

Intro to Machine Learning – Week 1



Supervised Learning

Hypothesis Testing
ML Concepts & Terminology
How to use Google Colab

What will we focus on?

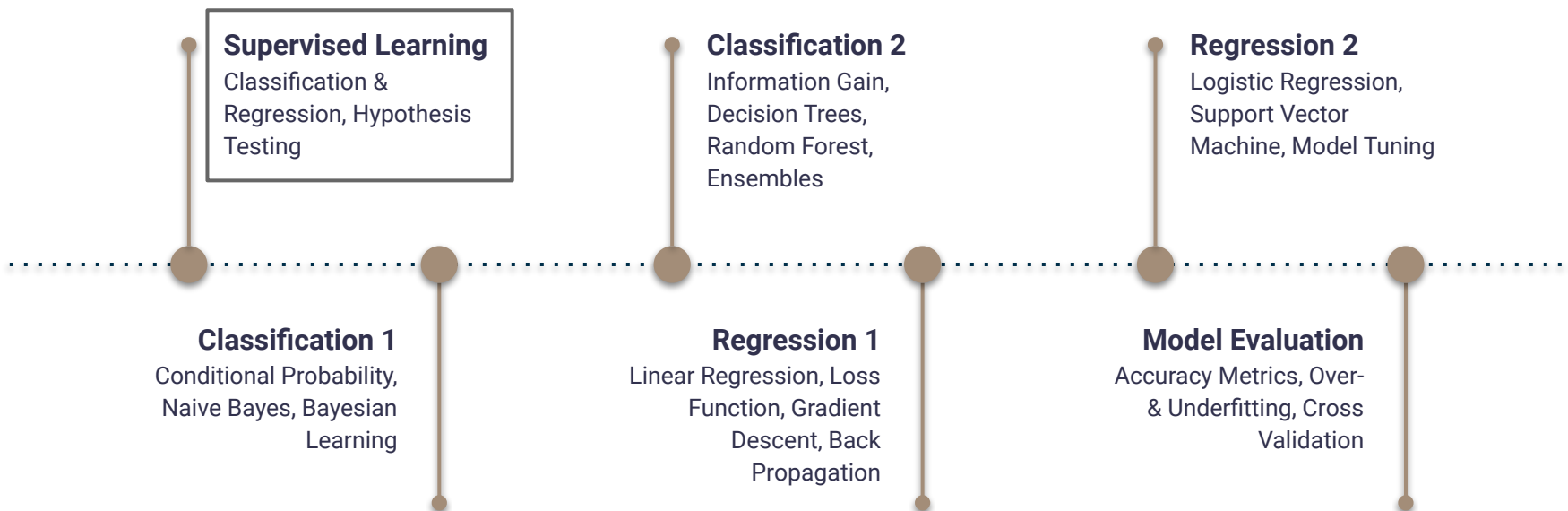
Concepts, Problems

1 hour

Google Colab Project

1 hour

Schedule



Machine Learning

What is Machine Learning?

Alan Turing proposed to change the question from "Can machines think?" to "Can machines do what we (as thinking entities) can do?"

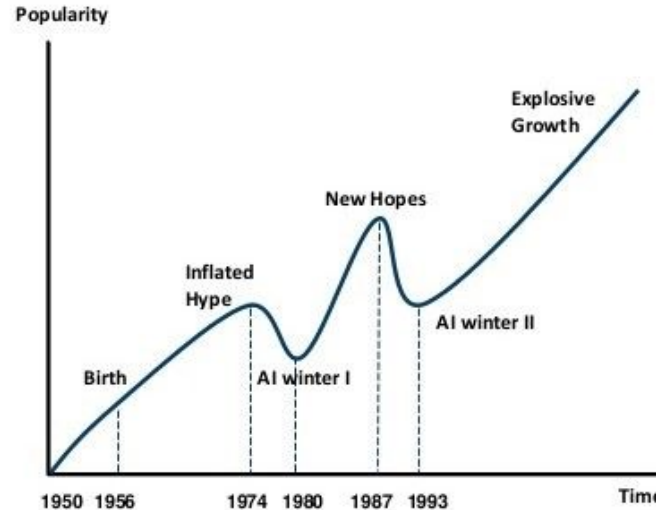
- Subset of **Artificial Intelligence**
- It is a study of computer algorithms that **improve their performance** at a task through “**experience**”
- Experience is a **formula** that is built to **generalize/fit** (most of) your data

What is Machine Learning?

AI Winters were collapses in the perception of the value of AI by government bureaucrats and venture capitalists.

Researchers continued to make advances despite the criticism.

AI HAS A LONG HISTORY OF BEING “THE NEXT BIG THING” ...



Timeline of AI Development

- **1950s-1960s:** First AI boom - the age of reasoning, prototype AI developed
- **1970s:** AI winter I
- **1980s-1990s:** Second AI boom: the age of Knowledge representation (appearance of expert systems capable of reproducing human decision-making)
- **1990s:** AI winter II
- **1997:** Deep Blue beats Gary Kasparov
- **2006:** University of Toronto develops Deep Learning
- **2011:** IBM's Watson won Jeopardy
- **2016:** Go software based on Deep Learning beats world's champions

What is Machine Learning?

Traditionally, software engineering has combined human-created rules with data to create answers to a problem. Instead of that, machine learning uses data and answers to discover the rules behind a problem

François Chollet

- ◆ With major advancements in computational power over the last decade, machine learning algorithms are being used in a wide **variety of applications**.
- ◆ Machine Learning is part of the Data Science process by **extracting** useful (non-trivial, potentially actionable) **knowledge** from large bodies of data.

Terminology

Example Data, x

- ◆ A dataset which contains information on certain qualities of objects/people.
- ◆ Eg. A .csv of pet information containing the columns name, age, type, owner name, contact.

Target Response, $f(x)$, y , \hat{y}

- ◆ A specific feature/column in the dataset which you would like to predict
- ◆ Eg. The age column from above .csv.

Types of Machine Learning

- ◆ Based on the **availability and data types** of **example data** and **target responses** machine learning can be broadly classified into
 - ▶ Supervised learning
 - ▶ Unsupervised learning
 - ▶ Semi-supervised learning
 - ▶ Reinforcement learning

Types of Machine Learning

- ◆ **Supervised Learning**
- ◆ *Unsupervised Learning*
- ◆ *Semi-supervised*
- ◆ *Reinforcement Learning*

- ◆ The algorithm learns from **example data and associated target responses**, in order to later predict the correct response when posed with new examples comes
- ◆ This approach is similar to a human learning **under the supervision** of a teacher.
 - ▶ The teacher provides good examples for the student to memorize
 - ▶ the student then derives general rules from these specific examples.

Types of Machine Learning

- ◆ *Supervised Learning*
- ◆ **Unsupervised Learning**
- ◆ *Semi-supervised*
- ◆ *Reinforcement Learning*

- ◆ The algorithm learns from **example data without any associated response**, leaving to the algorithm to **determine the data patterns** on its own.
- ◆ This is similar to methods humans use to figure out that certain objects or events are from the same class by observing the degree of similarity between objects.
 - ▶ Eg. When you recommend movies to friends

Types of Machine Learning

- ◆ *Supervised Learning*
- ◆ *Unsupervised Learning*
- ◆ **Semi-supervised**
- ◆ *Reinforcement Learning*

- ◆ Class of algorithms that are able to learn from **partially labeled data sets**
- ◆ Supervised Learning requires the dataset to be labeled either - this is a very **costly process**, especially when dealing with large volumes of data.
- ◆ Unsupervised Learning has the disadvantage of **limited application spectrum**.

Types of Machine Learning

- ◆ *Supervised Learning*
- ◆ *Unsupervised Learning*
- ◆ *Semi-supervised*
- ◆ **Reinforcement Learning**

- ◆ The algorithm learns from **example data** along with **positive/negative feedback** to each prediction the algorithm proposes.
- ◆ The algorithm must **make decisions** and the decisions **bear “consequences”**.
- ◆ In the human world, it is just like learning by **trial and error**.

Supervised Learning

Given: Training examples (x,y) for an unknown function f

Find: A good approximation of function f , which can be used for prediction if x is known

What is Supervised Learning?

- ◆ The algorithm learns from **example data and associated target responses**, in order to later predict the correct response when posed with new examples comes
- ◆ Based on the **data type of the target responses**, supervised learning can be broadly split into
 - ▶ **Classification** - for categorical data
 - ▶ **Regression** - for continuous data

Terminology

Categorical/Nominal

- ◆ Categorical values that have 2 or more categories/buckets.
- ◆ Eg. Colours like green, blue, red,.. etc.
- ◆ Do not have an explicit ordering from its values (eg. is green > blue?)

Continuous/Numerical

- ◆ Continuous values are numerical in nature.
- ◆ Eg. Weights like 0.5g, 10g, 2.5kg,.. etc
- ◆ Contain explicit ordering from its values (eg. $0.5 < 10 < 2500$)

Types of Supervised Learning

◈ Classification

◈ *Regression*

- ◈ Target responses are **categorical** in nature
- ◈ Eg. A bag contains 15 balls of colours **green**, **blue** and **red**. If we know the **diameters** and **colours** of 10 balls, is it possible to predict the colour of the other 5 balls based on their diameters?
 - ▶ Diameter is the example data
 - ▶ Colour is the target response

Types of Supervised Learning

◆ *Classification*

◆ **Regression**

- ◆ Target responses are **continuous/numerical** in nature
- ◆ Eg. If you know the height and weight of 20 people and the corresponding age of 15 of them. Can you predict the age for the remaining 5 people?
 - ▶ Height, weight is the example data
 - ▶ Age is the target response

Appropriate applications for Supervised Learning

- ◆ Situations where there is **no human expert**
 - ▶ x - bonding strength for a new molecule
 - ▶ $f(x)$ - binding strength to AIDS molecule
- ◆ Situations where humans can perform the task but **cannot describe how they do it**
 - ▶ x - image of a handwritten character
 - ▶ $f(x)$ - Corresponding ASCII character
- ◆ Situations where the desired function is **changing frequently**
 - ▶ x - stock prices in the last 10 days
 - ▶ $f(x)$ - recommended stock transactions
- ◆ Situations where each user needs a **customized function f**
 - ▶ x - Incoming email message
 - ▶ $f(x)$ - Importance for presenting to user

Hypothesis Testing

What is Hypothesis Testing?

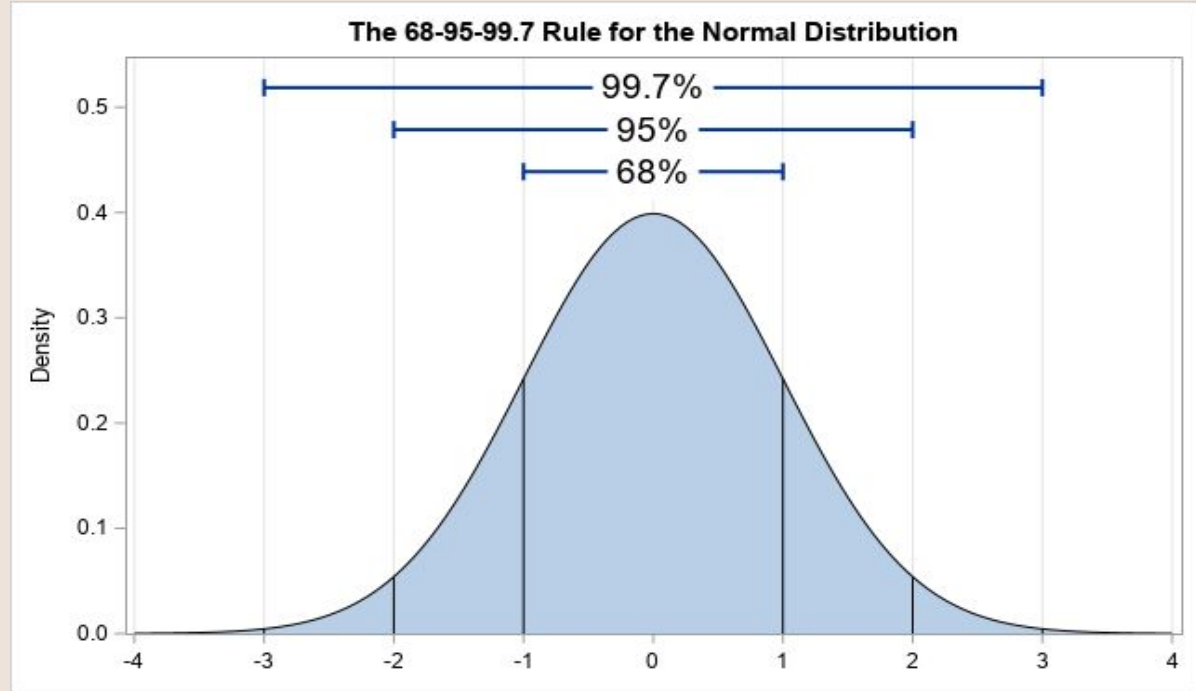
Data alone is not interesting. It is the **interpretation of the data** that we are really interested in.

- ◆ A hypothesis test calculates some quantity **under a given assumption**. The result of the test allows us to interpret whether the assumption holds or whether the assumption has been violated.
- ◆ Two concrete examples that we will **use a lot** in machine learning are:
 - Assume that data has a normal distribution.
 - Assume that two samples were drawn from the same underlying population distribution.

Terminology

Normal Distribution

- ◆ It is a type of **population distribution** that is commonly found.
- ◆ Eg. heights, blood pressure and IQ scores follow the normal distribution.



How to conduct Hypothesis Test?

Example – I have a bag with 15 balls in 3 colours

H_0 – There are equal number of balls in each colour

H_1 – H_0 is false

- ◆ The assumption of the statistical test is known as **null hypothesis**.
 - ▶ It is denoted by H_0
 - ▶ It is a **commonly accepted fact**
 - ▶ It is often called the **default assumption**, or the assumption that nothing has changed.
- ◆ A **violation of the assumption** is the first hypothesis
 - ▶ It is denoted by H_1
 - ▶ It really means “*some other hypothesis*,” as all we know is that the evidence suggests that the H_0 can be rejected.

p-value, α

Smaller α suggests a more
robust interpretation of
the null hypothesis

- ◆ **p-value** is a quantity that we can use to interpret or quantify the result of the test
- ◆ **α** is the significance level that is used to either accept or reject the hypothesis
 - ▶ It commonly is 5% or **0.05**.
 - ▶ **p-value > α** : Fail to reject the H_0 (not significant result).
 - ▶ **p-value $\leq \alpha$** : Reject the null H_0 (significant result).

Interpreting p-value, α

Confidence level can be
calculated by subtracting
significance level from 1

$$\text{confidence} = 1 - \alpha$$

H_0 = Data is normally distributed

H_1 = Data is not normally distributed

p-value = 0.7, α = 0.05

- ◆ The test found that the data sample was normal, failing to reject the null hypothesis at a 5% significance level.
- ◆ The test found that the data was normal, failing to reject the null hypothesis at a 95% confidence level

Errors in Statistical Tests

- ◆ The interpretation of a statistical hypothesis test is **probabilistic** - the evidence of the test may suggest an outcome and be **mistaken**.
- ◆ There are 2 types of such errors
 - ▶ **Type I Error**: The **incorrect rejection** of a true H_0 (ie. false positive)
 - ▶ **Type II Error**: The **incorrect failure of rejection** of a false H_0 (ie. false negative)

How to derive the p-value?

- ◆ **Based on the problem**, there are multiple test statistics that can be used to derive the p-value.
- ◆ Sampling distribution under the null hypothesis must be **calculable, either exactly or approximately**.
- ◆ H_0 also helps decide whether p-value should be **one-tailed or two-tailed**, depending on the values being testing

Types of Test Statistics

Name	Formula	Assumptions or notes
One-sample z-test	$z = \frac{\bar{x} - \mu_0}{(\sigma/\sqrt{n})}$	(Normal population or $n > 30$) and σ known. (z is the distance from the mean in relation to the standard deviation of the mean). For non-normal distributions it is possible to calculate a minimum proportion of a population that falls within k standard deviations for any k (see: Chebyshev's inequality).
Two-sample z-test	$z = \frac{(\bar{x}_1 - \bar{x}_2) - d_0}{\sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}}$	Normal population and independent observations and σ_1 and σ_2 are known
One-sample t-test	$t = \frac{\bar{x} - \mu_0}{(s/\sqrt{n})},$ $df = n - 1$	(Normal population or $n > 30$) and σ unknown
Paired t-test	$t = \frac{\bar{d} - d_0}{(s_d/\sqrt{n})},$ $df = n - 1$	(Normal population of differences or $n > 30$) and σ unknown
Two-sample pooled t-test , equal variances	$t = \frac{(\bar{x}_1 - \bar{x}_2) - d_0}{s_p \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}},$ $s_p^2 = \frac{(n_1 - 1)s_1^2 + (n_2 - 1)s_2^2}{n_1 + n_2 - 2},$ $df = n_1 + n_2 - 2$ ^[3]	(Normal populations or $n_1 + n_2 > 40$) and independent observations and $\sigma_1 = \sigma_2$ unknown

Terminology

z-test

- ◆ Used to determine whether **two population means are different** when
 - their variances are known
 - the sample size is large.
- ◆ It is assumed to follow normal distribution
- ◆ A z-statistic, or z-score, is a number representing how many standard deviations above or below the mean population a score derived from a z-test is.

Example #1

A principal at a certain school claims that the students in his school are above average intelligence. A random sample of thirty students IQ scores have a mean score of 112.5. Is there sufficient evidence to support the principal's claim? The mean population IQ is 100 with a standard deviation of 15.

Example #1

$$\blacklozenge \bar{x} = 112.5$$

$$\blacklozenge n = 30$$

$$\blacklozenge \mu = 100$$

$$\blacklozenge \sigma = 15$$

A principal at a certain school claims that the students in his school are above average intelligence. A random sample of thirty students IQ scores have a mean score of 112.5. Is there sufficient evidence to support the principal's claim? The mean population IQ is 100 with a standard deviation of 15.

Example #1

- ◆ $\bar{x} = 112.5$
- ◆ $n = 30$
- ◆ $\mu = 100$
- ◆ $\sigma = 15$

Step 1: Describe Null Hypothesis (H_0)

- ◆ The population mean is a common accepted fact
- ◆ H_0 is $\mu = 100$

Step 2: Describe First Hypothesis (H_1)

- ◆ The principal claims that the schools IQ is **above** average
- ◆ H_1 is $\mu > 100$
- ◆ One-tailed test, since we are looking at only “greater than”

Example #1

- ◆ $\bar{x} = 112.5$
- ◆ $n = 30$
- ◆ $\mu = 100$
- ◆ $\sigma = 15$
- ◆ $H_0: \mu = 100$
- ◆ $H_1: \mu > 100$

Step 3: Pick test statistic to determine p-value

Z-Test would be most suitable test statistic to use for p-value, since sample and population means are known

$$Z = \frac{\bar{x} - \mu}{\frac{\sigma}{\sqrt{n}}}$$

$$\begin{aligned} Z &= (112.5 - 100) / (15 / \sqrt{30}) \\ &= 12.5 / 2.7386 \\ &= 4.564 \end{aligned}$$

Example #1

- ◆ $H_0 - \mu = 100$
- ◆ $H_1 - \mu > 100$
- ◆ $p\text{-value} = 4.564$

Step 4: Determine α

Find the corresponding z-score for the α we want to use

- ◆ $\alpha = 0.05$
- ◆ You can use [online calculator](#) to do this conversion
- ◆ $\alpha = 0.05 \Rightarrow z = 1.6449$

Example #1

- ◆ $H_0 - \mu = 100$
- ◆ $H_1 - \mu > 100$
- ◆ $p\text{-value} = 4.564$
- ◆ $\alpha = 1.645$

Step 4: Compare p-value and α

From the values calculated, we can confidently say that **p-value** > α

Step 5: Draw a conclusion

Since $p\text{-value} > \alpha$, we fail to reject the H_0 . This means that we **cannot substantiate the principals claims**.

Terminology

t-test

- ◆ Used to evaluate if the means of two sets of data are statistically significantly different from each other
- ◆ There are **3 types** of t-tests,
 - ▶ **One-sample** - compare the mean of a population with a theoretical value.
 - ▶ **Unpaired Two-sample** - compare the mean of two **independent samples**.
 - ▶ **Paired** - compare the means between two related groups of samples.
- ◆ **Degrees of freedom** is 1 less than sample size ($n-1$)

Example #2

A study examined the effect of diet cola consumption on calcium levels in women. A sample of healthy women aged 18–40 were randomly assigned to drink 24 ounces of either diet cola or water. Their urine was collected for three hours after ingestion of the beverage and calcium excretion (in mg) was measured.

Variable	Drink	Samples (n)	Mean (\bar{x})	Std Dev (σ)
Calcium	Cola	8	56	4.93
	Water	8	49.1	3.64

- ◆ The researchers were investigating whether diet cola leaches calcium out of the system, which would increase the amount of calcium in the urine for cola drinkers. Assume there are no outliers.

Example #2

$$\blacklozenge \bar{x}_c = 56.0$$

$$\blacklozenge n_c = 8$$

$$\blacklozenge \sigma_c = 4.93$$

$$\blacklozenge \bar{x}_w = 49.1$$

$$\blacklozenge n_w = 8$$

$$\blacklozenge \sigma_w = 3.64$$

A study examined the effect of diet cola consumption on calcium levels in women. A sample of healthy women aged 18–40 were randomly assigned to drink 24 ounces of either diet cola or water. Their urine was collected for three hours after ingestion of the beverage and calcium excretion (in mg) was measured.

Variable	Drink	Samples (n)	Mean (\bar{x})	Std Dev (σ)
Calcium	Cola	8	56	4.93
	Water	8	49.1	3.64

- ◆ The researchers were investigating whether diet cola leaches calcium out of the system, which would increase the amount of calcium in the urine for cola drinkers. Assume there are no outliers.

Example #2

- ◆ $\bar{x}_c = 56.0$
- ◆ $n_c = 8$
- ◆ $\sigma_c = 4.93$
- ◆ $\bar{x}_w = 49.1$
- ◆ $n_w = 8$
- ◆ $\sigma_w = 3.64$

Step 1: Describe Null Hypothesis (H_0)

- ◆ μ is the mean calcium loss after drink
- ◆ Cola **does not** increase calcium content is a common known fact
- ◆ H_0 is $\mu_c = \mu_w$

Step 2: Describe First Hypothesis (H_1)

- ◆ Increased calcium content in cola drinkers
- ◆ H_1 is $\mu_c > \mu_w$
- ◆ One-tailed test, since we are looking at only “greater than” /upper tail test

Example #2

$$\diamond \bar{x}_c = 56.0$$

$$\diamond n_c = 8$$

$$\diamond \sigma_c = 4.93$$

$$\diamond \bar{x}_w = 49.13$$

$$\diamond n_w = 8$$

$$\diamond \sigma_w = 3.64$$

$$\diamond H_0 - \mu_c = \mu_w$$

$$\diamond H_1 - \mu_c > \mu_w$$

Step 3: Pick test statistic to determine p-value

Since there are no outliers in the data, we can use a t-test as the test statistic

$$t = \frac{\text{Sample statistics} - \text{Null parameter}}{SE}$$

$$= \frac{(\bar{x}_C - \bar{x}_W) - 0}{\sqrt{\frac{s_C^2}{n_C} + \frac{s_W^2}{n_W}}} = \frac{56.0 - 49.1}{\sqrt{\frac{4.93^2}{8} + \frac{3.64^2}{8}}}$$

$$= 3.18$$

Example #2

Step 3: Pick test statistic to determine p-value

Since there are no outliers in the data, we can use a t-test as the test statistic

$$t = \frac{\text{Sample statistics} - \text{Null parameter}}{SE}$$

$$= \frac{(\bar{x}_C - \bar{x}_W) - 0}{\sqrt{\frac{s_C^2}{n_C} + \frac{s_W^2}{n_W}}} = \frac{56.0 - 49.1}{\sqrt{\frac{4.93^2}{8} + \frac{3.64^2}{8}}}$$

$$= 3.18$$

- ◆ $H_0: \mu_c = \mu_w$
- ◆ $H_1: \mu_c > \mu_w$
- ◆ $t\text{-stat} = 3.18$

Example #2

Step 4: Convert t-statistic to p-value

You can use an [online calculator](#) to convert the t-value to p-value

Degree of Freedom is 7, since there are sample size is 8

p-value = 0.0078

- ◆ $H_0 - \mu_c = \mu_w$
- ◆ $H_1 - \mu_c > \mu_w$
- ◆ $t\text{-stat} = 3.18$

Example #2

Step 4: Compare p-value and α

From the values calculated, we can confidently say that **p-value** < α

Step 5: Draw a conclusion

- Since p-value < α , we reject the H_0 .
- There is **strong evidence** for that diet cola drinkers do lose more calcium, on average, than water drinkers ($\mu_c > \mu_w$)

$$\blacklozenge H_0: \mu_c = \mu_w$$

$$\blacklozenge H_1: \mu_c > \mu_w$$

$$\blacklozenge \text{p-value} = 0.0078$$

$$\blacklozenge \alpha = 0.05$$

Theory Recap

- ◆ **Machine Learning**
 - ▶ History
 - ▶ Types
- ◆ **Supervised Learning**
 - ▶ Types
 - ▶ Applications
- ◆ **Hypothesis testing**
 - ▶ p-value, α
 - ▶ Types of errors
 - ▶ Choosing a test statistic
 - ▶ Examples

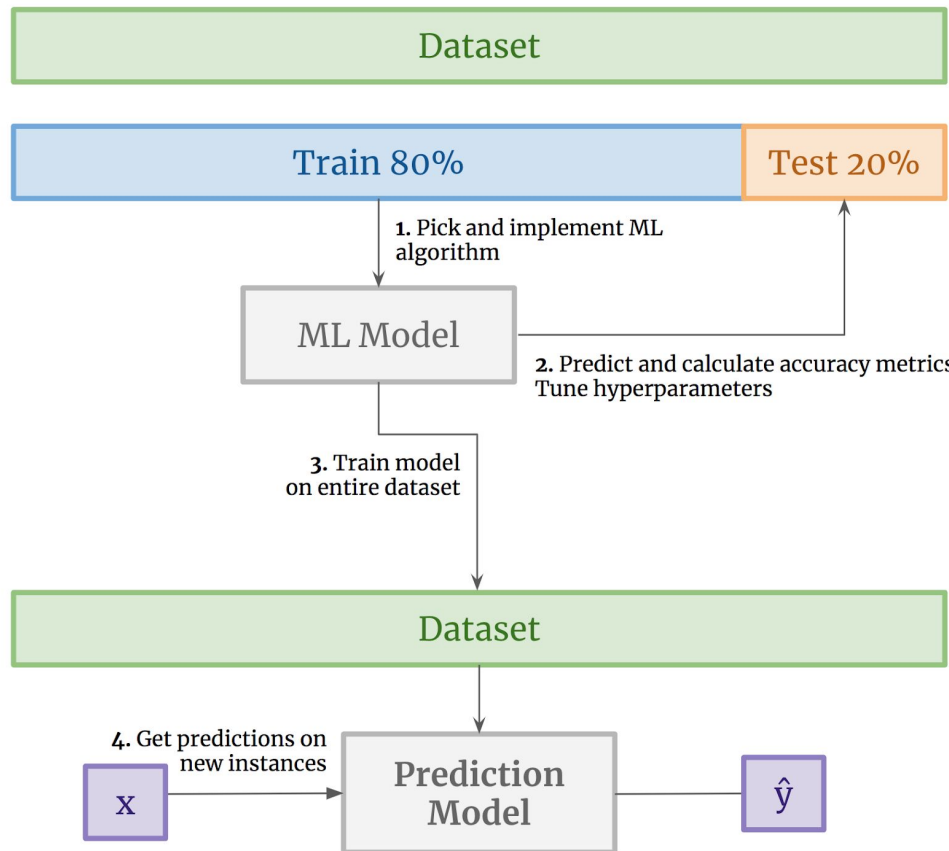
Google Colab Project

Example Dataset

- ◆ 14 rows
- ◆ Target - play
- ◆ Example Data - outlook, temp, humidity, windy

outlook	temperature	humidity	windy	play
overcast	cool	normal	TRUE	yes
overcast	hot	high	FALSE	yes
overcast	hot	normal	FALSE	yes
overcast	mild	high	TRUE	yes
rainy	cool	normal	FALSE	yes
rainy	cool	normal	TRUE	no
rainy	mild	high	FALSE	yes
rainy	mild	high	TRUE	no
rainy	mild	normal	FALSE	yes
sunny	cool	normal	FALSE	yes
sunny	hot	high	FALSE	no
sunny	hot	high	TRUE	no
sunny	mild	high	FALSE	no
sunny	mild	normal	TRUE	yes

Machine Learning Process



PlayGolf dataset (14 rows)

PlayGolf_Train (11), PlayGolf_Test(3)

Q: Will we play golf on a day that is (overcast, cool, high, FALSE)?

A: Yes

Google Colab & Github

Code at
[bit.ly/introtoml-
github](https://bit.ly/introtoml-github)

- ◆ You can use Google Colab to run **open source .ipynb** files.
- ◆ You can find a lot of cool machine learning algorithm implementations shared on Github.
- ◆ **No setup** required
- ◆ Google account **required**

Open Source Datasets

- **UCI Machine Learning Repository**
<https://archive.ics.uci.edu/ml/datasets.php>
- **Kaggle**
<https://www.kaggle.com/datasets>
- **US Government**
<https://www.data.gov/>
- **US Census**
<https://www.census.gov/data.html>

Homework #1

Learn how to load your own datasets in Google Colab by following along with this video

- <https://bit.ly/link-googlecolab-github>

Homework #2

Check out the **examples** listed for statistical tests on [Wikipedia](#)

- ◆ Are you able to understand these examples?

Homework #3

Hint: You can use the Z-test for this

Blood glucose levels for obese patients have a mean of 100 with a standard deviation of 15. A researcher thinks that a diet high in raw cornstarch will have a positive or negative effect on blood glucose levels. A sample of 30 patients who have tried the raw cornstarch diet have a mean glucose level of 140.

- ◆ Test the hypothesis that the raw cornstarch had an effect.

See you next week!

Questions?

Join us on slack
(bit.ly/wwcodedatascience-slack) and
post it on our **#help-me** channel.

Register?

Register for all sessions at
[linktr.ee/wwcodedatascience
_registration](https://linktr.ee/wwcodedatascience_registration)