No parameters at all.

Correct meaning(s):
Statistics that are not based on assumptions about underlying distributions (typically, Gaussian).

Statistical inference methods that generate the H_0 distribution from the data, not from an equation.

Master stats and ML — MX Cohen — sincxpress.com

Udemy

Statistical model complements

Parametric test	Nonparametric test
1-sample t-test	Wilcoxon sign-rank test
2-sample t-test	Mann-Whitney U test
Pearson correlation	Spearman correlation
ANOVA	Kruskal-Wallis test

Important application of nonparametric statistics:
Permutation testing and cross-validation

Master stats and ML — MX Cohen — sincxpress.com

Udemy

Advantages and limitations

Parametric statistics: <ul style="list-style-type: none">● Standard, widely used● Based on assumptions● Assumptions should be tested, (though rarely done)● Can be incorrect when assumptions are violated● Computationally fast● Analytically proven	Nonparametric statistics: <ul style="list-style-type: none">● Some are nonstandard● “No” assumptions necessary● Can be slow/intensive● Some are sensible algorithms rather than proven methods● Appropriate for non-numeric data● Appropriate for small sample sizes● Some methods give different results each time
-------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------	------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------

Conclusion:
Use parametric methods when possible.
Use nonparametric methods when necessary.

Master stats and ML — MX Cohen — sincxpress.com

Udemy

Multiple Comparisons and Bonferroni Correction

Multiple comparisons and Bonferroni correction

In this video, you will learn

- What “multiple comparisons” means and why it causes problems.
- How Bonferroni correction works.
- The difference between individual Type 1 error rate and familywise error rate.

Master stats and ML – MX Cohen – sincxpress.com

Udemy

Say we have three groups (blue, green, yellow) and we want to compare the means of these three groups.

In how many ways can we compare their means ?

→ 3 ($H_1 = \text{(blue,green)}$, $H_2 = \text{(green,yellow)}$, $H_3 = \text{(yellow,blue)}$)

(assuming that comparison is asymmetric, i.e comparison between blue and green is same as green and blue)

Mo' groups, mo' comparisons



$H_1 = \text{blue vs. teal}$

$H_2 = \text{blue vs. yellow}$

$H_3 = \text{teal vs. yellow}$

Master stats and ML – MX Cohen – sincxpress.com

- So there are three hypotheses which we want to perform with the threshold of 0.05.
- The problem is that probabilities are additive. The first comparison is independent of the second and so on.
- Therefore, even though for individual hypothesis we have 0.05 threshold (e.g $P(H_1 | H_0) = 0.05$), the total test will have the threshold of 0.15

— Mo' comparisons, mo' problems —

Probabilities are additive.

$$p(H_1 | H_0) = .05$$

$$p(H_2 | H_0) = .05$$

$$p(H_3 | H_0) = .05$$

$$p(H_1 | H_0) + p(H_2 | H_0) + p(H_3 | H_0) = .05 + .05 + .05 = .15$$

Problem: The probability of a false alarm (Type I error) is 15%! Unacceptably high!

Master stats and ML — MX Cohen — sincxpress.com

Udemy

ENG 6:50 PM

- Because they are all based on the same dataset, when we put all the tests in the set, the overall threshold becomes 0.15.
- Imagine, when we have many many comparisons to do, the overall p-value will be really high.

— Mo' comparisons, mo' problems —

Probabilities are additive.

$$p(H_1 | H_0) = .05$$

$$p(H_2 | H_0) = .05$$

$$p(H_3 | H_0) = .05$$

$$p(H_1 | H_0) + p(H_2 | H_0) + p(H_3 | H_0) = .05 + .05 + .05 = .15$$

Note: This is called the “Familywise error rate” (abbreviated FWE or FWER).

Master stats and ML — MX Cohen — sincxpress.com

Udemy

So by the Bonferroni correction, the threshold for the individual test will become alpha / n

The Italians come to the rescue!



Bonferroni correction:

Threshold = α/N

α = e.g., 0.05
 N = number of tests

Master stats and ML — MX Cohen — sincxpress.com

Udemy

Mo' comparisons, mo' problems

Probabilities are additive.

$$p(H_1 | H_0) = .05/3$$
$$p(H_2 | H_0) = .05/3$$
$$p(H_3 | H_0) = .05/3$$
$$p(H_1 | H_0) + p(H_2 | H_0) + p(H_3 | H_0) = .05/3 + .05/3 + .05/3 = .15/3 = .05$$

Problem solved! The probability of a false alarm (Type I error) is 5%!

Master stats and ML — MX Cohen — sincxpress.com

Udemy

Now individual test should have more evidence in-order to reject the H_0 in favour of H_1

Mo' comparisons, mo' problems

Probabilities are additive.

$$p(H_1 | H_0) = .05/3$$
$$p(H_2 | H_0) = .05/3$$
$$p(H_3 | H_0) = .05/3$$
$$p(H_1 | H_0) + p(H_2 | H_0) + p(H_3 | H_0) = .05/3 + .05/3 + .05/3 = .15/3 = .05$$

Note the difference between the individual α ($p=.0167$) and the familywise α ($p=.05$).

Master stats and ML — MX Cohen — sincxpress.com

Udemy

- Also note that Bonferroni correction is great and works well (when you have multiple tests and they are independent of each other).
- But there are many other correction methods which are more suitable for other specific situation like FDR (False discovery rate), cluster based correction (if you have spatially or temporally correlated data),...

Type of Significances

The significances

Statistical significance:
The probability of observing a test statistic this large given that the null hypothesis is true.

Theoretical significance:
A finding is relevant for a theory or leads to new experiments.
This has nothing to do with statistical significance.

Clinical (practical, societal, educational) significance:
A finding is relevant for diagnosing or treating a disease.

Master stats and ML — MX Cohen — sinoxpress.com

Udemy

Examples

Hypothesis and result:
MMR vaccine causes autism. $p = .79$.

No statistical significance.
Strong clinical and societal significance.

Master stats and ML — MX Cohen — sinoxpress.com

Udemy

Examples

Hypothesis and result:
Piles of sand with larger grains collapse sooner. $p = .001$.

Strong statistical significance.
Strong theoretical significance (scale-free dynamics)

Master stats and ML — MX Cohen — sincxpress.com



Cross Validation

Nice Explanation:  Machine Learning Fundamentals: Cross Validation

Uses of cross-validation

Variance of an estimate (e.g., confidence intervals)
Also called “jack-knifing.” Compare the variance of the parameter estimates over different training sets.

Avoid bias in analysis results
Applying the model to data it has never seen. Overfitting the training data will decrease performance on the test data.

Compute classification accuracy
Dominant application in machine learning and deep learning.

Master stats and ML — MX Cohen — sincxpress.com



Uses of cross-validation

Ideally, the test set is truly independent of the training set.
This avoids bias and overfitting in the test set.

Be mindful of whether this is the case in your data.
For example, a group of self-selected students from the same school.

8. The t-test Family

Table of Content

One Sample T-test:
How to get Statistically Significant Results
Two Sample t-test
Wilcoxon Signed Rank Test (also called Signed Rank Test) (Non parametric t-test):
Mann-Whitney U Test (Also called Mann-Whitney-Wilcoxon t-test OR Wilcoxon Rank-Sum test) (Non parametric t-test):
Permutation Testing for t-test:

One Sample T-test:

- One sample T-Test tests if the given sample could have been generated from a population with a specified mean.
- If it is found from the test that the means are statistically different, we infer that the sample is unlikely to have come from the population.
- **Example:** If you want to test a car manufacturer's claim that their cars give a highway mileage of 20 kmpl on an average. You sample 10 cars from the dealership, measure their mileage and use the one sample T-test to determine if the manufacturer's claim is true.
- **Purpose of One Sample T Test**

The purpose of the One Sample T Test is to determine if a sample observation could have come from a process that follows a specific parameter (like the mean). It is typically implemented on small samples.

- **How did we determine One sample T-test is the right test for this?**

Because, there is only one sample involved and you want to compare the mean of this sample against a particular (hypothesised) value.

- Null hypothesis will be that the given sample has been drawn from the specified population parameter, i.e there is no difference between sample statistic and population parameter.
- Depending on the how the problem is stated, the alternate hypothesis can be one of the following 3 cases:
- **Case 1:** $H_1 : \bar{x} \neq \mu$. Used when the true sample mean is not equal to the comparison mean. Use Two Tailed T-Tests.
- **Case 2:** $H_1 : \bar{x} > \mu$. Used when the true sample mean is greater than the comparison mean. Use the Upper Tailed T-Tests.
- **Case 3:** $H_1 : \bar{x} < \mu$. Used when the true sample mean is lesser than the comparison mean. Use Lower Tailed T Test.

The BASIC formula for the t-test is which gets modified based on different variations of t-test.

The screenshot shows a presentation slide with a dark background. The title 'T-test: the formula' is centered at the top. Below the title, the formula for a two-sample t-test is displayed in a large, light-colored font:

$$t_k = \frac{\bar{x} - \bar{y}}{s/\sqrt{n}} = \frac{\text{Difference of means}}{\text{Standard deviations}}$$

At the bottom left, the source is cited as 'Master stats and ML — MX Cohen — sincxpress.com'. At the bottom right, the Udemy logo is visible.

The assumptions of one sample t-test

One-sample t-test: assumptions

1. Data are numeric (not categorical), ideally interval or ratio (discrete is probably OK).
2. Data are independent from each other.
3. Data are randomly drawn from the population to which generalization should be made.
4. Mean and standard deviation are valid descriptors of central tendency and dispersion (i.e., data are approximately normally distributed).

Master stats and ML — MX Cohen — sincxpress.com

Udemy

One-sample t-test: formula

$$t_{n-1} = \frac{\bar{x} - \mu}{s/\sqrt{n}}$$

\bar{x} Sample mean
 μ H_0 value
 s Sample standard deviation
 n Number of data points

$n-1$ Degrees of freedom

Master stats and ML — MX Cohen — sincxpress.com

Udemy

Code

Since in one sample t-test we have only one sample, so substitute $n_1 = n_2$ and $s_1 = s_2$ and check if the formula matches with one sample t-test.

Formula for unequal N, unequal variances

$$t_{df} = \frac{\bar{x}_1 - \bar{x}_2}{\sqrt{\frac{(n_1-1)s_1^2 + (n_2-1)s_2^2}{n_1+n_2-2}} \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}}$$
$$df = n_1 + n_2 - 2$$

Master stats and ML — MX Cohen — sincxpress.com

How to get Statistically Significant Results

Biggie T

$$t_k = \frac{\bar{x} - \bar{y}}{s / \sqrt{n}} = \frac{(\bar{x} - \bar{y})\sqrt{n}}{s}$$

Increase the group differences
Reduce variances
Increase sample size

Master stats and ML — MX Cohen — sincxpress.com

Two Sample t-test

- It is a simple extension of the one-sample t-test.
- The numerator is pretty similar with the one sample t-test and the denominator depends on the characteristics of the data groups.
- Purpose: Test whether the two groups could have been drawn from the same distribution.
- i.e It is applied to compare whether the average (or any parameter) difference between two groups is really significant or if it is due instead to random chance.
- e.g It helps to answer questions like whether the average success rate is higher after implementing a new sales tool than before or whether the test results of patients who received a drug are better than test results of those who received a placebo.

i.e you collect data of group1 and group2 and compare them OR you take data before some effect (like we had drug) and then you again collect data after that effect. Now you want to compare them.

- There are several formulas for two sample t-tests. The numerator is always the same.
- *The denominator depends on whether the groups are paired or unpaired, groups have (ROUGHLY) same or different variances, groups have (ROUGHLY) same or different sample sizes.*
(unequal sample sizes apply only to the unpaired groups.)

Paired and unpaired groups:

- **Paired groups:** We compare two groups data collected from the same individual .(like data of health before taking drug and data of health after drug on same individual)
- **Unpaired groups:** We compare the two groups data collected from different individuals. data of effect of drug on men's and data of effect of drug on females.

Formula for unequal N, unequal variances

$$t_{df} = \frac{\bar{x}_1 - \bar{x}_2}{\sqrt{\frac{(n_1-1)s_1^2 + (n_2-1)s_2^2}{n_1+n_2-2}} \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}}$$

$$df = n_1 + n_2 - 2$$

Master stats and ML – MX Cohen – sinxpress.com

- $df = n_1 + n_2 - 2$
Because for the first group, $df = n_1 - 1$ and for second group, $df = n_2 - 1$, therefore overall $df = n_1 + n_2 - 2$
- Assumptions are the same as one Sample T-test.

Code

Wilcoxon Signed Rank Test (also called Signed Rank Test) (Non parametric t-test):

- It is a non-parametric alternative to the one sample and two sample PAIRED ONLY t-test.
- Mainly used when data is not close to gaussian AND For ordered (ranked) categorical variables without a numerical scale.
- Tests for the differences in the medians (instead of differences in the mean)

Algorithm to compute the Wilcoxon test

Step 1: Remove equal pairs

Step 2: Rank-transform diffs

Step 3: Sum ranks where $x > y$

Step 4: Convert to z

$$Z = \frac{W - n(n + 1)/4}{\sqrt{\frac{n(n+1)(2n+1)}{24}}}$$

n is the number of remaining pairs.

Z is normally distributed under H_0 and can be converted to a p-value.

Master stats and ML – MX Cohen – sincxpress.com

Udemy

Algorithm:

- **Remove the equal pairs**
 1. If we are doing two sample paired test, then remove pairs which are equal to H_0 value (null hypothesis value)
 2. If we are doing one sample test, then remove data points which are equal to H_0 value (null hypothesis value)
 3. Reasoning: Equal pairs do not contribute to the test either way.
- **Rank Transform diffs**
 1. $r = \text{rank}(|x - y|)$
 2. If it is two sample test, then x and y are pairs
 3. If we are doing one sample test, then y is replaced by H_0 (null hypothesis value)
- **Sum the rank where $x > y$,**
 $\text{sumR} = 0$
 for $i=1$ to n :
 $\text{sumR} += r_i$, if $x_i > y_i$

$$W = \text{sumR}$$

- **Convert to Z**

(Z is normally distributed under H_0 and can be used to find out p-values)

Null hypothesis: Both the data comes from the distribution (The null hypothesis for this test is that the medians of two samples are equal.)

Alternative hypothesis:

two tailed ($m1 != m2$)

one-tailed ($m1 < m2$ OR $m1 > m2$)

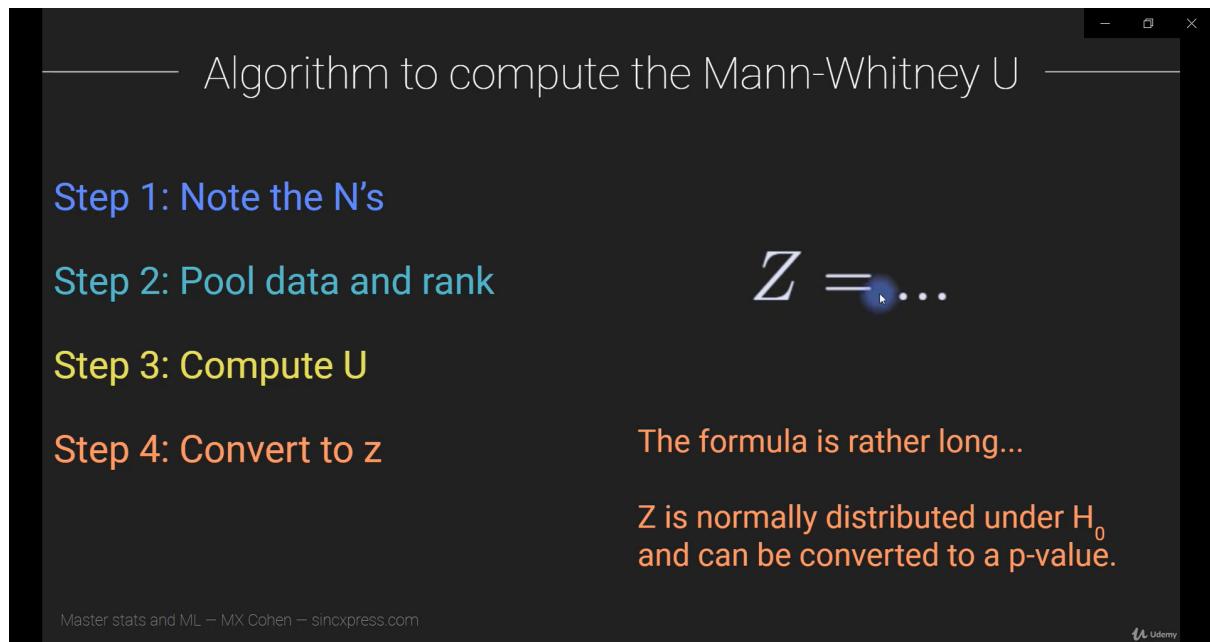
For paired non-parametric t-test example:

<https://www.statisticshowto.com/wilcoxon-signed-rank-test/>

Code

Mann-Whitney U Test (Also called Mann-Whitney-Wilcoxon t-test OR Wilcoxon Rank-Sum test) (Non parametric t-test):

- It is a non-parametric alternative to the two sample UNPAIRED t-test. (i.e two independent groups)
- Used when data is not normal and the dependent variable is either ordinal or continuous, but not normally distributed.
- Tests for the differences in the medians (instead of differences in the mean).
- It can have different sample sizes.



Algorithm:

- Do the assignments
 xf = dataset with fewer data points
 xm = dataset with more data points
 nf = number of data points in xf
 nm = number of data points in xm .
- Concatenate the xf and xm (note: xm should get concatenated after xf). After concatenation, compute rank
 $r = \text{rank}(xf.append(xm))$
- Compute U (sum ranks upto nf)
 $U \leftarrow \text{sum}(r[:nf])$
- Compute Z (Z is normally distributed under H_0 and can be used to find out p-values),

Null hypothesis: Both the data comes from the same distribution (The null hypothesis for this test is that the medians of two samples are equal.)

Alternative hypothesis:

two tailed ($m1 \neq m2$)

one-tailed ($m1 < m2$ OR $m1 > m2$)

[Code](#)

Permutation Testing for t-test:

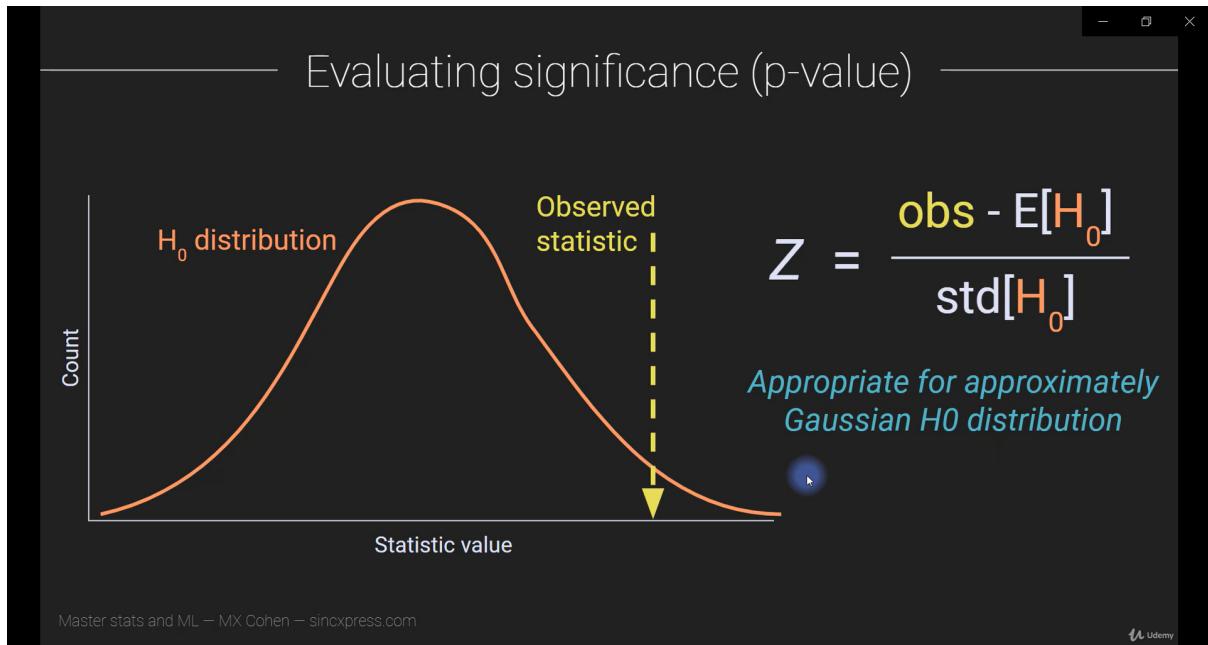
- The key difference between parametric and permutation based t-test significance is, in general parametric approach is, you evaluate the position of your test statistic relative to the analytic distribution (*based on your assumptions and this is not something that you compute by running some analysis i.e you are not generating the null hypothesis distribution on your own*) under the null hypothesis. That is standard parametric approach (like t-tests, correlation, etc)
- In contrast, in permutation testing, I don't want to make any assumptions like where the null hypothesis distribution comes from OR what it looks like. Instead, I want to compute the **null hypothesis distribution based on the data I have** and then interpret the value of the observed test statistic relative to the empirical distribution.

Q. How do we create empirical hypothesis distribution ?

- Say we had 7 data points, 4 from group1 and 3 from group2.
group1 = {1,2,3,4}
group2 = {5,6,7}
- **Step1** → Mix the data together so that there are no labels on them
pooledData = {1,2,3,4,5,6,7}
- **Step 2** → Randomly assign labels (group1 OR group2) to these data points, but the ratio of those labels should be the same (i.e there should be 4 points for group1 and 3 points for group2). Note: We are not modifying any value of the data points, just randomly relabeling them
group1 = {1,4,5,6}
group2 = {2,3,7}
- **Step 3** → calculate t value from these groups (t-test for two unpaired samples) and store it.
- Now repeat the steps 2 and 3, say N number of times.
- **Step 4** → You will get N t-values. You plot those t-values and this is your null hypothesis distribution.
- **Step 5** → Now, you go back to your original data and perform t-test between your two groups and evaluate your test statistic with respect to the null hypothesis distribution you calculated.

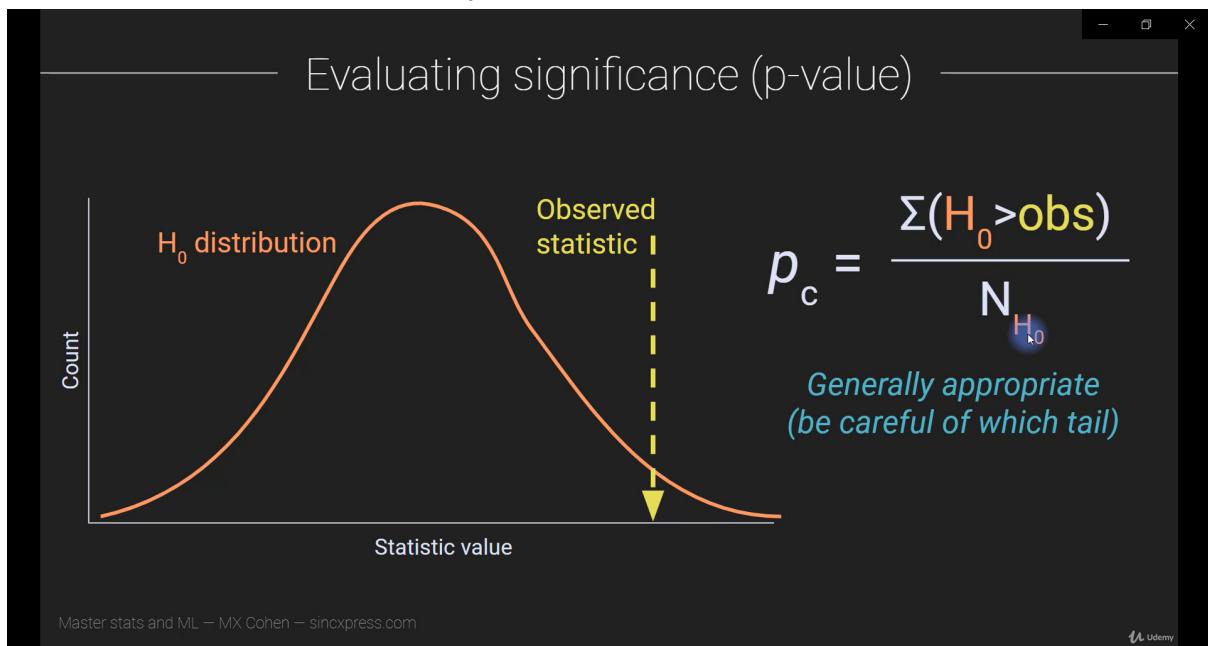
Now for evaluating the p-values, there are two ways,

This way is appropriate if your distribution is approximately gaussian.



Second Way:

p_c = probability counts, basically you take all the possible H_0 values which are greater than observed test statistic and divided by total possible values of H_0



Code

9. Confidence Intervals on Parameters

Table of Content

Confidence Intervals and why do we need them:

Factors influencing the confidence intervals:

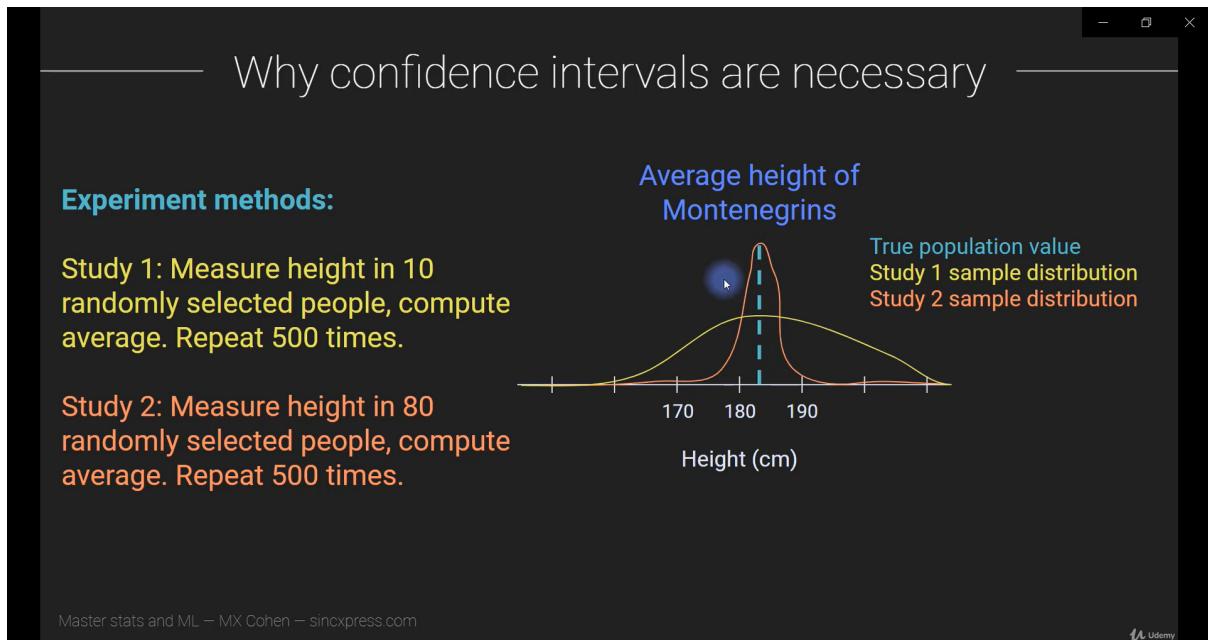
There are two methods to compute the confidence intervals:

Confidence Intervals using formula:

Confidence Intervals using Bootstrapping:

Confidence Intervals and why do we need them:

- say the population mean of heights of Montenegrins is 182.
- We perform 2 experiments now:
 1. Measure the height of 10 randomly selected Montenegrins and store the average. Repeat this experiment 500 times.
 2. Measure the height of 80 randomly selected Montenegrins and store the average. Repeat this experiment 500 times.
- plot the distribution of means for both the experiments.



→ Which distribution from the above two experiments will be more confident?

obviously (qualitatively visually) distribution 2, the distribution of means from experiment 2 because it will be narrower as compared to distribution from experiment 1 because we measured more people in experiment 2.

(see their mean of means is same but the spread is different)

This is the base idea of confidence intervals, in experiment 2, we measured more people data and got a narrow distribution and in experiment 1, we got a wider distribution (that's why we are less confident).

Naive Definition: The probability of unknown population parameter falls within the range of values after repeating the experiment multiple times.

NOTE:

1. *The confidence interval does not guarantee that the true mean (population parameter) will lie between sample estimate and confidence intervals.*
2. *In fact, it just tells our confidence on our estimation of our sample estimate, if we repeat this same experiment many times in future, about 95% of times, the true mean will lie between sample estimate and confidence intervals.*

Mathematically Confidence Interval:

$P(L < \mu < U) = c$

L = Lower bound

U = Upper Bound

c = Confidence

μ = population parameter

P = probability

Typical confidence intervals probabilities: 95%, 99%, 90%.

Factors influencing the confidence intervals:

- 1) **Sample size:** as sample size increases, the confidence intervals decreases
- 2) **Variance:** as variance decreases, the confidence interval also decreases.

There are two methods to compute the confidence intervals:

- 1) By applying formula
- 2) By using an empirical bootstrapping or resampling method.

Confidence Intervals using formula:

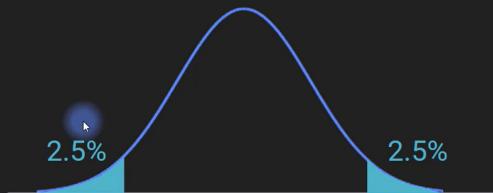
- Since, typically we are not able to repeat the experiments multiple times due to some reasons and we just have one sample. In that case, the formula method is used typically.
- Formula:
$$\text{Confidence Interval (CI)} = \bar{x} \pm t^*_{k-1} * \left(\frac{s}{\sqrt{n}} \right)$$

\bar{x} = sample mean
 s = sample std
 n = sample size
 t^*_{k-1} = t-value with $k-1$ degrees of freedom (we will talk more about this)

This t^* is not the same kind of value that we have been learning in the previous section, it does come from t-distribution, but this is not the test statistic.

The formula for C.I.

$$\text{C.I.} = \bar{x} \pm t^*(k) \frac{s}{\sqrt{n}}$$



Master stats and ML – MX Cohen – sinocpress.com

Udemy

t^* = t-value associated with the one tail of the confidence interval = $(1-c) / 2$

$t^* = t_{\text{inverse}}((1-c) / 2, n-1)$

t_{inverse} is basically some hypothetical function which takes the area under the curve (probability) and the degrees of freedom ($n-1$) and returns the corresponding t-value in the t-distribution.

e.g $c=95\%$, $n=20$

$t^* = t_{\text{inverse}}((1-0.95) / 2, 20-1) = 2.093$

The reason we are dividing by 2 is because we want the confidence intervals to be two tails. (both the sides).

Assumptions of this Formula:

1) s is an appropriate measure of dispersion in data. **That's why it is parametric way of computing confidence intervals.**

[Code](#)

Confidence Intervals using Bootstrapping:

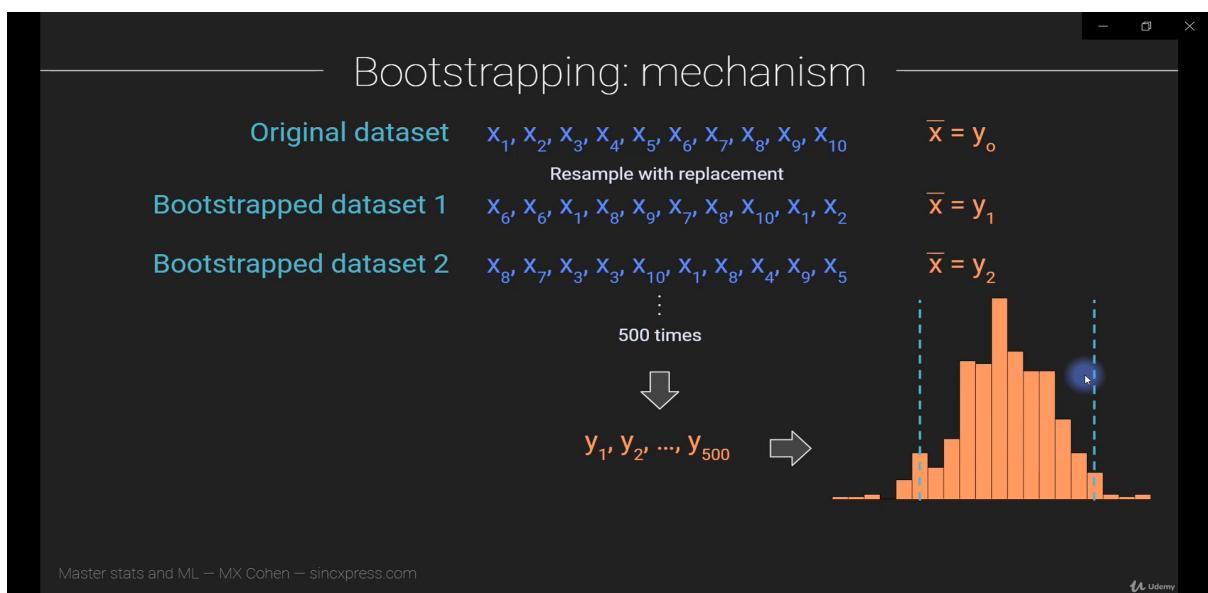
- Non parametric method for computing confidence intervals.
- Instead of using the formula for computing confidence intervals, we can compute them directly by using data.
- This is done by repeatedly randomly sampling the dataset.
- Thus, pretend that your sample is the population and your resampling is the sample.
- Mechanism:
say you had 10 points with mean (y_0)

You randomly sample with replacement data of size 10 (the size can be smaller than actual n OR equal to n) and compute the mean. say y_1 .

You repeat this experiment 500 times and you will get 500 means ($y_1, y_2, y_3, \dots, y_{500}$)

You plot those means and compare the mean of this distribution with y_0 .

- Now you can also find out 2.5 percentile and 97.5 percentile of distribution, which will give you the empirical estimation of confidence interval.



Pros and Cons of Bootstrapping:

Pros:

- 1) Works for any kind of parameter (mean, variance, correlation, median, etc.)
- 2) Do not assume any normality.
- 3) Useful for limited data.

Cons:

- 1) Gives (slightly) different results each time.
- 2) Can be time consuming for large datasets (even though you can use bootstrap sample size < N(size of dataset))
- 3) Sample must be a good representative of the population.

Code

10. Correlation

Table of Content

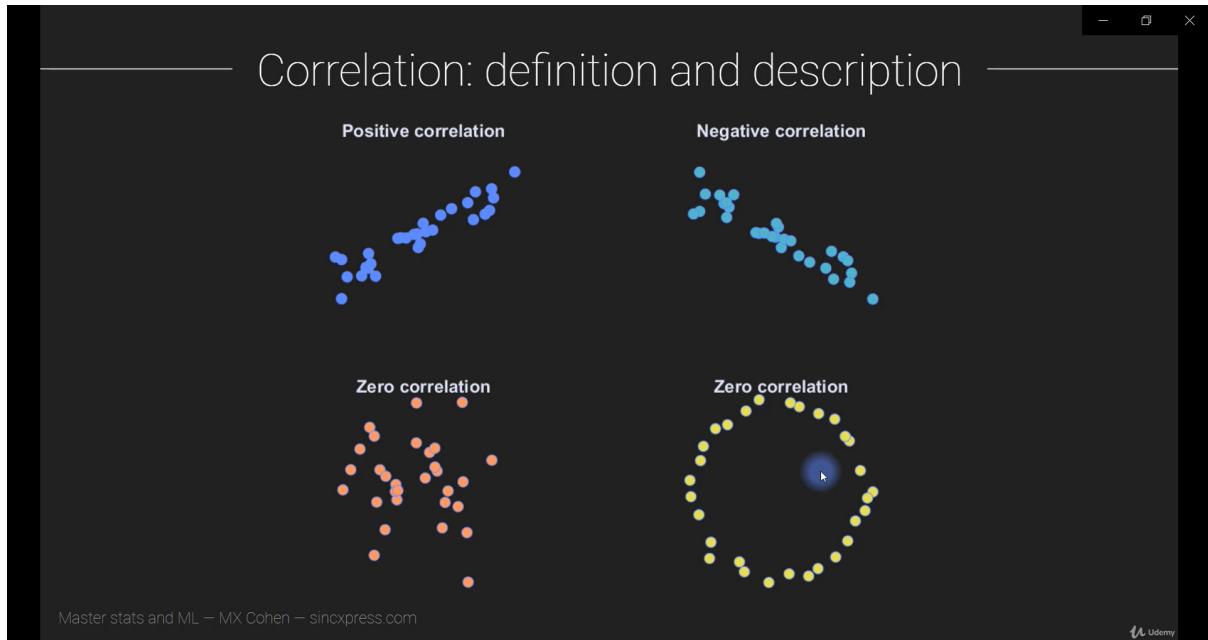
Motivation and Description of correlation.
What is Causation?
Correlation vs Causation:
Covariance and Correlation:
Correlation Matrix:
Partial Correlation:
Limitations of Pearson's Correlation Coefficient
Non-Parametric Correlation: Spearman Rank
Fisher Z-Transformation for Correlation:
Kendall's Correlation for Ordinal Data:
Cosine Similarity:
Point Biserial Correlation

Correlation:

Motivation and Description of correlation.

- In this, we basically check how the two groups are related to each other.
- Correlation analysis computes the correlation coefficient (r)
- It is standardised between -1 and 1. -1 indicating perfect inverse relationship, 0 indicating no relationship and 1 indicating perfect relationship.
- The r is a single number that shows the LINEAR RELATIONSHIP between two groups.
- WE SHOULD NOT JUST COMPUTE THE CORRELATION COEFFICIENT ONLY. WE SHOULD ALWAYS COMPUTE p-value TO INTERPRET ITS STATISTICAL SIGNIFICANCE.

Examples:



Why does the fourth figure have zero correlation even though the figure seems to have a relation between x and y ?

→ There is zero correlation because *correlation measures the linear relationship*. There is a non linear relationship.

- Null hypothesis: there is zero correlation.
Alternative Hypothesis: There is non-zero correlation.

What is Causation?

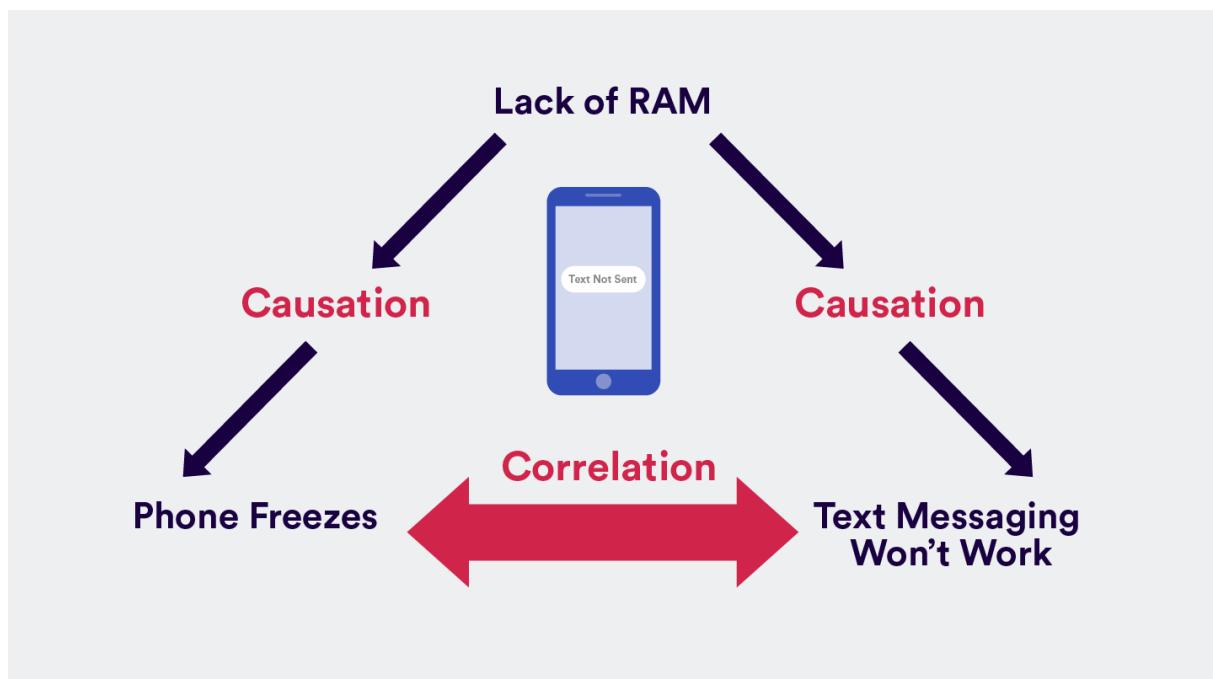
- Causation is implying that A and B have a cause-and-effect relationship with one another. You're saying A causes B. Causation is also known as causality.
- Firstly, causation means that two events appear at the same time or one after the other.
- And secondly, it means these two variables not only appear together, the existence of one causes the other to manifest.

Correlation vs Causation:

- Correlation just measures the relationships.
- It does not reveal or imply causation.
- Causality can be demonstrated by experimental manipulations.
- Example

Ref: <https://clevertap.com/blog/correlation-vs-causation/>

My mother-in-law recently complained to me: "Whenever I try to text message, my phone freezes." A quick look at her smartphone confirmed my suspicion: she had five game apps open at the same time plus Facebook and YouTube. The act of trying to send a text message wasn't causing the freeze, the lack of RAM was. But she immediately connected it with the last action she was doing before the freeze. She was implying a causation where there was only a correlation.



Covariance and Correlation:

- Covariance is a single number that measures the linear relationship between the two variables.
- Correlation is just normalised covariance. (from -1 to 1)
- Division by n-1 is not as critical, it depends on what you are doing with covariance.

— Formula for covariance and correlation —

$$c = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})$$
$$r = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2 \sum_{i=1}^n (y_i - \bar{y})^2}}$$

Master stats and ML – MX Cohen – sincxpress.com

Udemy

- But remember, we should not blindly trust the correlation value, we should check the p-value, is it really statistically significant ?

— P-value of correlation coefficient —

$$t_{n-2} = \frac{r \sqrt{n-2}}{1 - r^2}$$

Statistical significance is computed from a t-value that is based on the strength of the correlation and the number of data points.

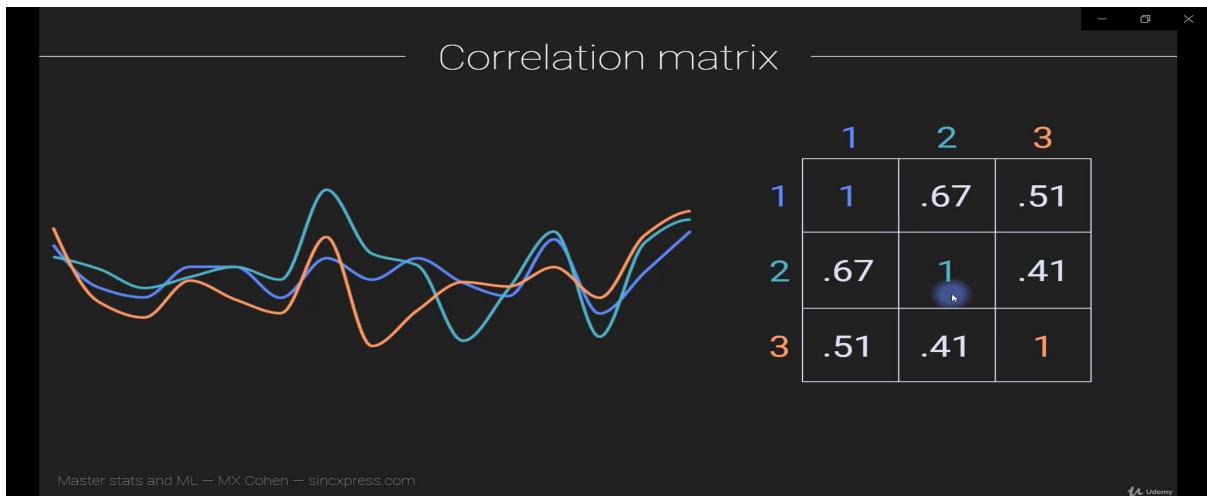
Master stats and ML – MX Cohen – sincxpress.com

Udemy

- Say, we are increasing the sample size, but assume that r will always be constant, then see the t-value will also increase i.e being more significant.
- Just an observation, when the mean of both the variables is zero and variance is 1, then the covariance formula is the same as correlation formula.

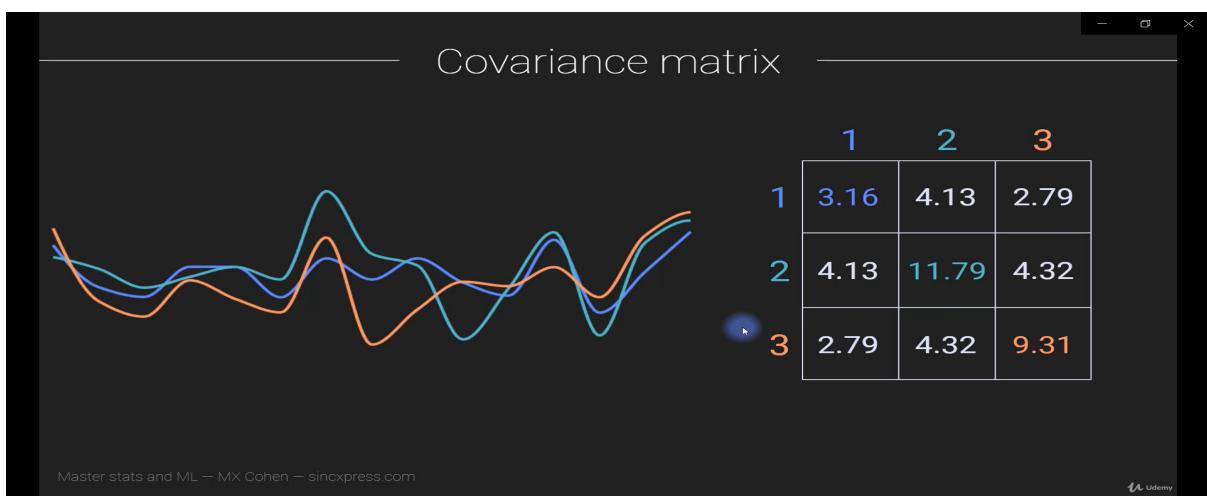
Code

Correlation Matrix:



Master stats and ML — MX Cohen — sincxpress.com

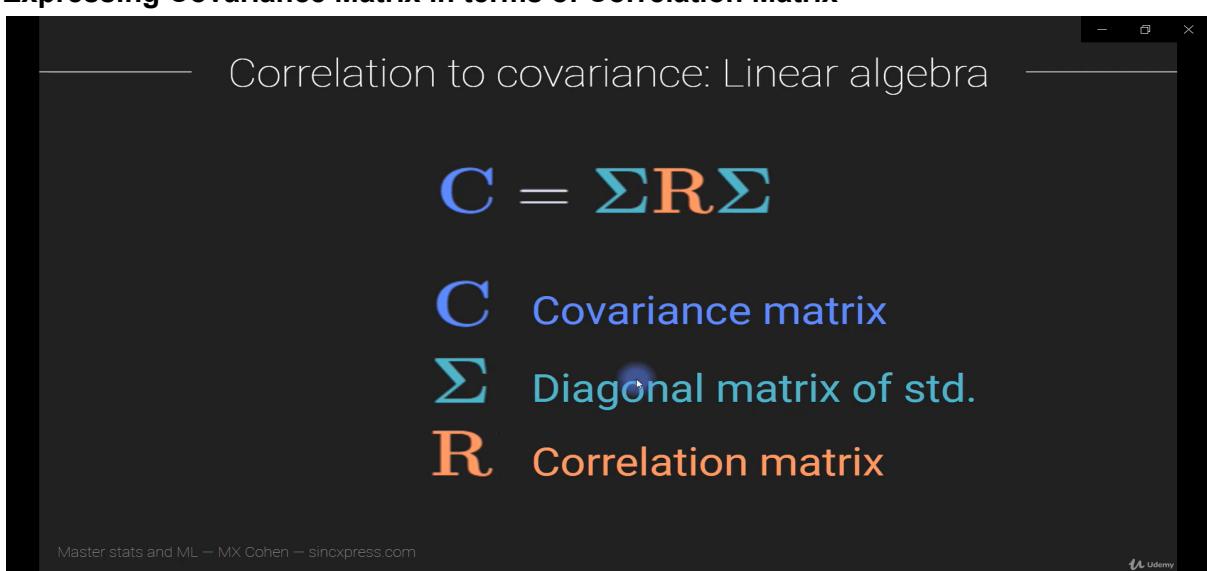
Udemy



Master stats and ML — MX Cohen — sincxpress.com

Udemy

Expressing Covariance Matrix in terms of Correlation Matrix



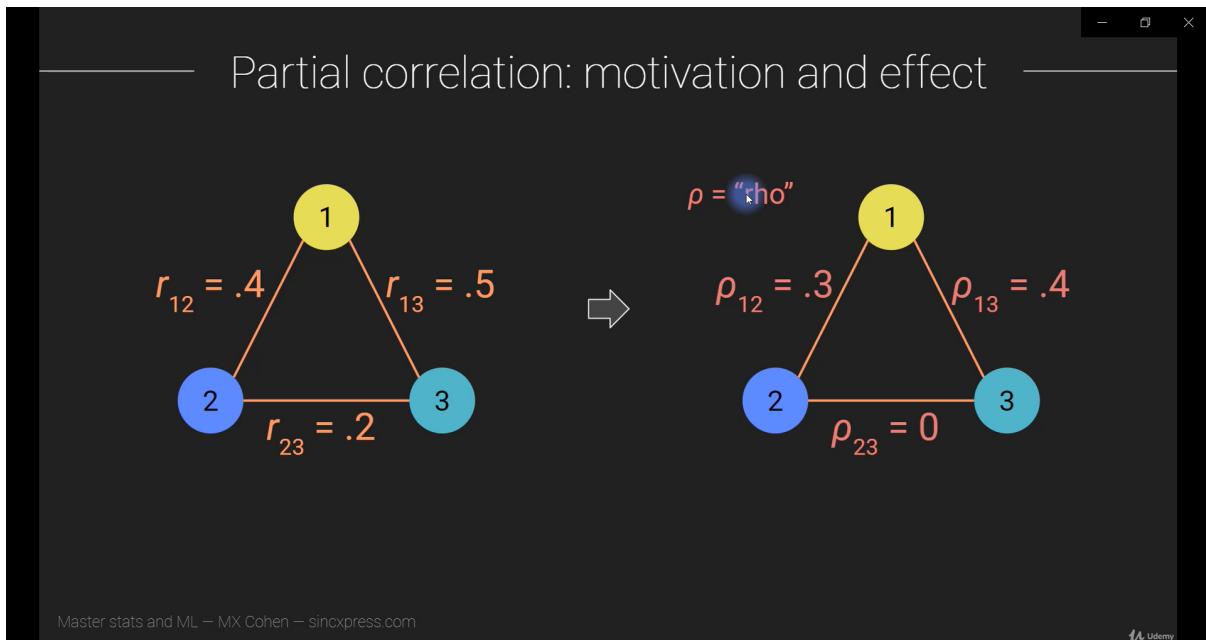
Master stats and ML — MX Cohen — sincxpress.com

Udemy

[Code](#)

Partial Correlation:

- Motivation: Say you had three variables, weather, ice cream consumption and shark attacks.
- From the data, it has been found that, as the ice cream consumption increases the shark attacks also increases. i.e there is some correlation between them. This sounds weird right.
- From the data (hypothetical) :
- The correlation between:
 1. weather and ice cream consumption = 0.4
 2. ice cream consumption and shark attacks = 0.2
 3. weather and shark attacks = 0.5



(check the above image, where node 1 = weather, 2 = ice cream consumption, 3 = shark attacks)

- *But actually, what is happening is, as the weather is hot, people consume more ice cream and go to the beach and as people go to the beach, more shark attacks happen.*
- *Therefore, the weather caused the increase in ice cream consumption which led to an increase in shark attacks.*

Now the **idea behind partial correlation** is, what is the correlation between variable A and variable B if I keep aside the effect of variable C.

Therefore, partial correlation between:

- 1) weather and ice cream consumption = 0.3
- 2) ice cream consumption and shark attacks = 0
- 3) weather and shark attacks = 0.4

Partial correlation: formula

$$\rho_{xy|z} = \frac{r_{xy} - r_{xz}r_{yz}}{\sqrt{1 - r_{xz}^2}\sqrt{1 - r_{yz}^2}}$$

Master stats and ML – MX Cohen – sincxpress.com

Udemy

- Intuitively, the formula for partial correlation makes sense because,
- say we wanted $p(xy|z)$ i.e we are interested in the correlation between x and y by removing the shared variance of z,
- In the numerator of formula, it is $r_{xy} - r_{xz}r_{yz}$
 r_{xy} = correlation between x and y, in which we are interested.
 r_{xz} = correlation between x and z.
 r_{yz} = correlation between y and z.
i.e removing(subtracting) the shared variance of z on x and y.
- And denominator is just some normalisation factor.

In the formula screenshot, we just partialled out the single variable, but it is also possible to waive out the many variables as well.

say we wanted, $p(xy | z,w)$

i.e we want correlation between x and y, by removing the shared variance of z and w, then

$$p(xy | z,w) = r_{xy} - r_{xz}r_{yz} - r_{xw}r_{yw}$$

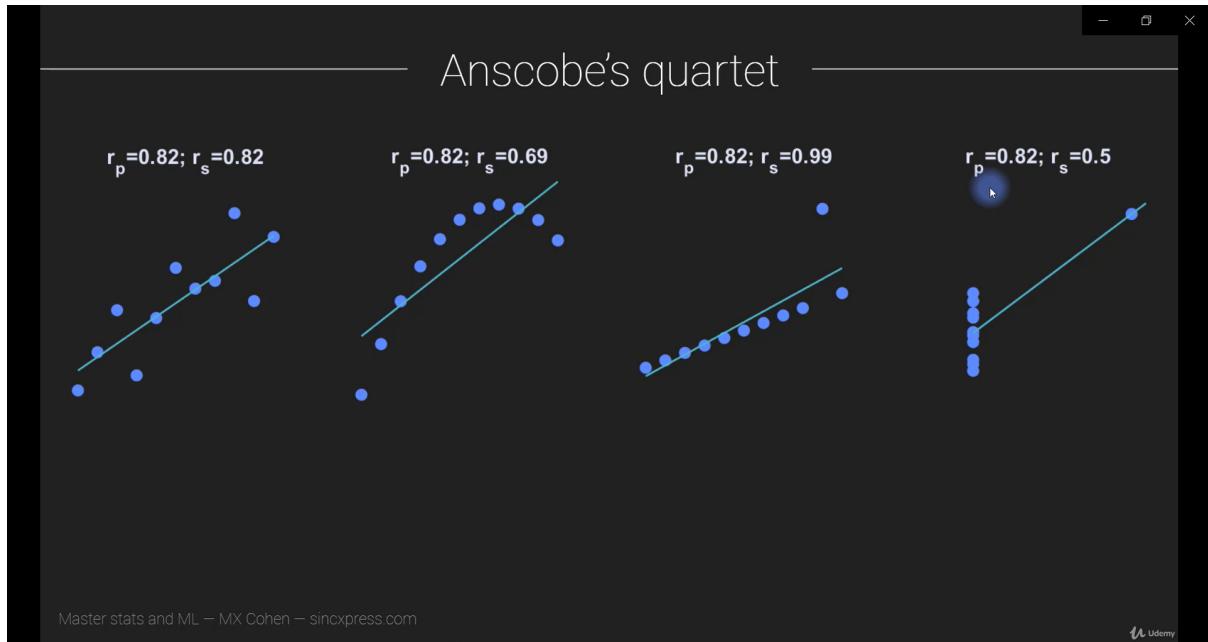
$\sqrt{1 - r_{xz}^{**2}} \cdot \sqrt{1 - r_{yz}^{**2}} \cdot \sqrt{1 - r_{xw}^{**2}} \cdot \sqrt{1 - r_{yw}^{**2}}$
and you can extend this to many variables.

[Code](#)

Limitations of Pearson's Correlation Coefficient

r_p = pearson correlation, r_s = spearman correlation

Notice: The pearson correlation is same in all the four diagrams but the nature of data is different in all 4 cases.



- In the first diagram, the data seems to have a linear relationship.
- In the second diagram, the data has non-linear relation, but Pearson correlation is trying to capture its linear component.
- In the third diagram, the data seems to have a perfect linear relationship but with one outlier, hence pearson correlation gets strongly affected by it.
- In the fourth diagram, the data is almost vertical with one outlier due to which the $r_p = 0.82$. Also spearman correlation also gets affected by that outlier but has much lower value than r_p .

Conclusion:

- The Pearson correlation can over-represent or under-represent relationships if the data is non-linear or contains outliers.
- Therefore, the Pearson coefficient is appropriate if x and y are normally distributed and their relationship is linear without outliers.

Note:

The other correlation measures which we will talk about later, also finds linear relationships in the data just as they are differently sensitive to outliers.

Way to generate data with the correlation coefficient we want: [Code](#)

Non-Parametric Correlation: Spearman Rank

- Pearson and spearman correlation converge when the x and y are normally distributed.
- The main benefit of spearman correlation is that they are robust to outliers.
- The spearman correlation tests for the MONOTONIC RELATIONSHIP, regardless of whether the relationship is linear or nonlinear.
- **Algorithm:**
 - 1) Transform both x and y into ranks.
e.g reminder of rank: (874387438, -40, 1, 0) --> (4,1,3,2)
 - 2) Compute the Pearson correlation on ranks.
 - 3) p-value calculated the same as in Pearson correlation.
- See how spearman correlation deals with the outliers, instead of considering the distance of data points with mean, it simply transforms the variable into ranks which totally reduces the effect of outliers.
- see the above example, how the outlier 874387438 is converted into rank and this will not affect the correlation process.

Code

Fisher Z-Transformation for Correlation:

- Let's say you have two populations x and y. You sample from x (say s_x) and sample from y (say s_y), you compute the correlation and store it.
- You repeat the above step N(large) times, you will get a correlation list and when you plot it, say it will be a uniform distribution bounded between -1 and 1.
- Now, many analysis methods use normality as the assumption and say you want to do some analysis on that correlation list and that analysis requires that correlation should be normally distributed, so you need to do some transformation on it to make it normal.
- Fisher-z transformation is used to transform that correlation list into (roughly) normally distributed values.

— How does the Fisher-Z transform work? —

$$z_r = \frac{1}{2} \ln \left(\frac{1+r}{1-r} \right) = \operatorname{arctanh}(r)$$

Master stats and ML – MX Cohen – sincxpress.com

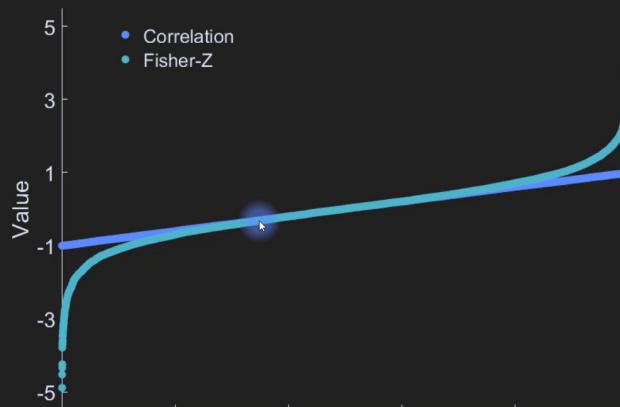
Udemy

arctanh means inverse hyperbolic tangent of correlation coefficient.

See the correlations are bounded between -1 and 1.

But the transformed z-values differs, specially at the tails, but in the middle the transformed values are quite similar to actual correlation coefficient,

— The effect of the Fisher-z transform —



Master stats and ML – MX Cohen – sincxpress.com

Udemy

- If you want to know the relationship between two samples, then you don't need to worry about fisher-z transform ,but if you have many many correlation values and you want to do some analysis, then you might use this.

Kendall's Correlation for Ordinal Data:

- Used for the ordinal data (i.e the order in the data is important and some numeric values can be assigned)
e.g size: small, medium, high -> 0,1,2 (1.5 doesn't make sense in this case).
- How it works:

Kendall's Rank Order Correlation | Kendall's Tau - τ | Numerical Example | Non-Parametric ..   

KENDALL'S RANK ORDER CORRELATION (KENDALL'S TAU) - τ

- Kendall's Tau (τ) is a **non-parametric** measure of relationships between columns of ranked data.
- The Tau correlation coefficient returns a value of **0 to 1**, where:
 - 0** is no relationship,
 - 1** is a perfect relationship

$T = 2S / (N(N - 1))$

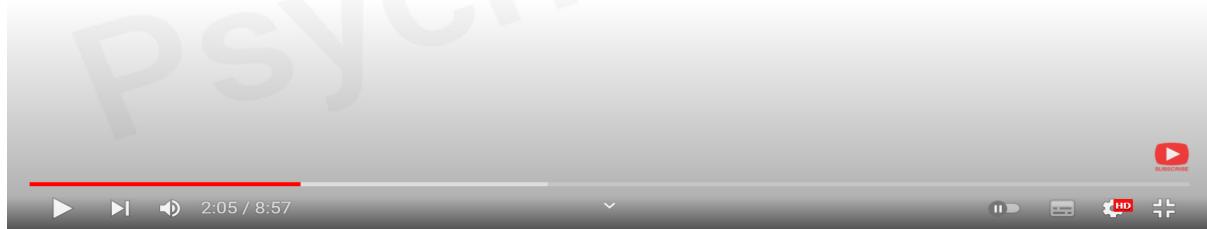
Where:

S = (score of agreement – score of disagreement on X and Y)

N = Number of objects or individuals ranked on both X and Y



Kendall's Rank Order Correlation | Kendall's Tau - τ | Numerical Example | Non Parametric ... i t →
Suppose we ask X and Y to rate their preference for four objects and give points out of 10. Now to see whether their preferences are related to each other we may use the following steps:



Kendall's Rank Order Correlation | Kendall's Tau - τ | Numerical Example | Non Parametric ... i t →
10. Now to see whether their preferences are related to each other we may use the following steps:

Data:

	A	B	C	D
X	6	8	5	2
Y	8	4	9	6



	A	B	C	D
X	6	8	5	2
Y	8	4	9	6

Step 1: Ranking the data of X and Y

	A	B	C	D
X	3	4	2	1
Y	3	1	4	2

Step 2: Rearrange the data of X in order of 1 to N (4 in this case)

X	D	C	A	B
	4	2	3	4

Step 2: Rearrange the data of X in order of 1 to N (4 in this case)

X	D	C	A	B
	1	2	3	4

Step 3: Put the corresponding score of Y in order of X and Determine number of agreements and disagreements

	D	C	A	B
X	1	2	3	4
Y	2	4	3	1

To calculate S we need number of agreements and disagreements. This can be calculated by

Using the Y scores, starting from left and counting the number of ranks to its right that are larger, these are agreements in order. We subtract from this the number of ranks to its right that are smaller- these are the disagreements in order. If we do this for all the ranks and then sum the results we obtain S:

Using the Y scores, starting from left and counting the number of ranks to its right that are larger, these are agreements in order. We subtract from this the number of ranks to its right that are smaller- these are the disagreements in order. If we do this for all the ranks and then sum the results we obtain S:

Y	2	4	3	1	Total
	2	+	+	-	+1
		4	-	-	-2
			3	-	-1
				1	0
				Grand Total= S	- 2

Step 4: Calculate T

$$T = 2S / (N(N - 1))$$

$$T = 2(-2) / (4(4 - 1))$$

$$T = -4 / 12$$

$$T = -0.33$$

Thus $T = -0.33$ is a measure of the agreement between the preferences of X and Y.

- Kendall tau-b has an adjustment for the ties and it is most often used for ordinal data.
- kendall-tau-a do not have adjustment for ties.
- Interpretation is identical to pearson and spearman correlation (i.e bounded by -1 and 1)
- **The assumptions for Kendall's Tau include:**
 1. Continuous or ordinal
 2. Monotonicity: Your two variables should have a monotonic relationship. This means that the direction of the relationship between the variables is consistent. For instance, when one variable goes up, the other goes up (in general). The relationship would also be monotonic if when one variable goes up, the other goes down (in general).

Code

What is the difference between Spearman's Rho and Kendall's Tau?

→ Spearman's Rho and Kendall's Tau are very similar tests and are used in similar scenarios.

We recommend using Kendall's Tau first and Spearman's Rho as a backup.

Cosine Similarity:

- 1) Very very similar to Pearson correlation.

The slide has a dark background with white text. At the top, it says "Cosine vs. Pearson". Below that are two mathematical formulas. The first formula is for Pearson correlation: $r = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2} \sqrt{\sum_{i=1}^n (y_i - \bar{y})^2}}$. The second formula is for cosine similarity: $\cos(\theta) = \frac{\sum_{i=1}^n x_i y_i}{\sqrt{\sum_{i=1}^n x_i^2} \sqrt{\sum_{i=1}^n y_i^2}}$. At the bottom left, it says "Master stats and ML — MX Cohen — sincxpress.com". At the bottom right, there is a small Udemy logo.

- where x is a point and y is point in n-dimensional space.(and there is line between origin and x, origin and y)
- The Pearson and cosine similarity are identical when in cosine similarity, the data is already mean-centred.
- The interpretation of cosine similarity and Pearson correlation is the same.

Code

Point Biserial Correlation

- Point-biserial correlation is used to understand the strength of the relationship between two variables. Your variables of interest should include one continuous and one binary variable.
- Assumptions for Point-Biserial correlation
 1. Continuous and Binary
 2. Normally Distributed: Only use Point-Biserial Correlation on your data if the continuous variable is normally distributed.
 3. No Outliers: The continuous variables that must not contain outliers.
 4. Equal Variances (variance in continuous variable for category1 should be similar to variance in continuous variable for category 2)

11. ANOVA

Table of Content

ANOVA:
Some Important Terminologies in ANOVA:
Assumptions of ANOVA:
When ANOVA is not appropriate:
Now mathematics
Important: Correct interpretation of p-value:
Interpretation of Main Effects and Interactions:
One-Way ANOVA Example:
Two-Way ANOVA Example:

ANOVA:

- It stands for ANalysis Of VAriance.
- The goal of an ANOVA is to determine the effects of discrete independent variables (categorical) on a continuous dependent variable.
- Setting up steps for ANOVA:
 - Step 1) Review the experiment design and make sure ANOVA is an appropriate method.
 - Step 2) Identify the independent and dependent variables.
 - Step 3) Create a table of factors and levels (when possible) (more on this later)
 - Step 4) Compute the model and interpret the results.

Let's start with each steps:

Step 1)

Research Goal: Test whether the new covid-19 medications are effective for different groups of people.

Experiment: Randomly assign patients to receive medication A, B and placebo. Measure the disease severity after 10 days. Separate older people(>50 years) from younger people(<= 50)

Step 2:

Dependent Variable: The variable for which you are trying to explain variance.

Independent variables: The variables which you HOPE will explain the Dependent variable.

Therefore, for the above the experiment:

IVs: Medication and age group (Both are categorical)

DV: Severity after 10 days.

Step 3:

Create a table of factors and levels

Factors: The dimensions in the IVs

Levels: The specific groups or manipulations within each dimension (basically possible values in each dimension)

In the above experiment,

Factors: medication, age

Levels: A, B, placebo (for medication)
older, younger (for age)

Now create table like this

	Medication A	Medication B	Placebo
Younger people			
Older People			

We have written WHEN POSSIBLE why?

- Say if we had a third categorical variable as well, then we had to create a 3d table and not 2d table like above, and if we have n variables, then we cannot visualise this table of factors and levels.
- We can compute ANOVA values with n variables, and it is okay if we cannot create tables.

Step 4 (Very Important)

We will get into the maths later, but for now there are 3 main interpretations for ANOVA:

- **Main effect:**
How the one factor influences the DV even when ignoring all the other factors.
e.g The younger people's symptoms improve faster as compared to the older people, regardless of medication type.
(i.e main effect of age on severity independent of medication type)
- **Interactions:**
The effect of one factor on DV depends on the levels of other factors.
e.g The medication A works better for older people and medication B works better for younger people.
(i.e effect of medication type depends on levels of age)
- **Intercept:**
The average value of DV is different from zero. The intercept of ANOVA is usually ignored.
e.g the symptoms improve for almost everyone after 10 days.

We are mainly concerned with main effects and interactions between the factors.

Some Important Terminologies in ANOVA:

- $<x>$ -way ANOVA, where x = number of factors
- e.g one-way ANOVA: i.e number of factors considered = 1
(it can have any number of levels.)
e.g Determine the influence of day-of-week on the purchase of iphones.
There is only one factor = day-of-week and number of levels = 7
- Two-way ANOVA: number of factors considered = 2
e.g Determine the influence of day-of-week and gender on iphone purchases.

Repeated-measures ANOVA:

- If atleast one factor involves multiple measurements from the same individual.
- e.g

Research question: Determine the effect of snack-type on mood.

Experiment: Volunteers eat chocolate for 2 days, potato chips for 2 days and icecream for 2 days. (randomised order)

There is only one factor called snack-type and dependent variable mood.

This is repeated-measures ANOVA, because the SAME INDIVIDUAL will eat chocolate for 2 days, potato chips for 2 days and icecream for 2 days and their mood will be noted after every 2 days.

Balanced vs unbalanced ANOVA:

Balanced vs. unbalanced ANOVA

Balanced: The same number of data points in each cell.
Unbalanced: Different number of data points across cells.
This could happen because of data collection or cleaning.

Balanced ANOVA

		Medication		
		A	B	placebo
Age group	Younger	20	20	20
	Older	20	20	20

Master stats and ML – MX Cohen – sincxpress.com

Balanced vs. unbalanced ANOVA

Balanced: The same number of data points in each cell.
Unbalanced: Different number of data points across cells.
This could happen because of data collection or cleaning.

Unbalanced ANOVA

		Medication		
		A	B	placebo
Age group	Younger	20	23	21
	Older	18	20	20

Master stats and ML – MX Cohen – sincxpress.com

ANOVA vs MANOVA:

- anova: only one dependent variable (you can have many independent variables)
- manova: multivariate dependent variable (as many independent variables as appropriate)
e.g effect of medication type and age on disease severity and medical expenses.

Fixed and Random Effects ANOVA:

- **Fixed effects:** The number of levels of a factor is always fixed. (e.g operating systems: They are always three: mac, windows, linux)
- **Random effects:** The number of levels of a factor are not fixed. (e.g say age is factor, which is not fixed, so we can discretize into specific number of levels.)
- **Mixed effects:** Some of the factors are fixed effects and some factors are random effects.

Assumptions of ANOVA:

- Independence: The data points should be sampled independently from the population.
- Normality: The residuals (unexplained variance after fitting the model) are normally distributed.
- Homogeneity of variance: The variance in the factor-levels table should be roughly the same, otherwise there might be some bias.

Note:

If these assumptions are strongly NOT met, then you can use Non-parametric ANOVA alternatives:

Kruskal-Wallis test (KW-ANOVA): Alternative to the one-way ANOVA on the rank-transformed data.

Note:

The anova are generally robust to violations of the assumptions. non-parametric anovas are rarely used and may cause more of a headache than they are worth.

When ANOVA is not appropriate:

—— Examples of when ANOVAs are inappropriate ——

Research goal: Test whether people with more Facebook friends have higher self-reported extraversion.

Variables: DV: number of Facebook friends. IV: scores on a personality questionnaire.

ANOVA? No categorical factors.

T-test? No group(s) to compare.

Correlation? Yes, because we are looking for a linear relationship between two continuous variables.

Master stats and ML — MX Cohen — sincxpress.com

Udemy

—— Examples of when ANOVAs are inappropriate ——

Research goal: Test whether RSI (repetitive stress injury) is decreased for a group of meditators vs. non-meditators.

Variables: IV: group (meditate or not). DV: scores on an RSI index.

ANOVA? Only one factor with two levels.

T-test? Yes, because there are two groups to compare.

Correlation? No, because the IVs are discrete groups, not continuous variables.

Master stats and ML — MX Cohen — sincxpress.com

Udemy

Now mathematics

(Just building some blocks)

- Remember we had a table of factors vs levels, let n_i be the number of data points in each cell and μ_i be the mean of each cell.
- Null hypothesis: all the cells have the same mean i.e $\mu_1 = \mu_2 = \dots = \mu_k$ i.e means all the groups are the same.
- Alternative hypothesis: the mean of at least two cells differ from each other. This is a very generic alternative hypothesis. So to get exactly which two cells have different means, we have to do some more investigation. i.e mean of at least one group differs from at least one other group.
- Anova is based on the sum of squares i.e $\sum (x_i - \bar{x})^2$
This formula is very similar to variance.
- But in anova, why don't we have a denominator of $n-1$?
→ Because in anova we calculate multiple sum of squares, we are not actually interested in sum of squares values, we calculate their ratios, therefore when we do ratios, the denominators cancel out.
- The total variation in the dataset is expressed as the sum of variations across different groups and the variation within each group.

We are basically going to calculate the f-statistic and it is ratio of two things

$f = \text{"Explained" variance (due to factors or levels in each factor) (sum of squares between groups)}$

"Unexplained" variance" (due to natural variation) (sum of squares within groups)

Omnibus f-test means: your p-value associated with the f-statistic is significant but you don't know exactly the means of which two groups differ from each other, so that f-value is called omnibus. So you have to do post-hoc analysis i.e comparison of multiple groups.

Now check this article

<https://medium.com/data-science-in-your-pocket/anova-mathematics-explained-c85a36172ac1>

This is called the ANOVA table (k = number of levels in a factor)

Source of variance	Sums of squares	Degrees of freedom	Mean square	F	P-value
Between groups	SS_B	k-1	MS_B	MS_B/MS_W	p
Within groups	SS_W	N-k	MS_W		
Total	SS_T	N-1			

Remember:

$$H_0 : \mu_1 = \mu_2 = \dots = \mu_k$$

$$H_A : \mu_i \neq \mu_j$$

Master stats and ML — MX Cohen — sincxpress.com

Udemy

Important: Correct interpretation of p-value:

Source of variance	Sums of squares	Degrees of freedom	Mean square	F	P-value
Between groups	SS_B	k-1	MS_B	MS_B/MS_W	p
Within groups	SS_W	N-k	MS_W		
Total	SS_T	N-1			

Correct interpretation of p<.05: At least one level (group) is statistically significantly different from at least one other level. Determining which groups differ requires data visualization follow-up t-tests.

Master stats and ML — MX Cohen — sincxpress.com

Udemy

Example for correct interpretation of p-value:

	Sums of squares	Degrees of freedom	Mean square	F	P-value
Between groups	1850.47	3 (k-1)	616.82	13.25	.0002
Within groups	697.95	15 (N-k)	46.53		
Total	2548.42	18 (N-1)			

One-way ANOVA with four levels and 19 data points.

Conclusion: The mean of at least one level is statistically significantly different from the mean of at least one other level.

Master stats and ML — MX Cohen — sincxpress.com

Udemy

Example ANOVA table

	Sums of squares	Degrees of freedom	Mean square	F	P-value
Between groups	1850.47	3 (k-1)	616.82	13.25	.0002
Within groups	697.95	15 (N-k)			
Total	2548.42	18 (N-1)			

“Omnibus” F-test

One-way ANOVA with four levels and 19 data points.

Conclusion: The mean of at least one level is statistically significantly different from the mean of at least one other level.

Master stats and ML – MX Cohen – sincxpress.com

Udemy

omnibus f-test means: your p-value associated with the f-statistic is significant but you don't know exactly the means of which two groups differ from each other, so that f-value is called omnibus. So you have to do post-hoc analysis i.e comparison of multiple groups

Say for the independent category there are 4 levels, then we can have 6 possible comparisons with 4 levels.

Example graphed results

Problem: Which conditions are different from which?

A bar chart with four blue bars. The y-axis is labeled "Dependent variable" and the x-axis is labeled "Levels of the factor". The bars are of equal height, suggesting no significant difference between the levels.

Dependent variable

Levels of the factor

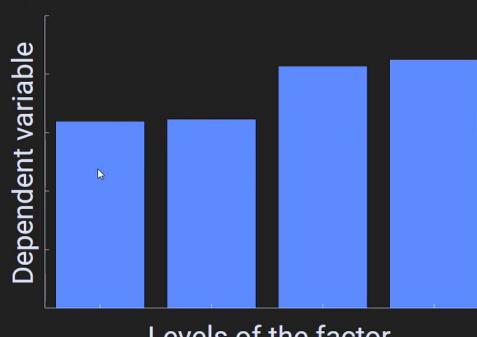
Master stats and ML – MX Cohen – sincxpress.com

Udemy

All 6 possible comparisons with $p \leq 0.05$ will lead to such high combined type-1 error rate.

— Mo' comparisons, mo' problems —

Problem: All possible comparisons (each at $p < .05$) leads to $6 * .05 = .3$ Type I error rate.



Master stats and ML — MX Cohen — sincxpress.com

Udemy

— Thanks, Tukey, for the test —

Solution: The Tukey test allows for post-hoc comparisons while controlling the familywise error rate.

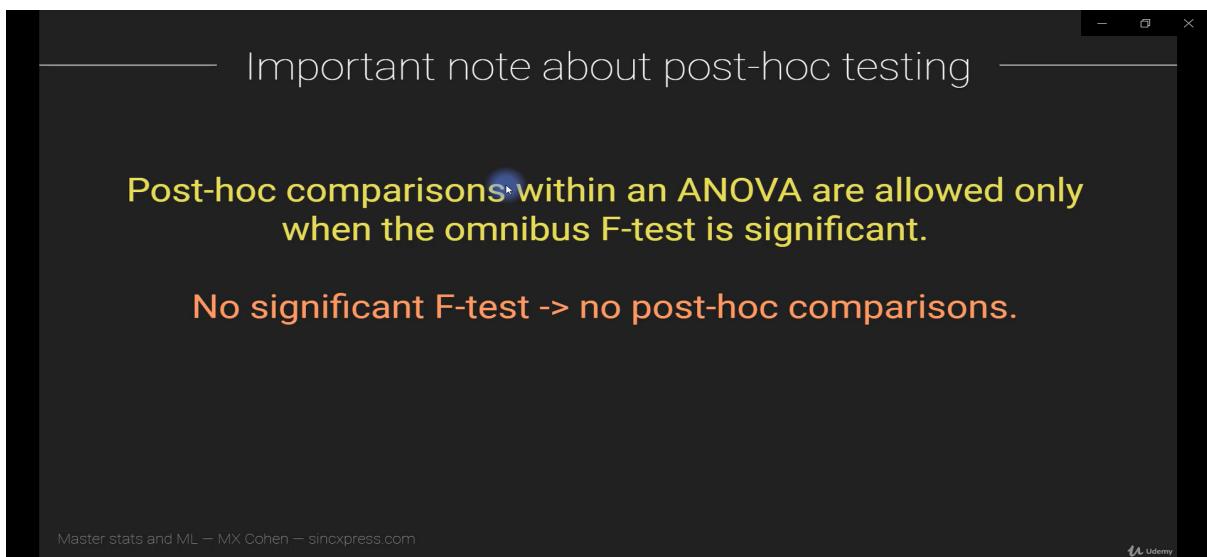
$$q = \frac{\bar{x}_b - \bar{x}_s}{\sqrt{MS_{\text{Within}}} \sqrt{2/n}}$$

q is evaluated with $(j, n-j)$ degrees of freedom.

j is the number of comparisons.
 n is the total number of data values.

Master stats and ML — MX Cohen — sincxpress.com

Udemy



Interpretation of Main Effects and Interactions:

— How to interpret the main effects and interactions —

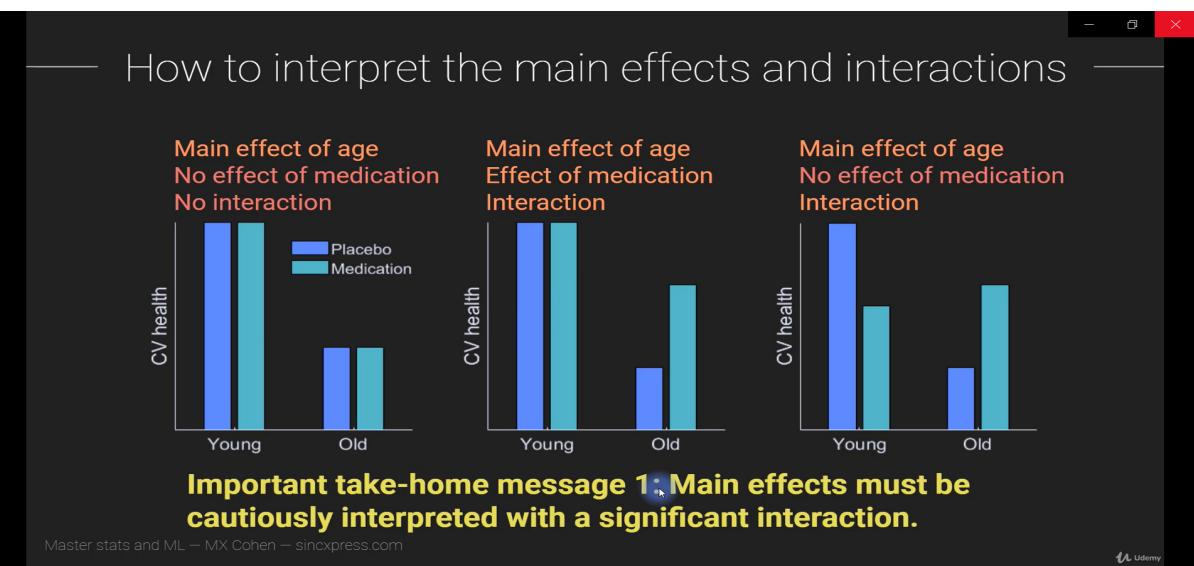
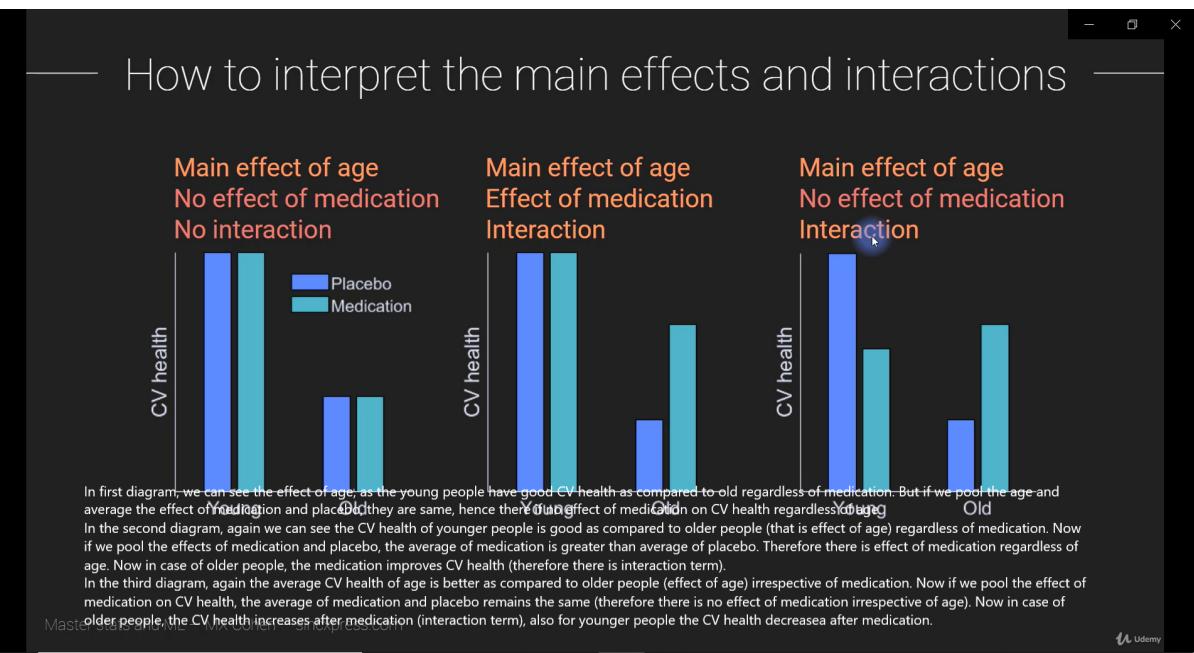
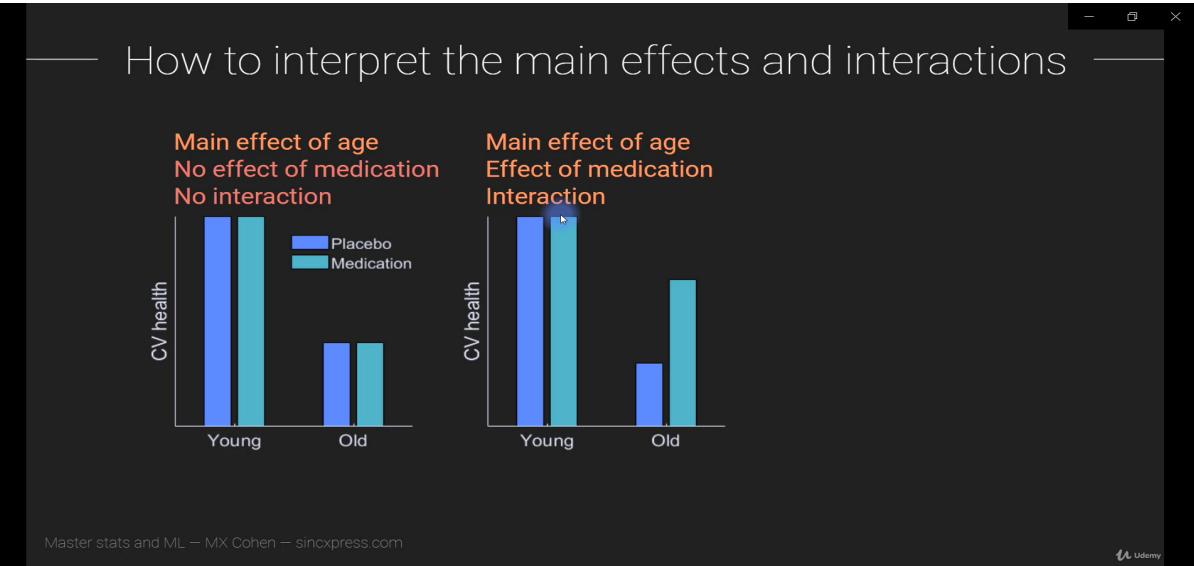
Example: The effects of vitamin supplement (vs. placebo) and age (30-40 vs. 50-60) on cardiovascular health.

Experiment design table:

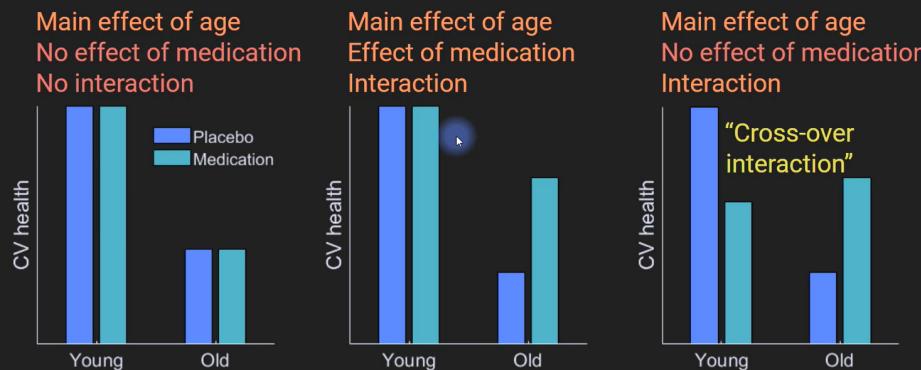
		Supplement	
		real	placebo
Age group	30-40	X	X
	50-60	X	X

Master stats and ML - MX Cohen - sinxpress.com

Udemy



— How to interpret the main effects and interactions —



Important take-home message 2: Data must always be visualized for proper interpretation!

Master stats and ML — MX Cohen — sincxpress.com

Udemy

One-Way ANOVA Example:

(This is unbalanced-ANOVA)

The experiment

Research goal: Test self-reported happiness after watching different genres of movies.

Variables: IV: movie genre (horror, romcom, documentary, sci-fi).
DV: happiness rating (1-100).

		Movie genre			
		Horror	Romcom	Docu	Sci-fi
N per group		5	5	4	5

Master stats and ML — MX Cohen — sincxpress.com

Udemy

The data

		Movie genre			
		Horror	Romcom	Docum.	Sci-fi
Happiness rating	61	68	82	90	
	57	78	84	87	
	70	65	90	75	
	65	57	75	88	
	67	55			85

* Note: Made-up data!

Master stats and ML — MX Cohen — sincxpress.com

After running the ANOVA on the above data, we get the following ANOVA table

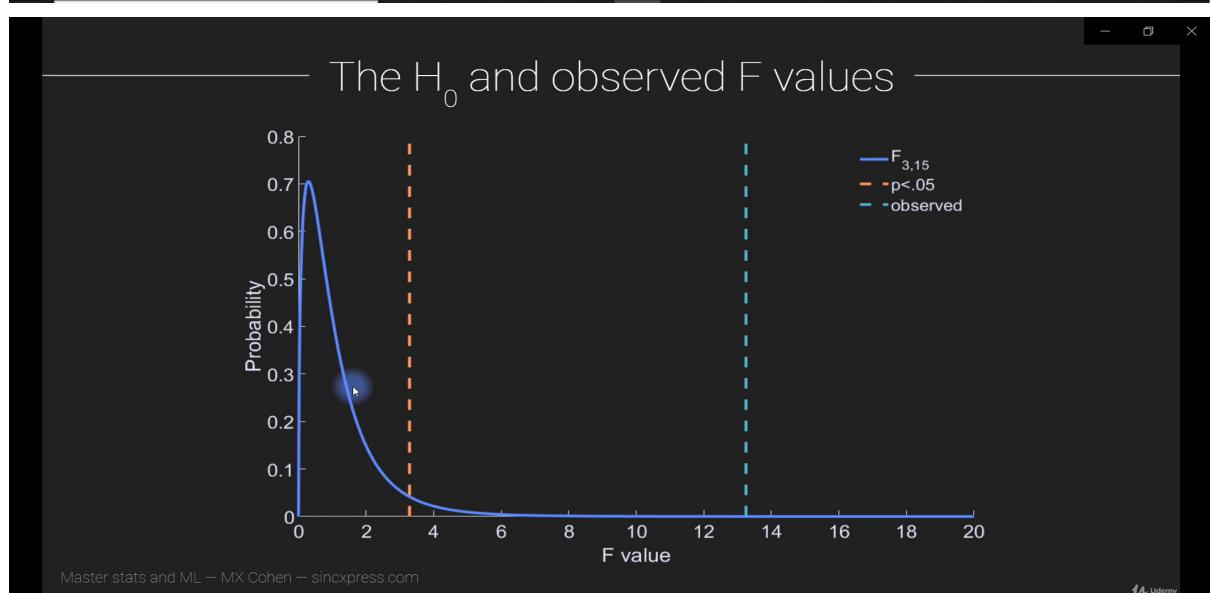
The ANOVA table

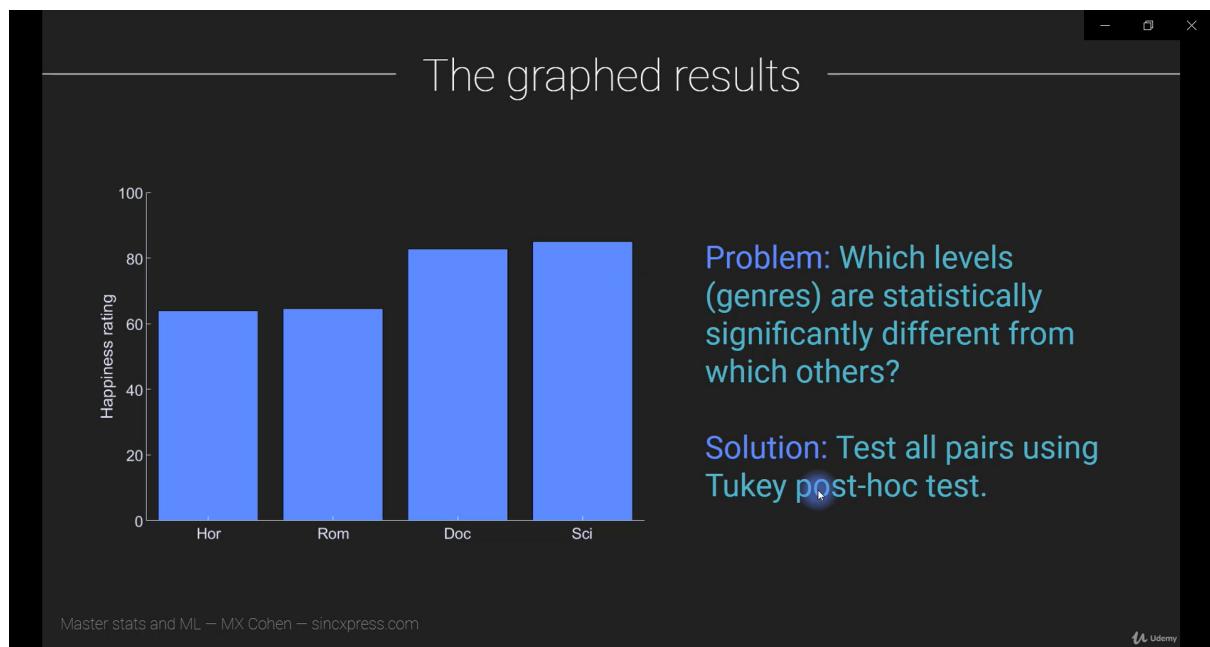
	Sums of squares	Degrees of freedom	Mean square	F	P-value
Between groups	1850.47	3 (k-1)	616.82	13.25	.0002
Within groups	697.95	15 (N-k)	46.53		
Total	2548.42	18 (N-1)			

Conclusion: The mean of at least one level is statistically significantly different from the mean of at least one other level.

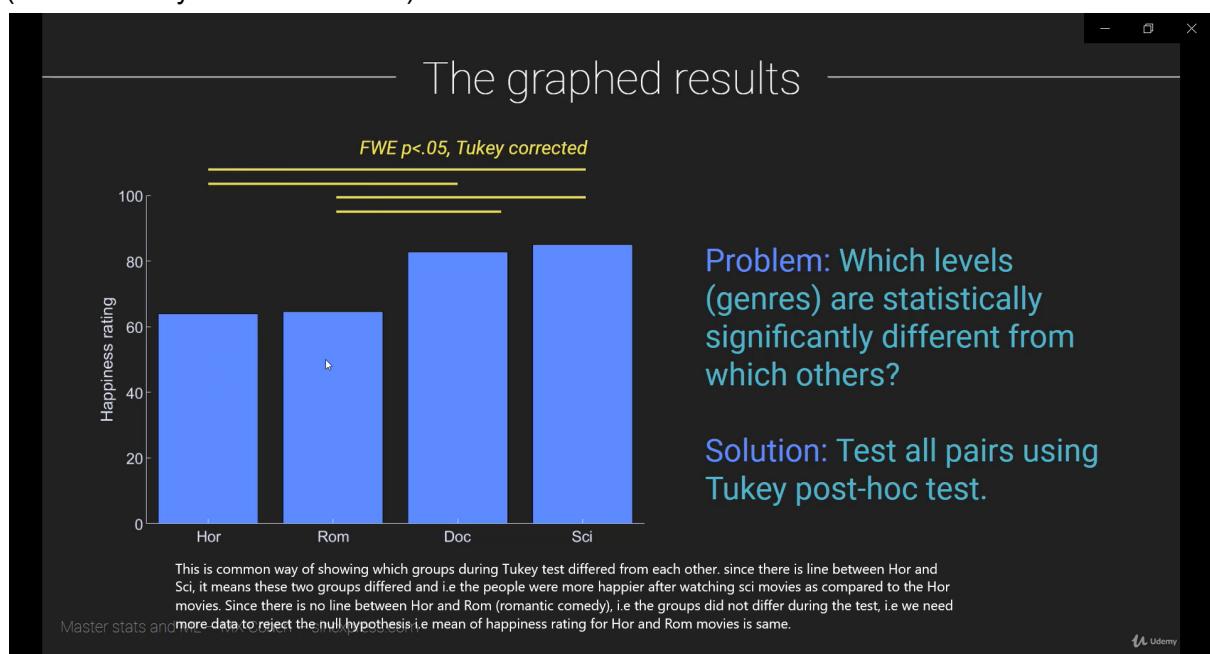
People were happier after at least one genre compared to at least one other genre.

Master stats and ML — MX Cohen — sincxpress.com





(FWE: Family-Wise Error rate)



Code

One-way ANOVA with repeated measures: Code

Two-Way ANOVA Example:

Experiment: Are girls bad at math?

Research goal: Test whether girls differ from boys in STEM and non-STEM subjects.

Variables: IV: class gender (boys, girls, mixed-boys, mixed-girls), subject (math, history).
DV: exam scores.

Experiment design: Three classes of 20 students each study the same material from the same teacher. An exam is given after two months.

*note: made-up numbers but conclusions are based on actual research.

Master stats and ML — MX Cohen — sinxpress.com

Udemy

The ANOVA table

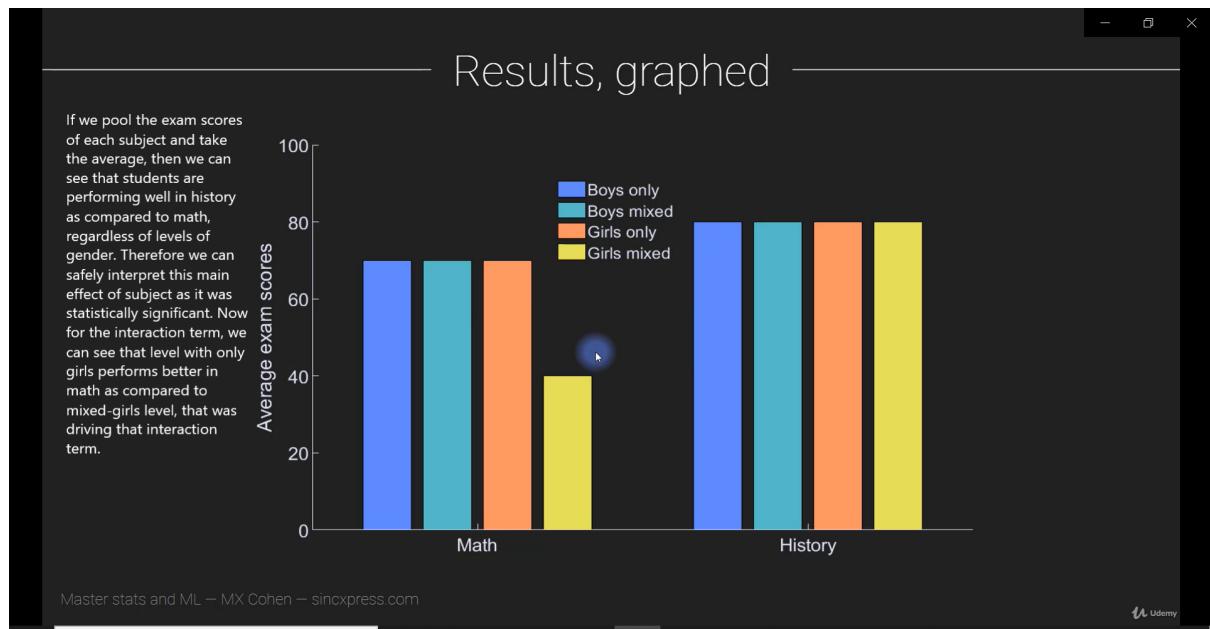
Source of variance	Sums of squares	Degrees of freedom	Mean square	F	P-value
Factor group		3 (4-1)			.456
Factor subj.		1 (2-1)			.028
G×S interact.		3 (4-1)(2-1)			.018
Within (error)		52 (60-4*2)			
Total		59 (60-1)			

Conclusion: There is a statistically significant interaction between Group and subject. The interaction makes the main effect of subject difficult to interpret without subsequent visualization and possible post-hoc tests.

Note the various ways of indicating factor names.

Master stats and ML — MX Cohen — sinxpress.com

Udemy



Two-way mixed effects ANOVA: [Code](#)

12. Regression

Table of Content

For detailed lectures on Regression
Evaluating Regression Model using f-test:
Standardising Regression Coefficients
Polynomial Regression
Bayes Information Criterion

For detailed lectures on Regression

check out this playlist

<https://www.youtube.com/playlist?list=PLTNMv857s9WUI1Nz4SssXDKAELESXz-bi>

Evaluating Regression Model using f-test:

▶ F-Test in regression analysis (Hypothesis test using F-Statistic)

Standardising Regression Coefficients

The scales of the β 's

$$y = \beta_0 + \beta_1 s + \beta_2 h + \beta_3 (s \times h) + \epsilon$$

Beta parameters and IV scales

	Hours	Minutes	Seconds
β_1	.167	10	600

If we used the time studied in hours, then value of beta1 = 0.167, if it is in minutes, then beta1 = 10, if it is in seconds, then beta1 = 600. Therefore, the raw outputs (beta values) from regression are called unstandardized betas, as they are dependent on the scale of the data.

Master stats and ML – MX Cohen – sincxpress.com

Udemy

Difficulties with interpreting β 's

Unstandardized β coefficients change depending on the scale of the IV.

Unstandardized β coefficients can be difficult (or impossible) to compare across variables (and studies...).

These difficulties motivate a standardization of β coefficients.

Master stats and ML – MX Cohen – sincxpress.com

Udemy

The scales of the β 's

$$y = \beta_0 + \boxed{\beta_1 s} + \beta_2 h + \beta_3 (s \times h) + \epsilon$$

Standardized beta parameters

	Hours	Minutes	Seconds	Calories
β_1	.6	.6	.6	
β_2				.8

say s=time you sleep and h=calories you intake before exam. Now after standardization of betas, see the value of beta1 is same irrespective of what the unit of s is. Also, from this we can interpret that the effect of calories you intake has more effect on exam score as compared to amount of time you sleep before exam

Master stats and ML — MX Cohen — sincxpress.com

Udemy

The scales of the β 's

Unstandardized β coefficients reflect the scales of the data (IV and DV). This can facilitate interpretation but can also stymie comparisons across variables or models.

Standardized β coefficients are in standard deviation units, unrelated to the scales of the data.

Both are correct and neither is better; sometimes one is more natural or easier to interpret than the other.

Importantly, standardization has no effect on the statistics!

Master stats and ML — MX Cohen — sincxpress.com

Udemy

How to standardize regressors

Basic idea: normalize β so that its variance is 1.

Method 1: z-normalize DV and IVs before the regression. All β 's will be in the units of the data, which are already standard deviation units.

Method 2: Scale the unstandardized β by the standard deviations of the IV and corresponding DV. $b_k = \beta_k \frac{s_{x_k}}{s_y}$

Master stats and ML — MX Cohen — sincxpress.com

Udemy

Interpreting standardized β 's

$$y = \beta_0 + \beta_1 s + \boxed{\beta_2 h} + \beta_3 (s \times h) + \epsilon$$

Interpretation: β_2 reflects the effect of a one-standard deviation change in h on standard deviation changes in y , when all other variables are held constant.

Polynomial Regression

Polynomial regression

$$y = \beta_0 + \beta_1 x + \epsilon$$

$$y = \beta_0 + \beta_1 x + \beta_2 x^2 + \epsilon$$

:

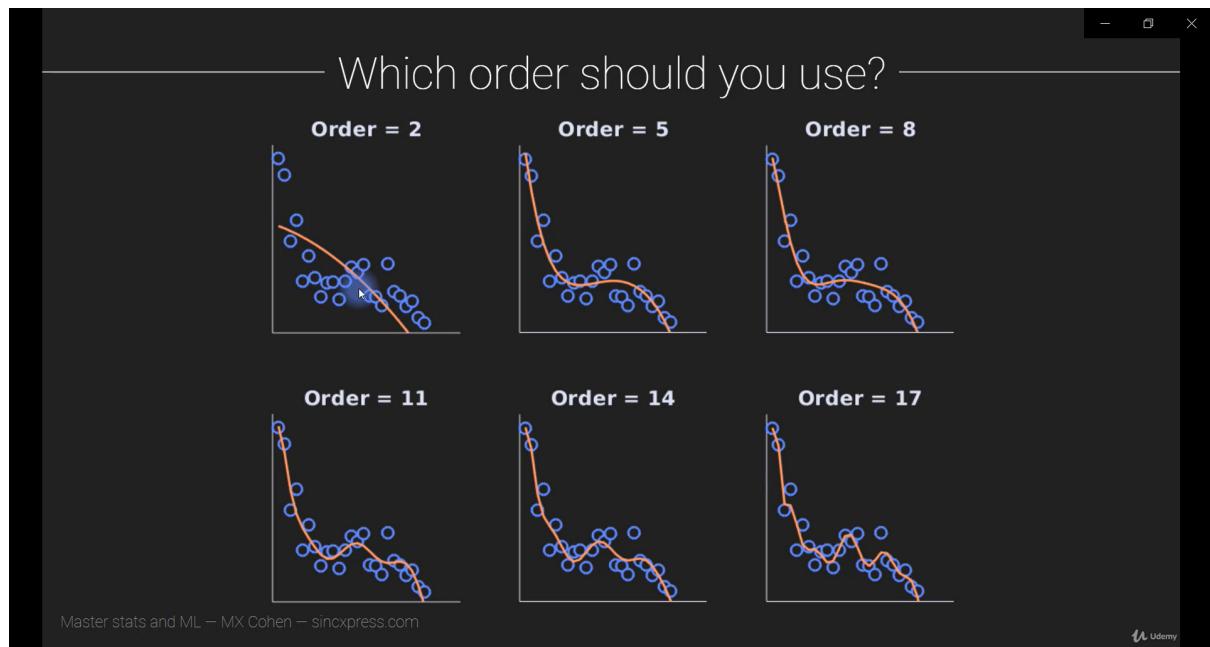
$$y = \beta_0 x^0 + \beta_1 x^1 + \dots + \beta_k x^k + \epsilon$$

Polynomial regression

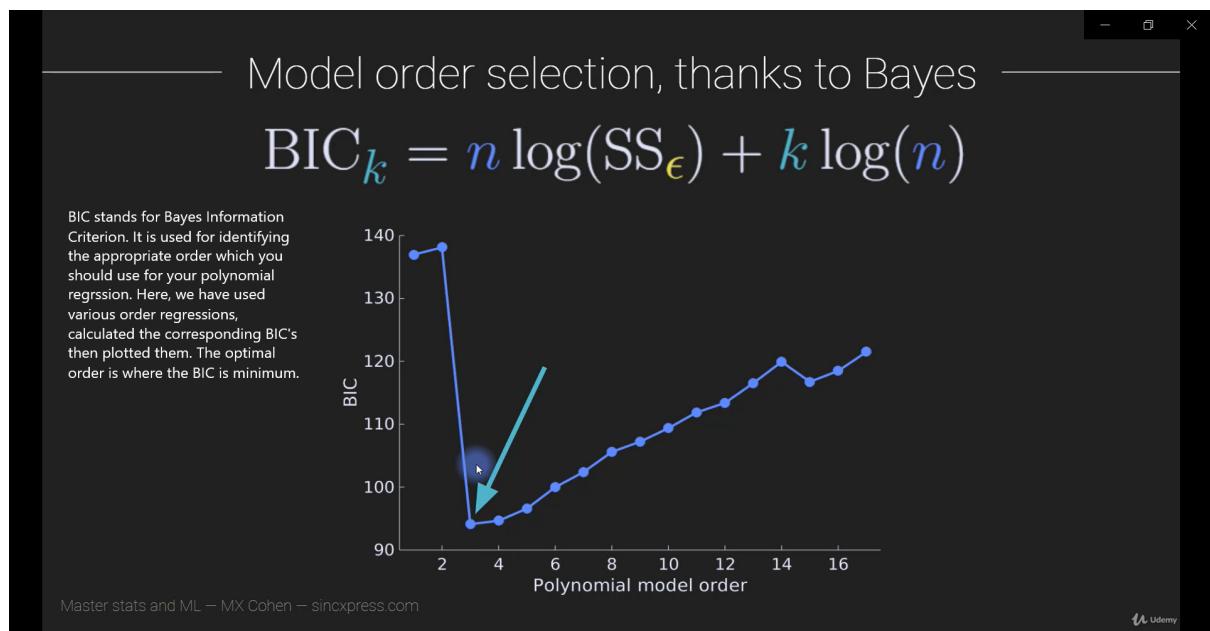
But how can we fit this kind of model with the nonlinearities?

The coefficients are all linear! This is a standard linear model!

$$y = \beta_0 x^0 + \beta_1 x^1 + \dots + \beta_k x^k + \epsilon$$



Bayes Information Criterion

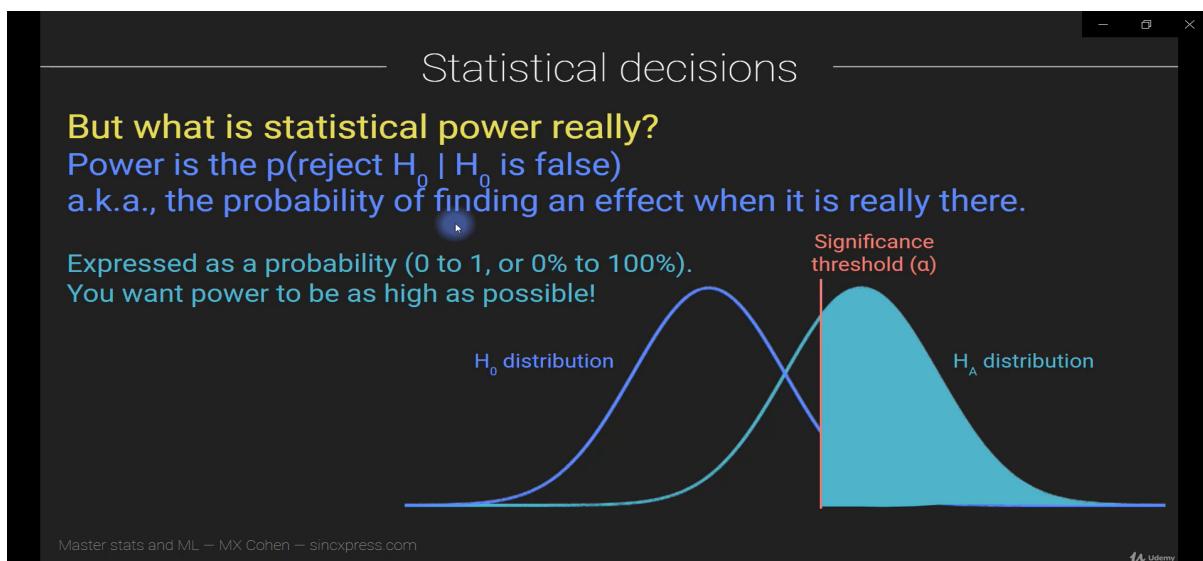
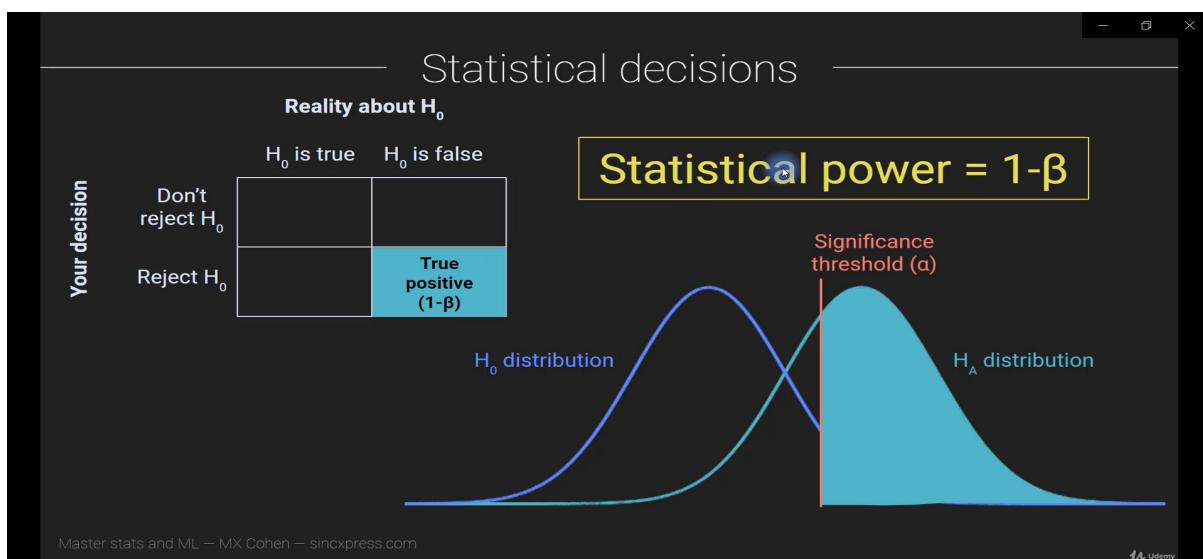
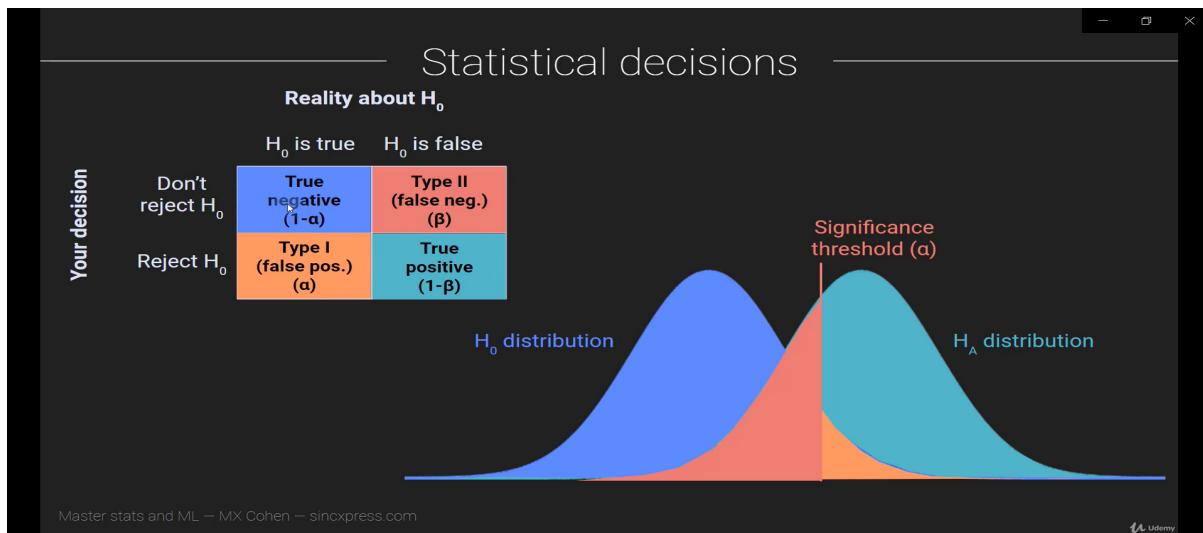


13. Statistical Power and Sample Sizes

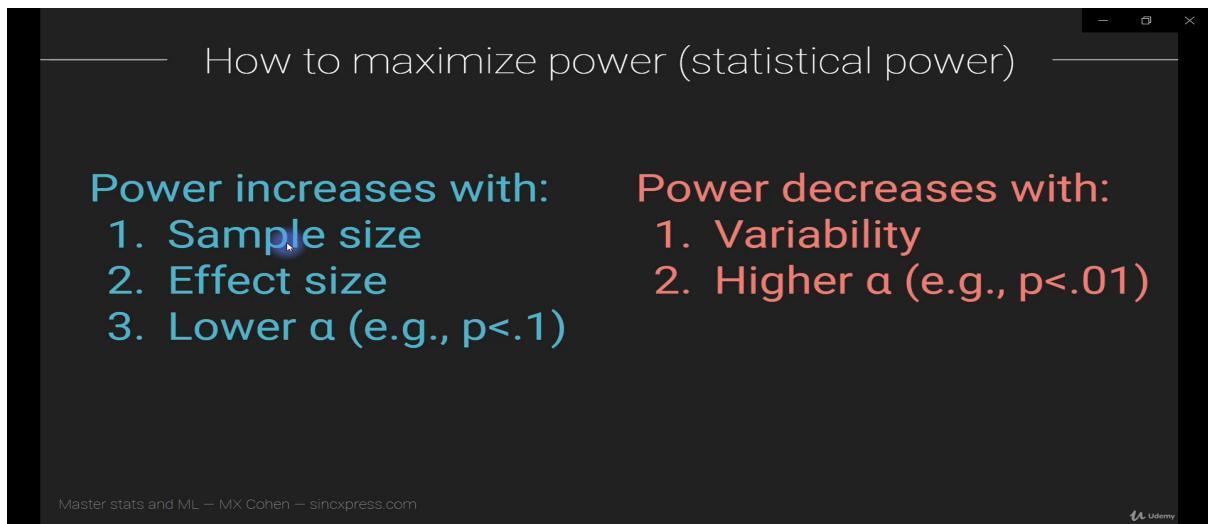
Table of Content

What is Statistical Power ?
How to maximise power ?
What is Effect Size ?
Effect of effect size on Power:
Problems with Power
Estimating Power and Sample Sizes
Calculating sample size
A Priori Power vs post-hoc power
Great Online Tool for Calculating Power

What is Statistical Power ?



How to maximise power ?



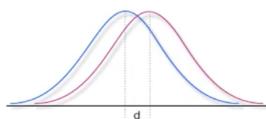
The slide has a dark background with white text. The title 'Power' is at the top. Below it is a bulleted list of five points. At the bottom left, it says 'RECORDED WITH SCREENCASTOMATIC'. At the bottom right, there is a small number '3'.

- Jacob Cohen is the father of power analysis
- Power is described as $1 - \beta$
- Acceptable power is .80 or higher – some say $> .70$ is adequate and $> .90$ is excellent
- Power analysis is usually done *a priori*, but can also be done afterward - controversial

The slide has a dark background with white text. The title 'Power' is at the top. Below it is a bulleted list of five points. At the bottom left, it says 'RECORDED WITH SCREENCASTOMATIC'. At the bottom right, there is a small number '5'.

- Power is a function of
 - 1) alpha level
 - 2) sample size (most dependent on this)
 - 3) effect size
 - 4) the type of statistical test being conducted
 - 5) the type of design used

What is Effect Size ?



Effect size

- Effect size is a quantitative measure of the *strength of a phenomenon*.
- Effect size emphasizes the **size** of the difference or relationship
- Examples:
 - the correlation between two variables (specifically r^2)
 - $r=.1$ weak, $r=.5$ moderate, $r=.7$ strong, $r=.9$ very strong
 - the regression coefficient in a regression (B_0, B_1, B_2)
 - Relative to model and field
 - the mean differences in t tests (use Cohen's D)
 - $d = .2$ is small; $r = .5$ is medium; $r = .8$ is large
 - The mean differences in ANOVA (use eta)
 - $.01$ is small, $.06$ medium, $.14$ large

RECORDED WITH SCREENCASTOMATIC

6

Effect of effect size on Power:

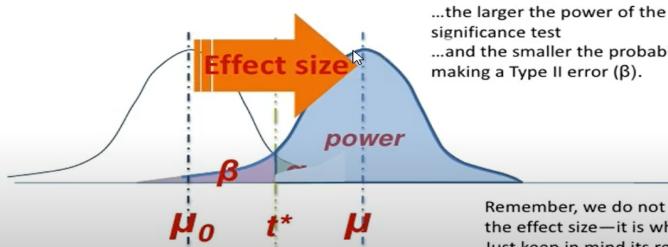
Greater the effect size (in the below example, greater the difference between the means, greater the effect size), greater will be the power.

Power & Effect Size

Types of errors and their probabilities

- How does effect size relate to power and β ?

The larger the effect size...
...the larger the power of the significance test
...and the smaller the probability of making a Type II error (β).



Remember, we do not control the effect size—it is what it is. Just keep in mind its relation to both power and β .

RECORDED WITH SCREENCASTOMATIC

00:07:04 / 11:04

How they are related

- If the effect size is small to medium, then you will need more subjects to find a significant result
- If the effect size is large, you do not need as many subjects (even as little as 15 per group)

Problems with Power

The problems with power

- Increasing alpha increases power but also increases $p(\text{Type I})$.
- Some factors that influence power you can control; others you cannot (sample size, effect size, variability).
- Power can differ for different analyses in the same experiment.
- Computing “true power” requires knowing the true effect size.
- Power calculations from published studies are unreliable due to “publication bias” (non-significant findings are not reported).

Conclusion:

Statistical power is more of a useful guideline than a precise and trustworthy numerical value.

Estimating Power and Sample Sizes

Formula for calculating power

$$z = \frac{\bar{x} - \mu_0}{\sigma / \sqrt{n}} = \frac{(\bar{x} - \mu_0) \sqrt{n}}{\sigma}$$

z "Z" value for power (next slide)
 \bar{x} Effect size
 μ_0 H_0 value
 σ / \sqrt{n} Standard error

Master stats and ML — MX Cohen — sincxpress.com

Formula for calculating power

$$z = \frac{\bar{x} - \mu_0}{\sigma / \sqrt{n}} = \frac{(\bar{x} - \mu_0) \sqrt{n}}{\sigma}$$

$1 - \beta = p(Z > z)$

Master stats and ML — MX Cohen — sincxpress.com

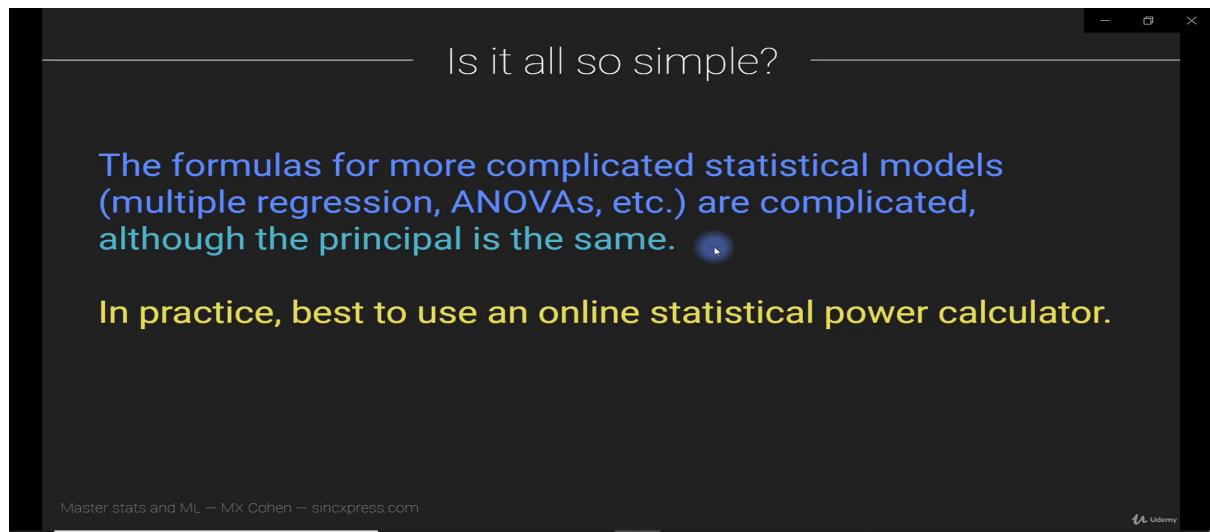
Calculating sample size

Formula for calculating sample size

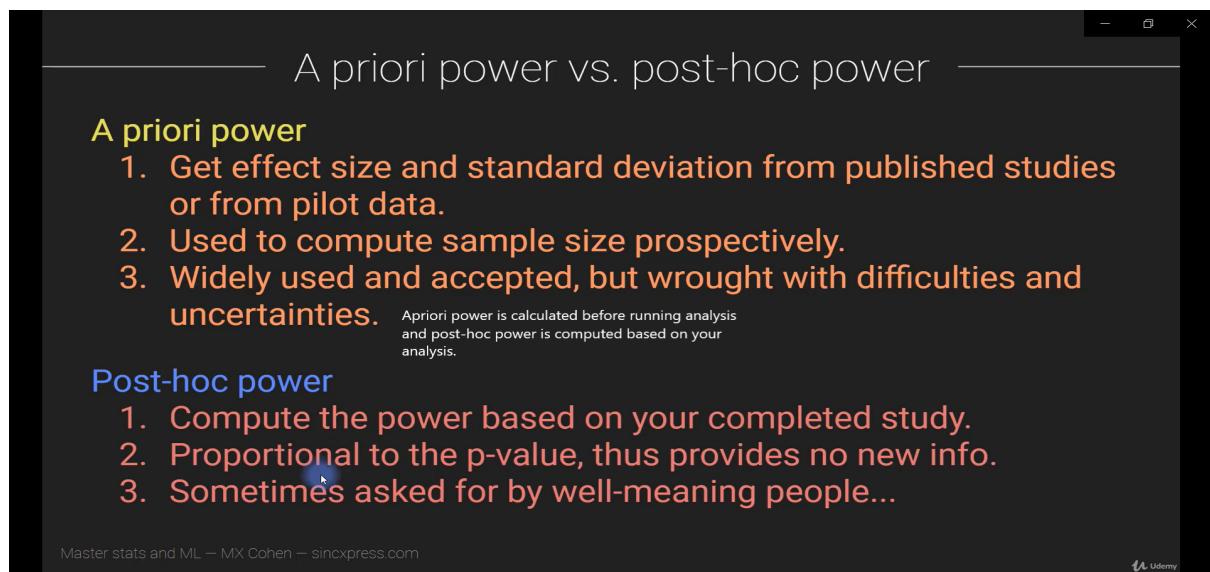
say we fix the statistical power, then we can find different quantities in this formula. Here, for the given power, std, and given effect size, we compute the corresponding sample size.

$$z = \frac{(\bar{x} - \mu_0) \sqrt{n}}{\sigma}$$
$$n = \left(\frac{z \sigma}{\bar{x} - \mu_0} \right)^2$$

Master stats and ML — MX Cohen — sincxpress.com



A Priori Power vs post-hoc power



Great Online Tool for Calculating Power

