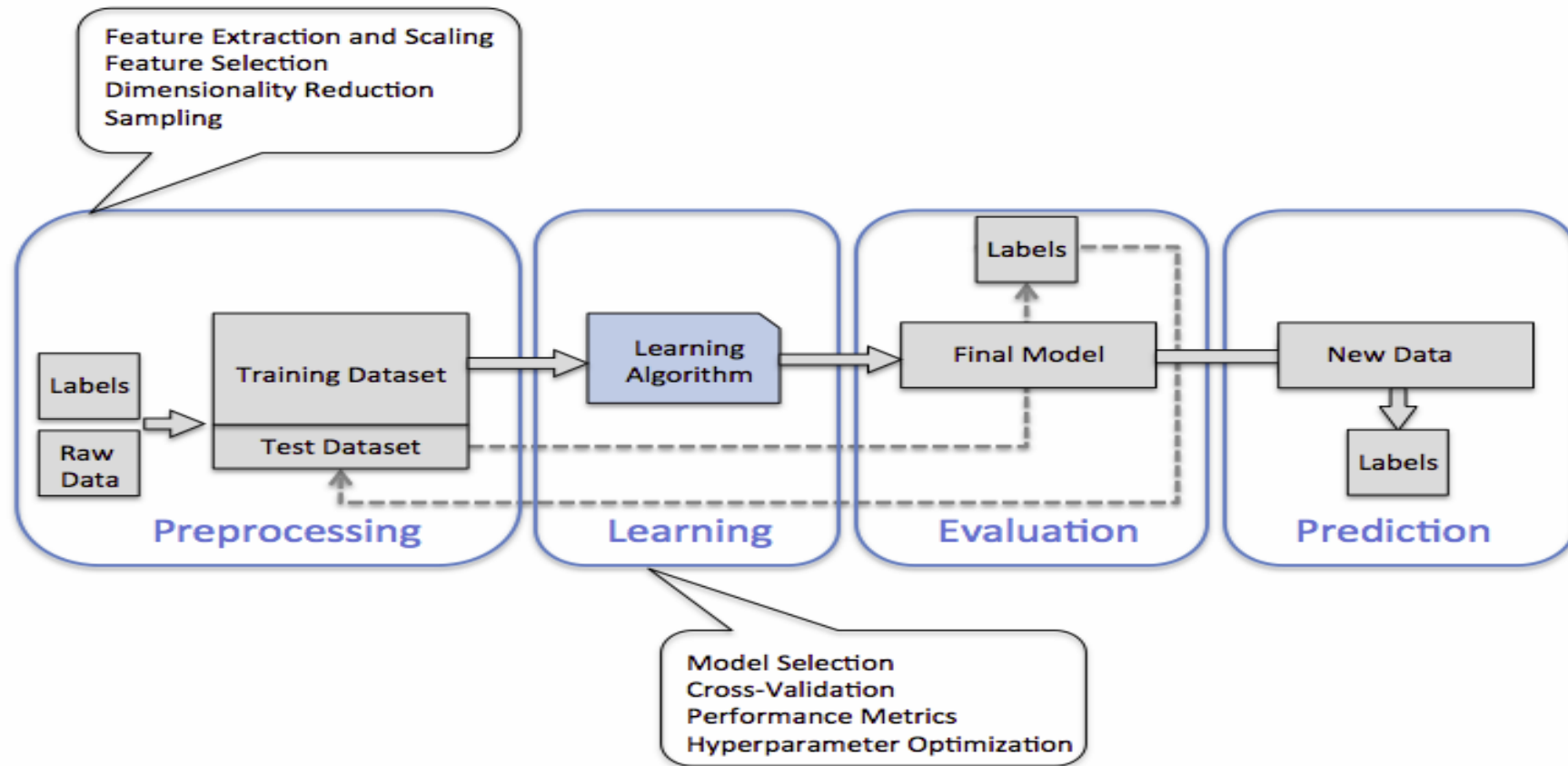


# NLP for SE and AI Techniques

# Agenda

- Advanced Evaluation Techniques (cont.)
- Supervised learning
  - KNNs
  - Linear Regression (OLS)
  - Ridge Regression
  - Lasso Regression

# A roadmap for building machine learning systems

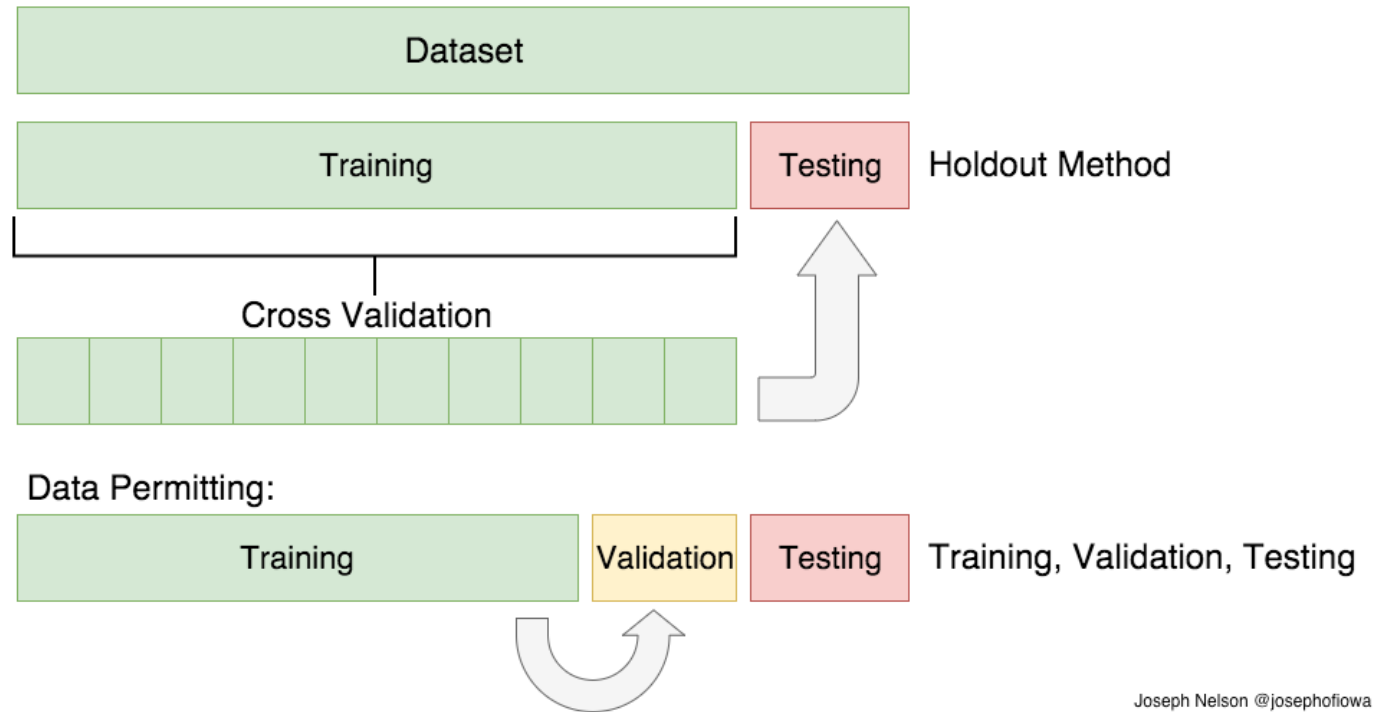


# ADV. Evaluation Techniques (when we don't have enough or quality dataset)

- Cross Validation
- Stratified Cross Validation
- Leave-One-Out Cross validation
- Bootstrapping

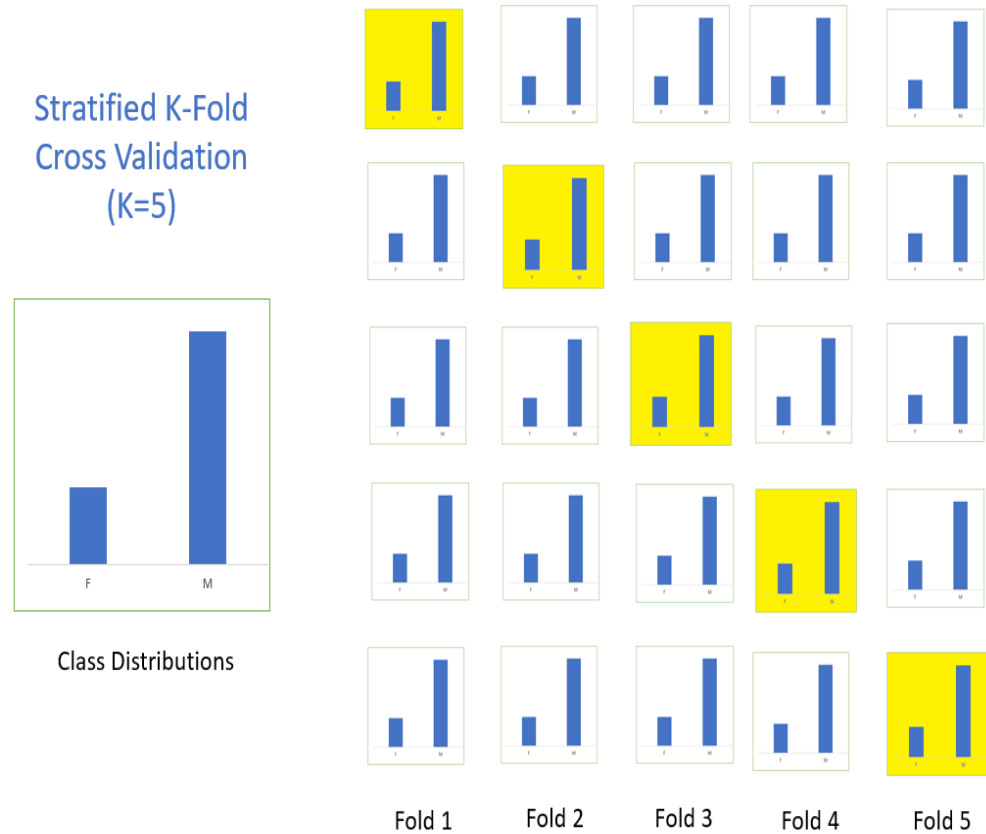
# K-folds Cross-validation Method

- AKA. Rotation estimation
- Use to estimate a performance of the model (i.e. mean of accuracy rate)



# Stratified Cross-validation Method

- Same as Cross-validation but here we ensure that each fold is representative of all strata of the class.



# Leave-one-out Cross-validation Method

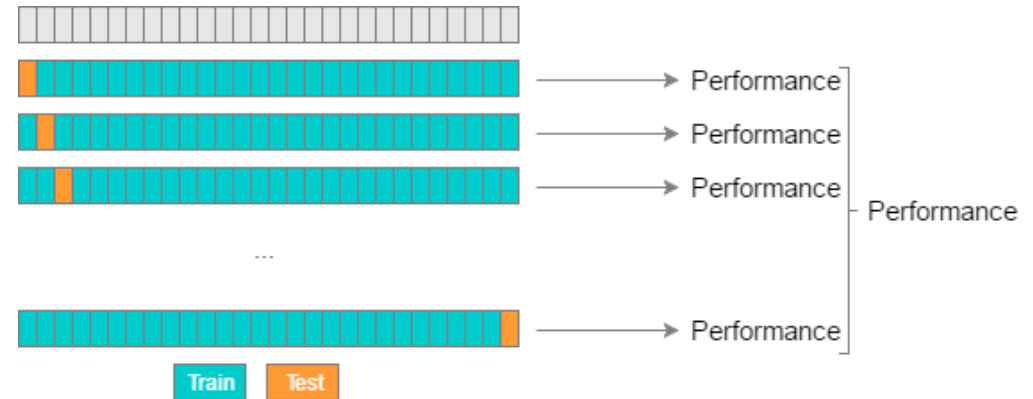
- Cross-validation for small sample size.
- The number of folds is the same as the number of training instances.

- Advantages:

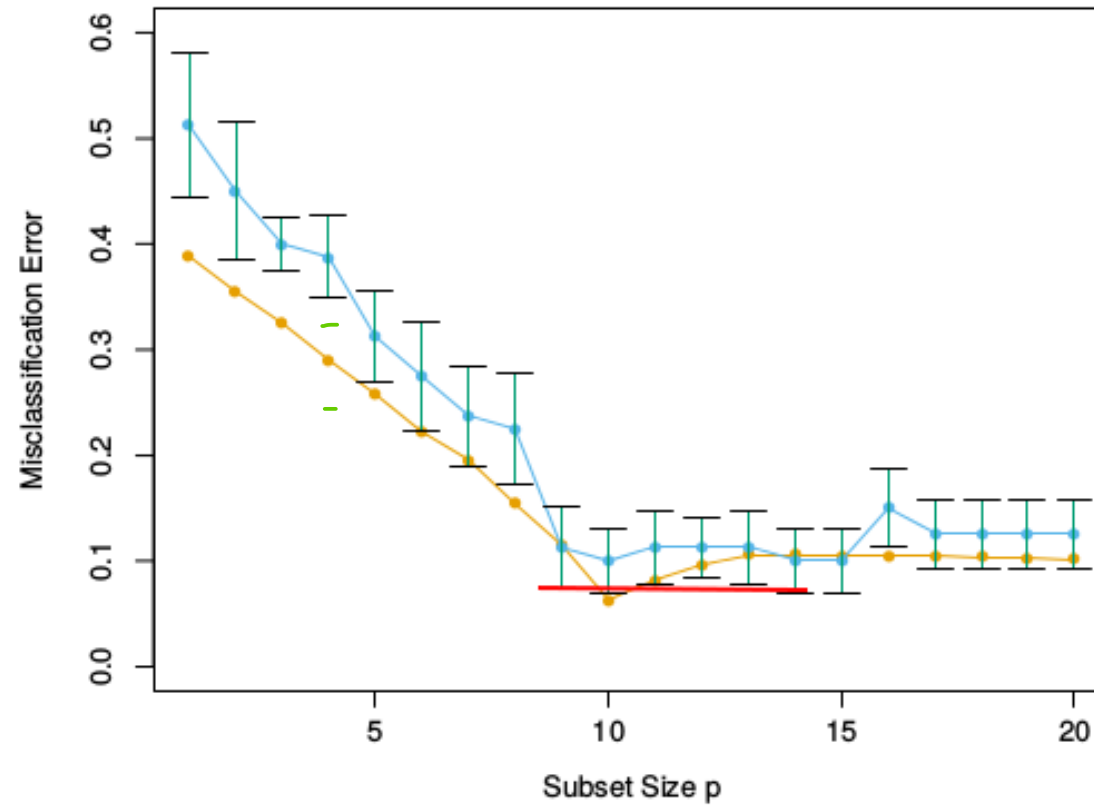
- Makes the best use of the data
- Involve no random sampling

- Disadvantages:

- Took long time to run, computationally expensive



# K-fold cross validation (cont.)

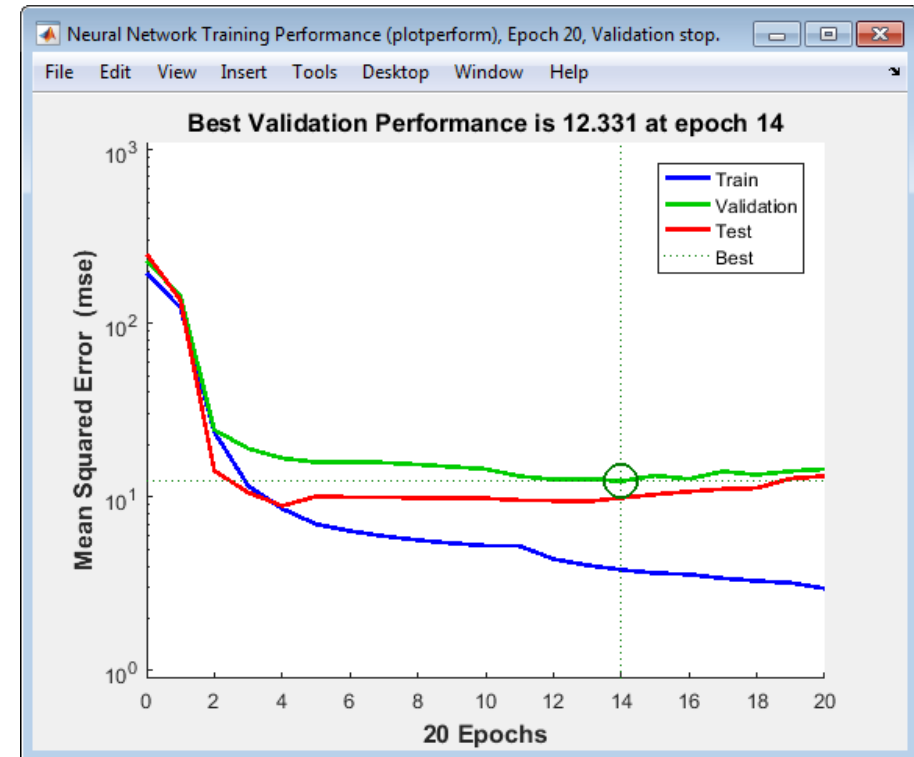




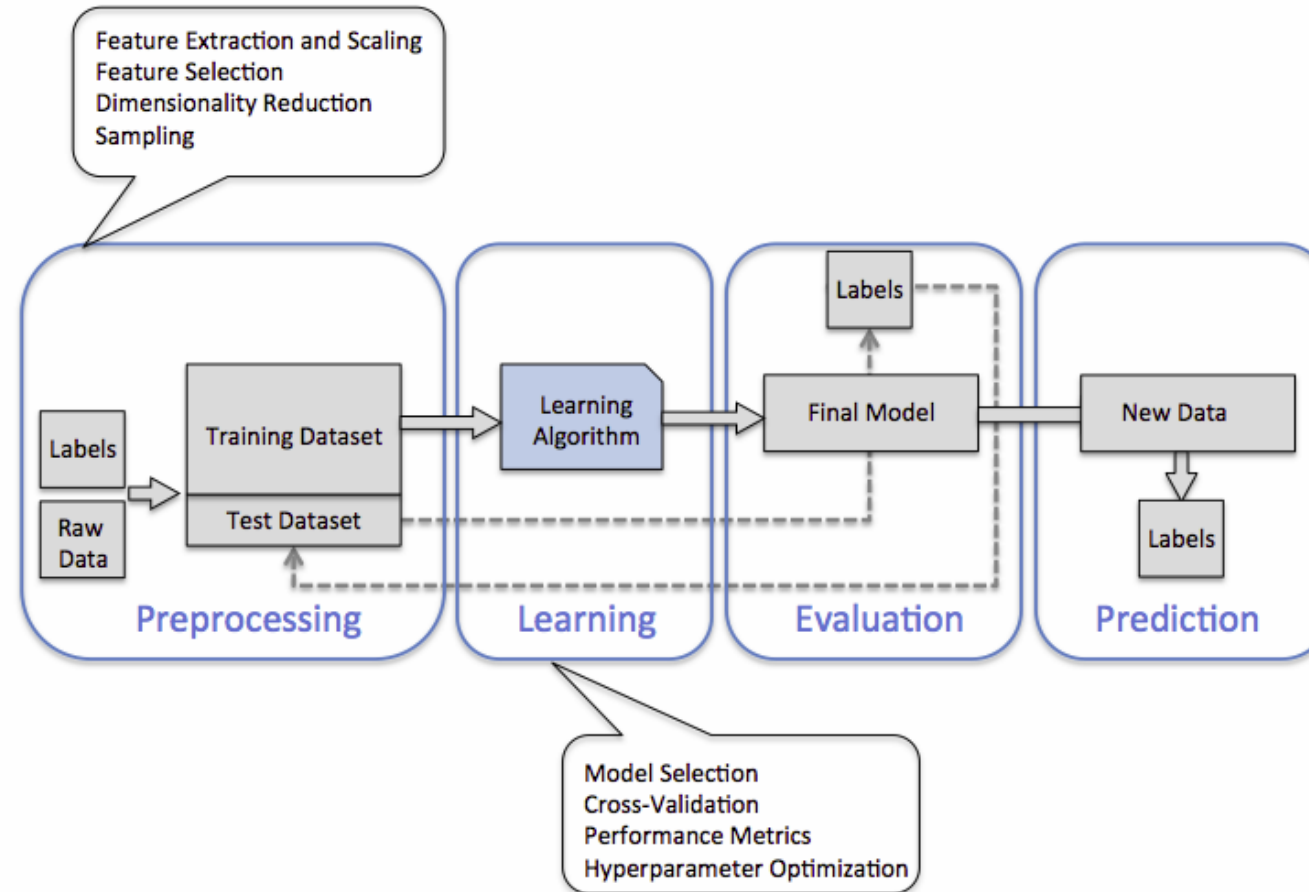
# Three-way sampling method



Validation set is use for tuning parameters of your Model.



# Recall: Road map to Data Mining/ Machine learning



# Intro to ML

- Previous, “Intelligent” applications, many systems used hand coded rules of “if ” and “else” decisions to process data or adjust to user input.
  - Changing the task even slightly might require a rewrite of the whole system.
  - Designing rules requires a deep understanding of how a decision should be made by a human expert.

# AI vs. machine learning vs. deep learning

	AI	Machine learning	Deep learning
Optimal data volumes	Varying data volumes	Thousands of data points	Big data: millions of data points
Outputs	Anything from predictions to recommendations to decision-making	Numerical value, like a classification or score	Anything from numerical values to free-form elements, like free text and sound
How it works	Machines are programmed to mimic human activity with human-like accuracy	Uses various types of automated algorithms that learn to model functions and predict future actions from data	Uses neural networks that pass data through many processing layers to interpret data features and relationships
How it's managed	Algorithms require human oversight in order to function properly	Algorithms are directed by data analysts to examine specific variables in data sets	Algorithms are largely self-directed on data analysis once they're put into production

<b>Machine learning algorithm</b>	<b>Data processing tasks</b>	<b>Section</b>	<b>Representative references</b>
K-Nearest Neighbors	Classification	5.1.1	[58] [59]
Naive Bayes	Classification	5.1.2	[60] [61]
Support Vector Machine	Classification	5.1.3	[62] [63] [64] [65]
Linear Regression	Regression	5.2.1	[66] [66] [67] [68]
Support Vector Regression	Regression	5.2.2	[69] [70]
Classification and Regression Trees	Classification/Regression	5.3.1	[71] [72] [73]
Random Forests	Classification/Regression	5.3.2	[74]
Bagging	Classification/Regression	5.3.3	[75]
K-Means	Clustering	5.4.1	[76] [77] [78]
Density-Based Spatial Clustering of Applications with Noise	Clustering	5.4.2	[79] [80] [81]
Principal Component Analysis	Feature extraction	5.5.1	[82] [83] [84] [85] [86]
Canonical Correlation Analysis	Feature extraction	5.5.2	[87] [88]
Feed Forward Neural Network	Regression/Classification/ Clustering/Feature extraction	5.6.1	[89] [90] [91] [92] [93] [57]
One-class Support Vector Machines	Anomaly detection	5.8.1	[94] [95]

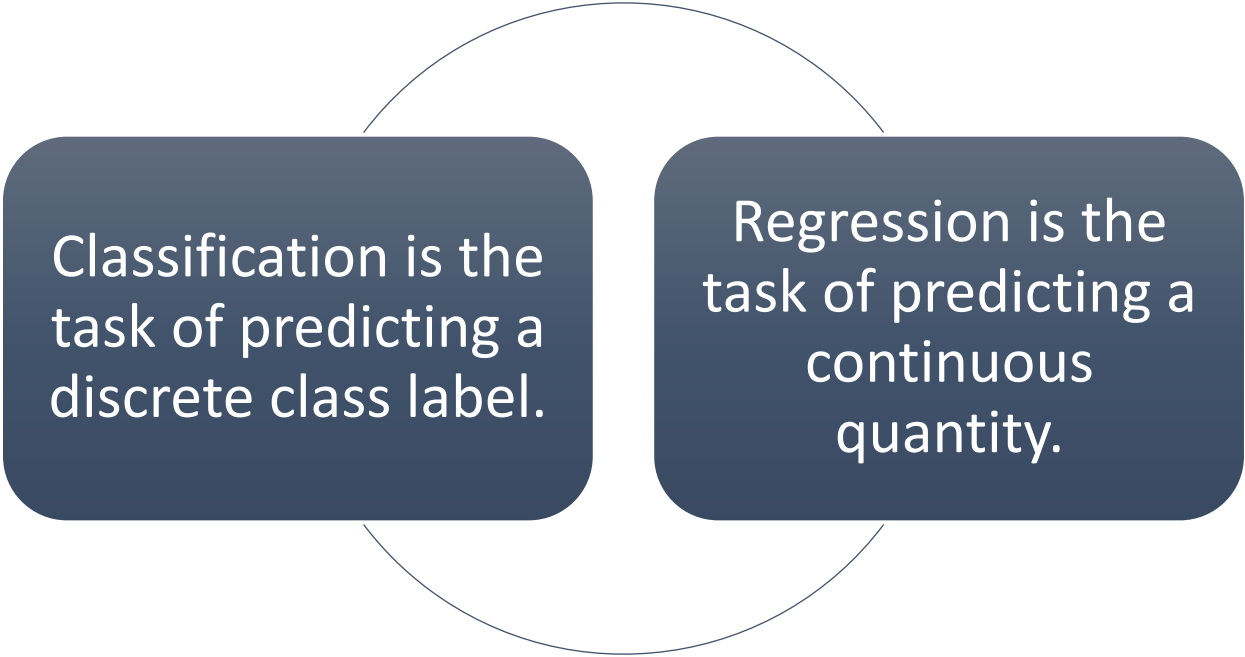
(Mahdavinejad et al., 2018)

The real challenge in using ML is to find the algorithm whose learning bias is the best match for a particular data set.

# Predictive Model (in short)

- Construct a model from historical data to make a prediction on unseen data (e.g. we don't know the answer)
- The job of machine learning functions is to find the optimal mapping function.
- There are two ways
  - Classification
  - Regression

# Classification vs. Regression



Classification is the task of predicting a discrete class label.

Regression is the task of predicting a continuous quantity.



# Classification vs. Regression

## However,

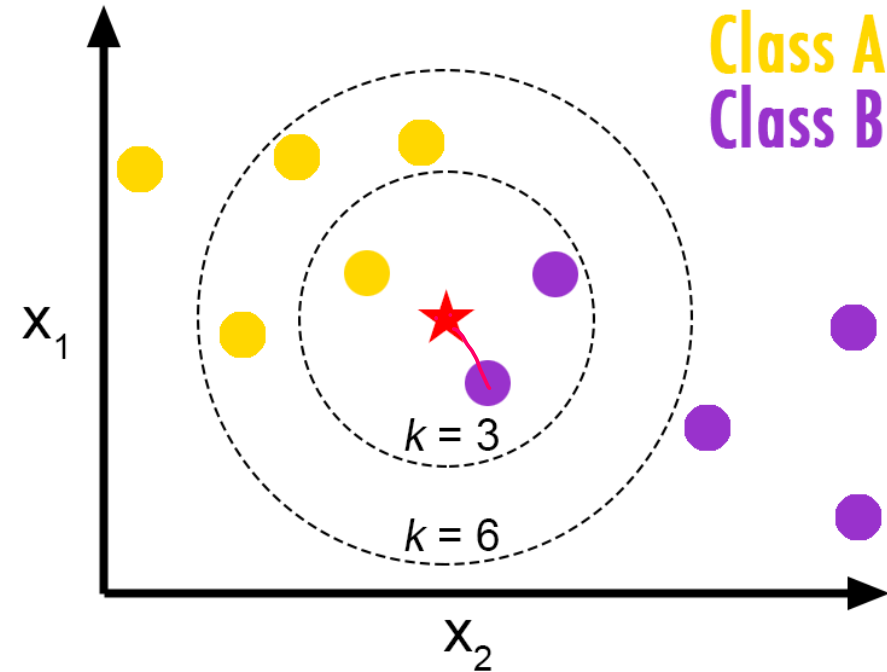
- A classification algorithm **may/can** predict a continuous value, **but** the continuous value is in the form of a probability for a class label.
- A regression algorithm **may/can** predict a discrete value, **but** the discrete value in the form of an integer quantity.

# K-Nearest Neighbors

- One of the **most simplest** classification algorithm
- Its belonged to **supervised learning** algorithm.
- One of the most used and produce good performance
- Use similarity measurements
- Can be use in **regression** and **classification** problems
- Non-parametric (don't need to follow specific distribution of data)

# K-Nearest Neighbors (cont.)

The object is assigned a class of its nearest neighbor

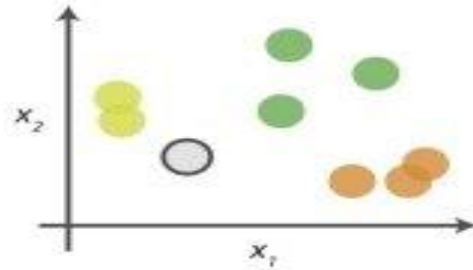


# K-Nearest Neighbors (cont.)

1. Compute a distance value between the item to be classified and every item in the training data-set
2. Pick the  $k$  closest data points (the items with the  $k$  lowest distances)
3. Conduct a “majority vote” among those data points — the dominating classification in that pool is decided as the final classification

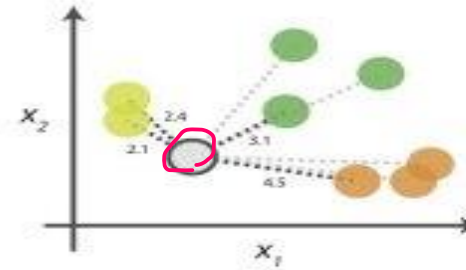
# kNN Algorithm

## 0. Look at the data











Say you want to classify the grey point into a class. Here, there are three potential classes - lime green, green and orange.

## 1. Calculate distances









Start by calculating the distances between the grey point and all other points.

## 2. Find neighbours

Point Distance		
 ... 	2.1	→ 1st NN
 ... 	2.4	→ 2nd NN
 ... 	3.1	→ 3rd NN
 ... 	4.5	→ 4th NN

Next, find the nearest neighbours by ranking points by increasing distance. The nearest neighbours (NNs) of the grey point are the ones closest in dataspace.

## 3. Vote on labels

Class	# of votes	
	2	→ Class  wins the vote! Point  is therefore predicted to be of class  .
	1	
	1	

Vote on the predicted class labels based on the classes of the k nearest neighbours. Here, the labels were predicted based on the k=3 nearest neighbours.

Image source from Kdnugget websites

# How to get the k value?

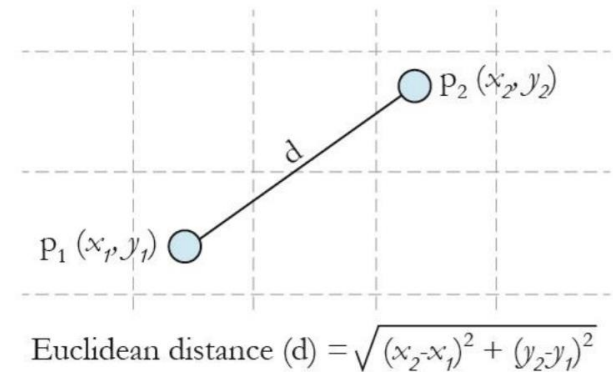
- Euclidean distance

$$E(x, y) = \sqrt{\sum_{i=0}^n (x_i - y_i)^2}$$

- Cosine similarity

Measurement between two non-zero vectors

$$\text{similarity} = \cos(\theta) = \frac{A \cdot B}{\|A\| \|B\|} = \frac{\sum_{i=1}^n A_i B_i}{\sqrt{\sum_{i=1}^n A_i^2} \sqrt{\sum_{i=1}^n B_i^2}}$$



# K-Nearest Neighbors Algorithm Impl.(cont.)

1. **Load** the data
2. Initialize the value of ***k***
3. To get the predicted class, iterate from 1 to total number of training data points
  1. **Calculate the distance** to all training data. Here, we will use Euclidean distance as our distance metric since it's the most popular method.
  2. **Sort** the calculated distances in **ascending order** based on distance values
  3. Get **top *k* rows** from the sorted array
  4. Get **the most frequent class (Majority vote)** of these rows
  5. Return the predicted class

# Scikit-learn library for python

- **Install Scikit-learn from (<http://scikit-learn.org/>)**
- **from sklearn.neighbors import KNeighborsClassifier**
- **from sklearn.metrics import accuracy\_sc**
- **from sklearn.cross\_validation import train\_test\_split**

<https://scikit-learn.org/stable/modules/generated/sklearn.neighbors.KNeighborsClassifier.html>



# KNN from scratch

```
1 # a Counter is a collection where elements are stored as dictionary keys,  
2 # and the key's counts are stored as dictionary values. The example below illustrates this.  
3 from collections import Counter  
4  
5 import numpy as np  
6  
7  
8 def euclidean_distance(x1, x2):  
9     return np.sqrt(np.sum((x1 - x2) ** 2))  
10  
11 class KNN:  
12     def __init__(self, k=3):  
13         self.k = k  
14  
15     def fit(self, X, y):  
16         self.X_train = X  
17         self.y_train = y  
18  
19     def predict(self, X):  
20         y_pred = [self._predict(x) for x in X]  
21         return np.array(y_pred)  
22  
23     def _predict(self, x):  
24         # Compute distances between x and all examples in the training set  
25         distances = [euclidean_distance(x, x_train) for x_train in self.X_train]  
26  
27         # Sort by distance and return indices of the first k neighbors  
28         k_idx = np.argsort(distances)[: self.k]  
29  
30         # Extract the labels of the k nearest neighbor training samples  
31         k_neighbor_labels = [self.y_train[i] for i in k_idx]  
32  
33         # return the most common class label  
34         most_common = Counter(k_neighbor_labels).most_common(1)  
35         return most_common[0][0]
```

```
1 if __name__ == "__main__":  
2  
3     from sklearn import datasets  
4     from sklearn.model_selection import train_test_split  
5  
6  
7     def accuracy(y_true, y_pred):  
8         accuracy = np.sum(y_true == y_pred) / len(y_true)  
9         return accuracy  
10  
11     iris = datasets.load_iris()  
12     X, y = iris.data, iris.target  
13  
14     X_train, X_test, y_train, y_test = train_test_split(  
15         X, y, test_size=0.2, random_state=1234  
16     )  
17  
18     k = 3  
19     clf = KNN(k=k)  
20     clf.fit(X_train, y_train)  
21     predictions = clf.predict(X_test)  
22     print("KNN classification accuracy", accuracy(y_test, predictions))
```

# KNN Demo with sklearn

- [https://github.com/preenet/961701\\_65/blob/main/Classification\\_demo\\_with\\_KNN.ipynb](https://github.com/preenet/961701_65/blob/main/Classification_demo_with_KNN.ipynb)

# Workshop (Term project Proposal(mini))

- Maximum number of 3 people
- Maximum number 3-5 pages
  - NLP related problem
  - (dataset should contain at least one feature that is textual)
- Proposal content
  - Introduction to your problem
  - Activities, output(s), outcome, identify
  - Stakeholders
  - Data sources
  - DS tasks
  - Solution as software
- For master students (Replicate conference paper that is related to NLP to get 100%)

# Agenda

- Term project due a week before final exam.
- One more workshop (ML topic)

# Pros and Cons

- Pros
  - Simple
  - Can apply to both regression and classification
  - No data assumption needs
- Cons
  - Sensitive to scale of data
  - Curse of demission
  - Outlier sensitive
  - Missing value treatment

# More KNN

- How to handle categorical variables?
  - Use dummy variables
- Optimal k value?
  - Perform model selection

# KNN for regression

- KneighborsRegressor
- Use the same distance function like the classifier version

Mean absolute error :

Overview

Formula

---

Formula

$$\text{MAE} = \frac{\sum_{i=1}^n |y_i - x_i|}{n}$$

**MAE** = mean absolute error

$y_i$  = prediction

$x_i$  = true value

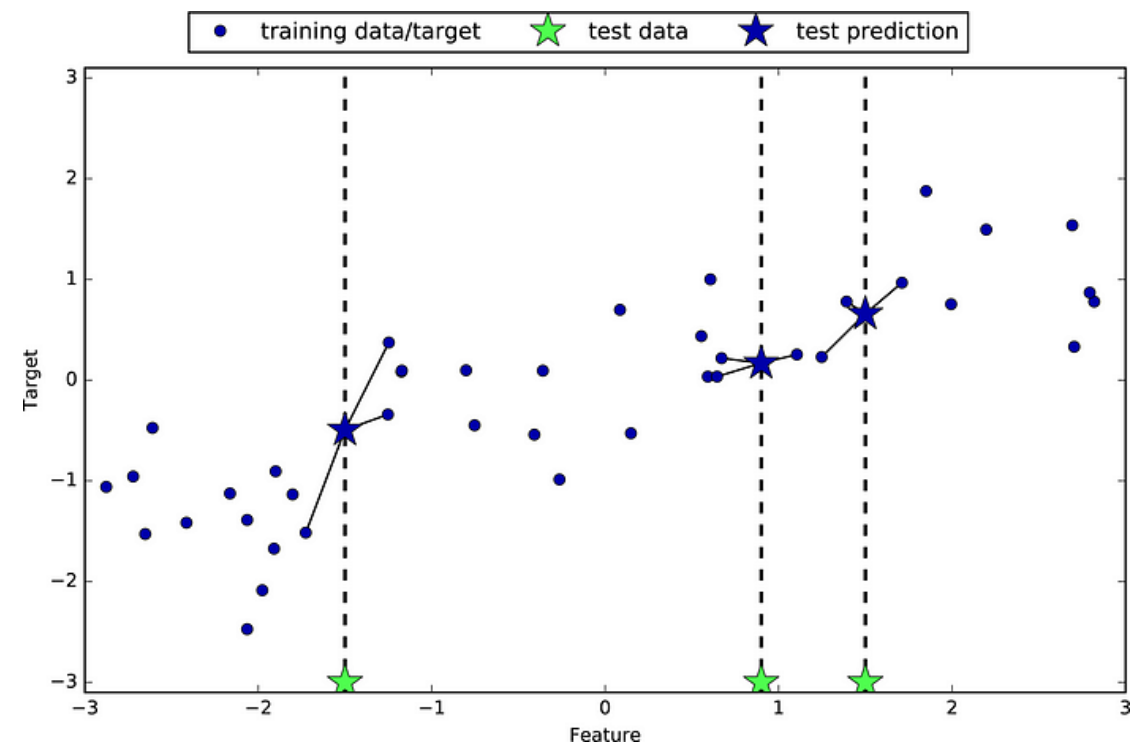
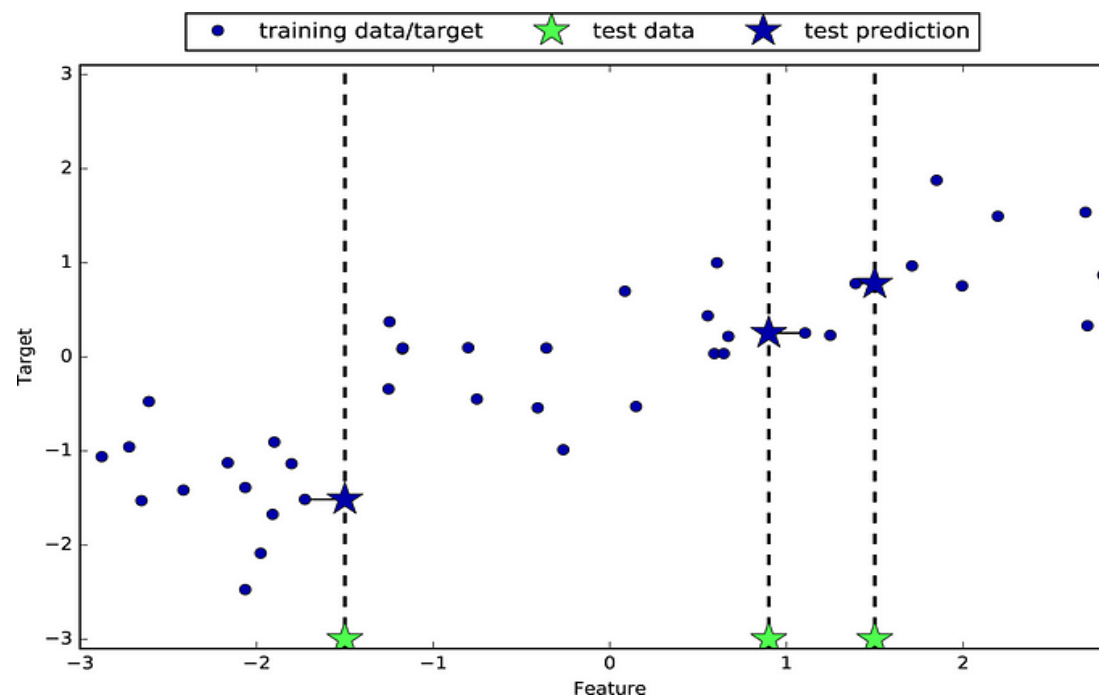
$n$  = total number of data points

# KNN regressor example

```
1  # Import the necessary libraries
2  import numpy as np
3  from sklearn.model_selection import train_test_split
4  from sklearn.neighbors import KNeighborsRegressor
5  from sklearn.metrics import mean_squared_error, r2_score
6
7  # Generate some sample data
8  X = np.random.rand(100, 1) # Feature
9  y = 2 * X + np.random.randn(100, 1) # Target
10
11 # Split the data into training and testing sets
12 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)
13
14 # Create a KNN regressor with a specified number of neighbors (e.g., 5)
15 k = 5
16 knn_regressor = KNeighborsRegressor(n_neighbors=k)
17
18 # Fit the model to the training data
19 knn_regressor.fit(X_train, y_train)
20
21 # Make predictions on the test data
22 y_pred = knn_regressor.predict(X_test)
23
24 # Evaluate the model
25 mse = mean_squared_error(y_test, y_pred)
26 r2 = r2_score(y_test, y_pred)
27
28 print(f"Mean Squared Error: {mse}")
29 print(f"R-squared: {r2}")
```

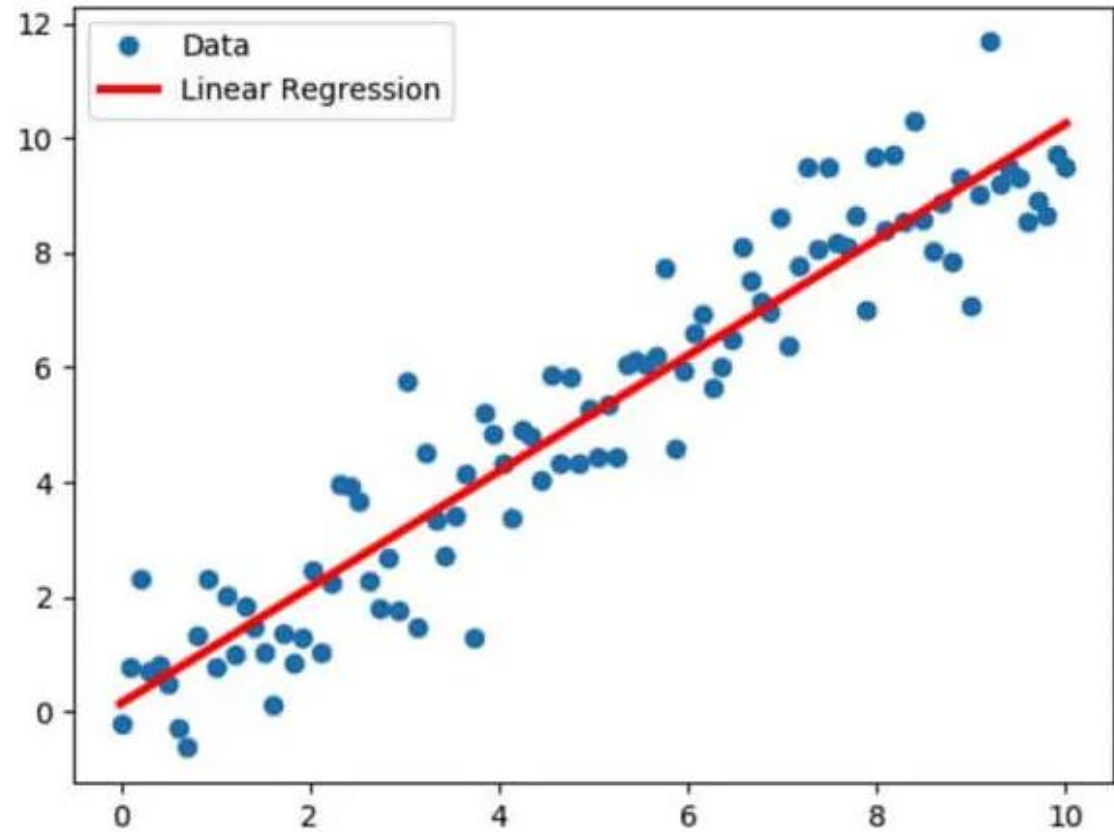


# KNN regressor



# Linear regression

- When we need model interpretability
- When the relationship tends to be linear



# Regression introduction

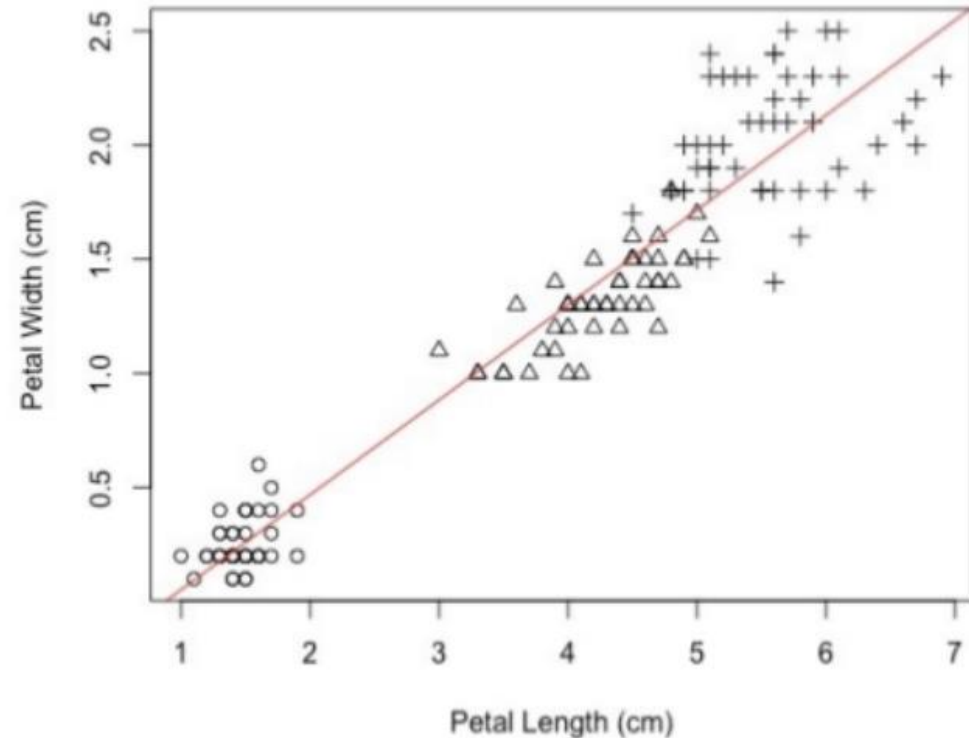
- Regression – a task of approximating a mapping function ( $f$ ) from input variables ( $X$ ) to a continuous output variable ( $y$ ).
- A continuous output variable real number( $R$ ), (e.g. int or float values.)
  - These are often quantities (e.g. amounts and sizes.)
- A problem with multiple input variables aka. multivariate regression problem.

# Regression introduction

- Classification predictions can be evaluated using accuracy, whereas regression predictions cannot.
- $\text{Accuracy} = \text{correct\_pred} / \text{total\_pred} \times 100$
- Regression predictions can be evaluated using root mean squared error, whereas classification predictions cannot.
- $\text{RMSE} = \text{sqrt}(\text{avg}(\text{error}^2))$

# Linear Regression for Iris dataset

- We can use linear regression as machine learning algorithm to predict petal\_width given petal\_length.



# Least Squares - Linear Regression

- Linear regression function
  - $\hat{y} = b_0 + b_1x$  or you may know from basic math such as  $(y = mx+b)$
- Slope or Gradient (how steep the line is)
  - $b_1 = \frac{\sum(x-\bar{x})(y-\bar{y})}{\sum(x-\bar{x})^2}$
- Y-intercept (where the line cut the Y-axis)
  - $b_0 = \bar{y} - b_1\bar{x}$

# Linear regression- intro

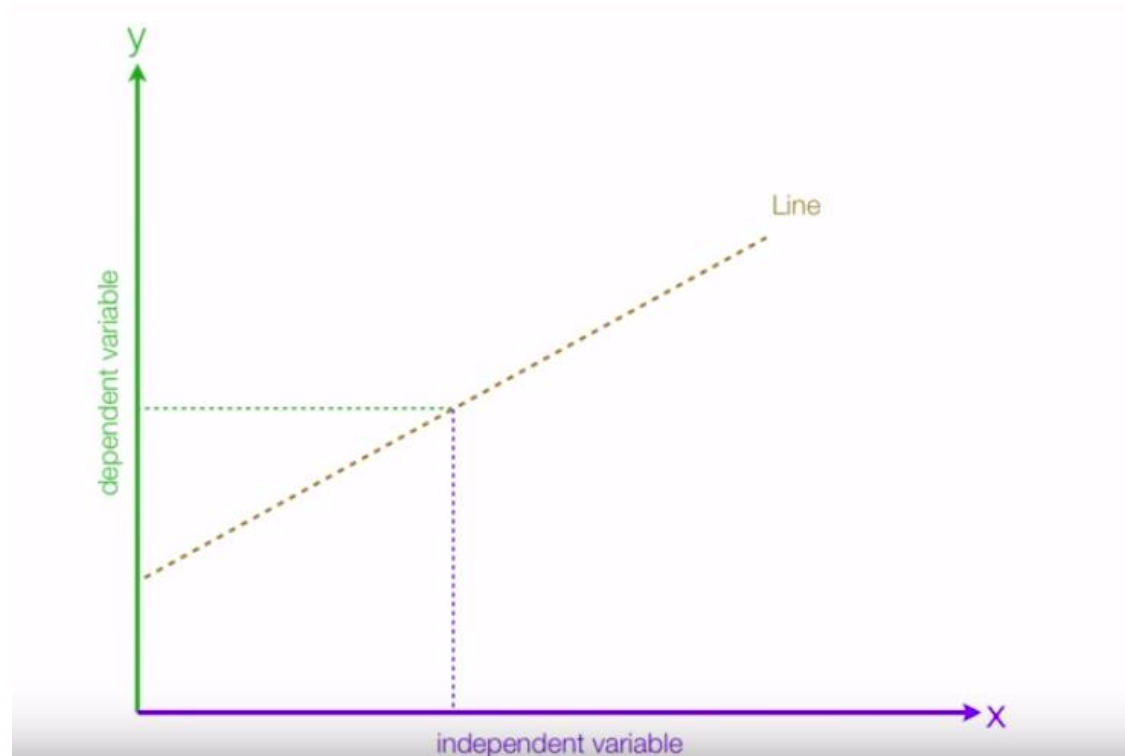
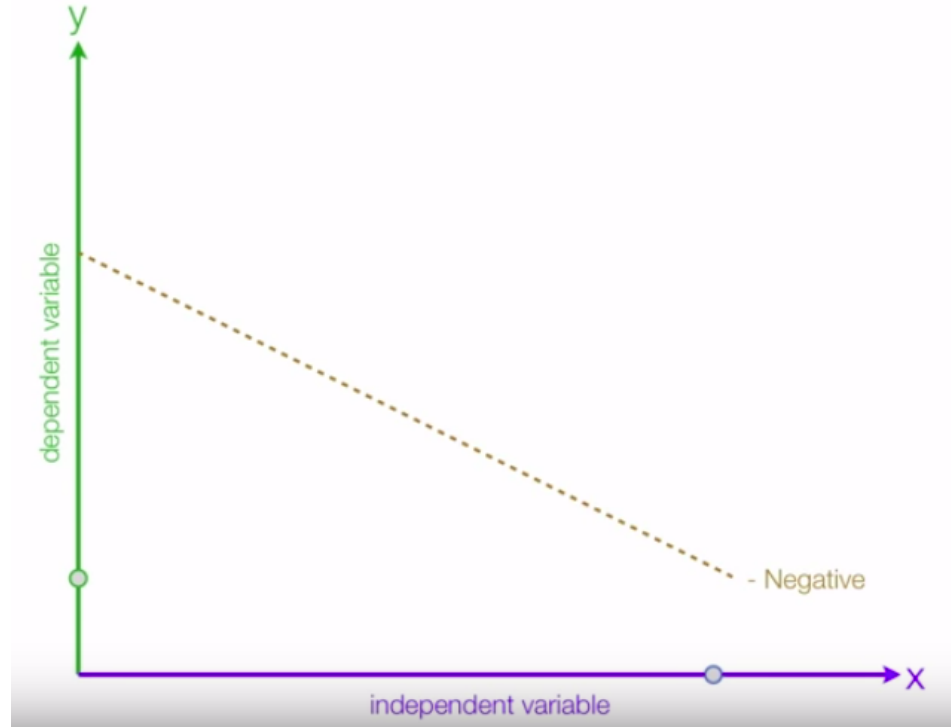
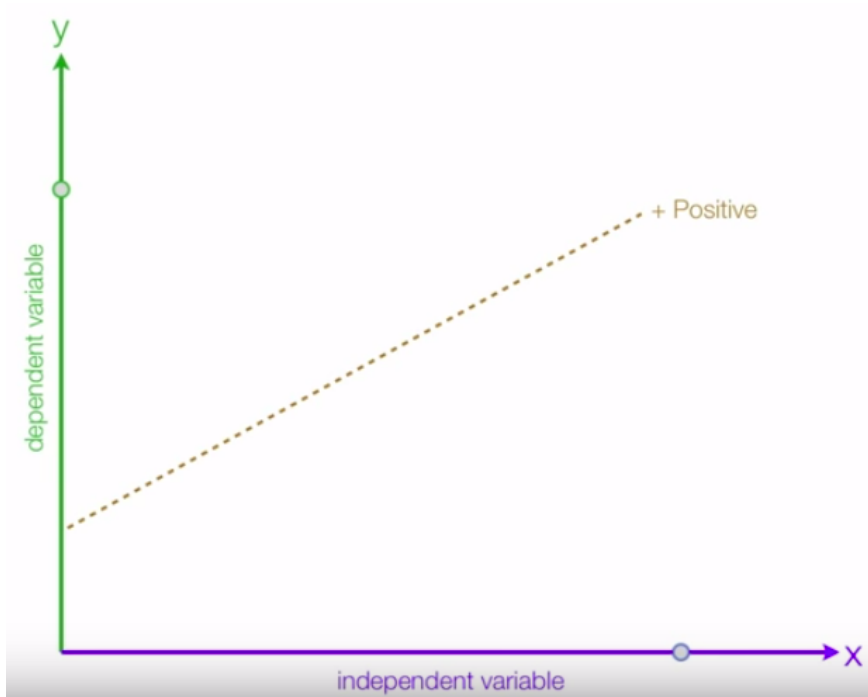


Image source: David Longstreet, Professor of the Universe, MyBookSucks, (Accessed, August 2018)

# Linear regression- intro2

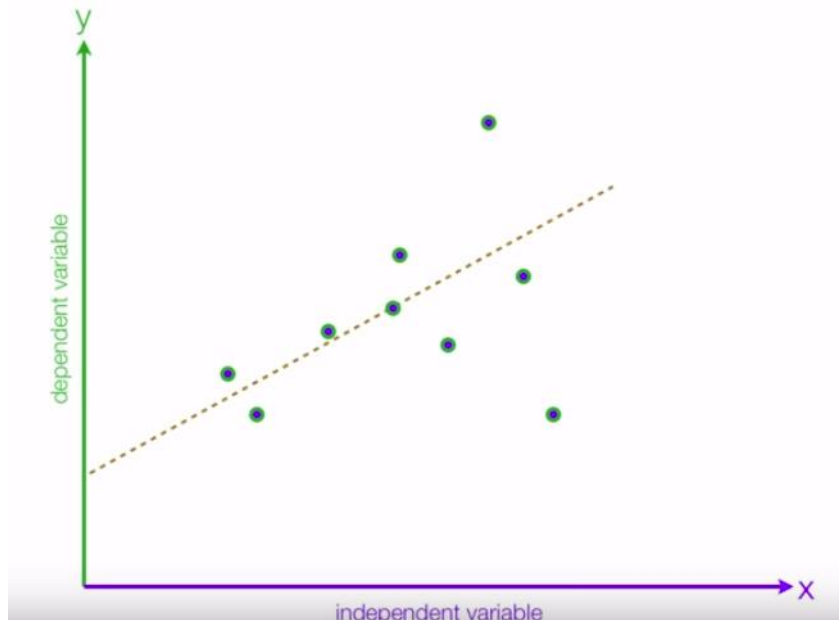


*Image source: David Longstreet, Professor of the Universe, MyBookSucks, (Accessed, August 2018)*



# Linear regression- intro3

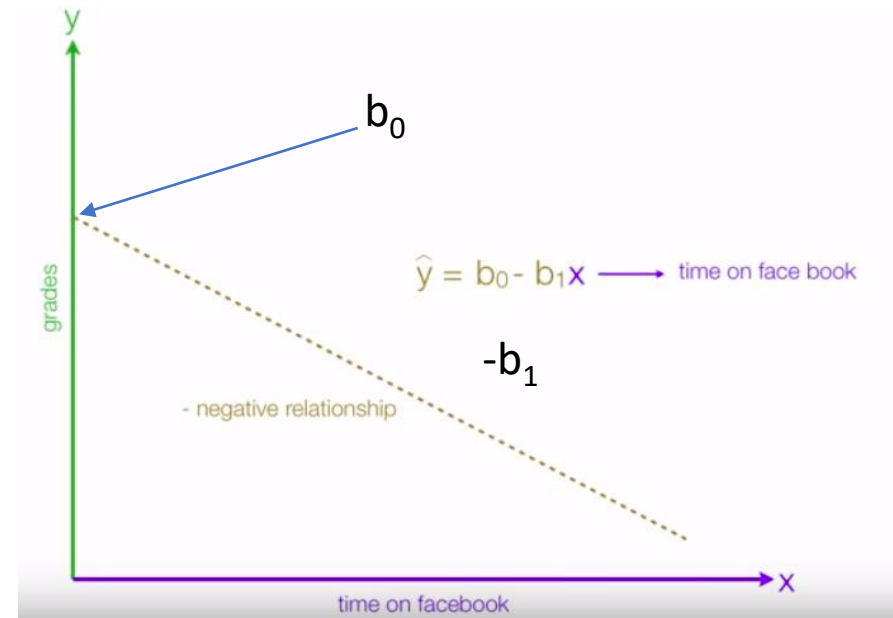
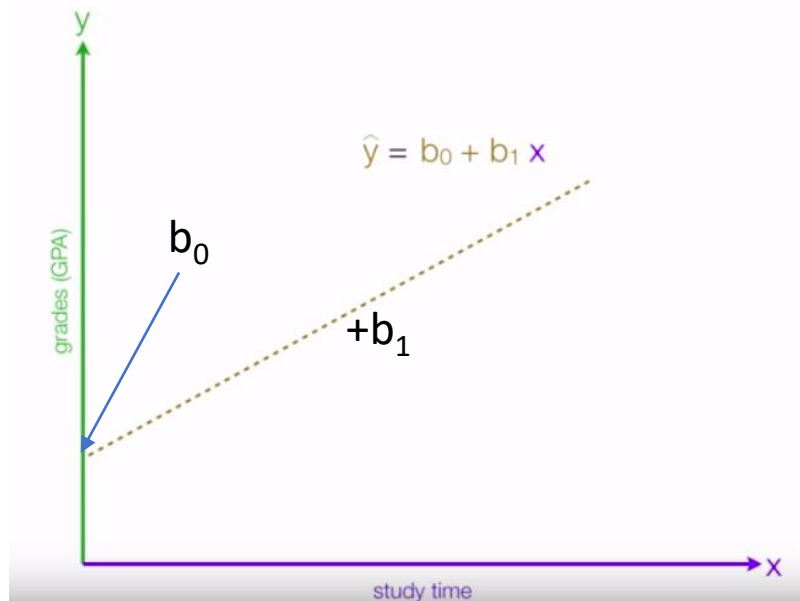
- We try to find the line that can fit the observation points  
the line is called regression line
- We try to draw the line that minimize the errors



*Image source: David Longstreet, Professor of the Universe, MyBookSucks, (Accessed, August 2018)*

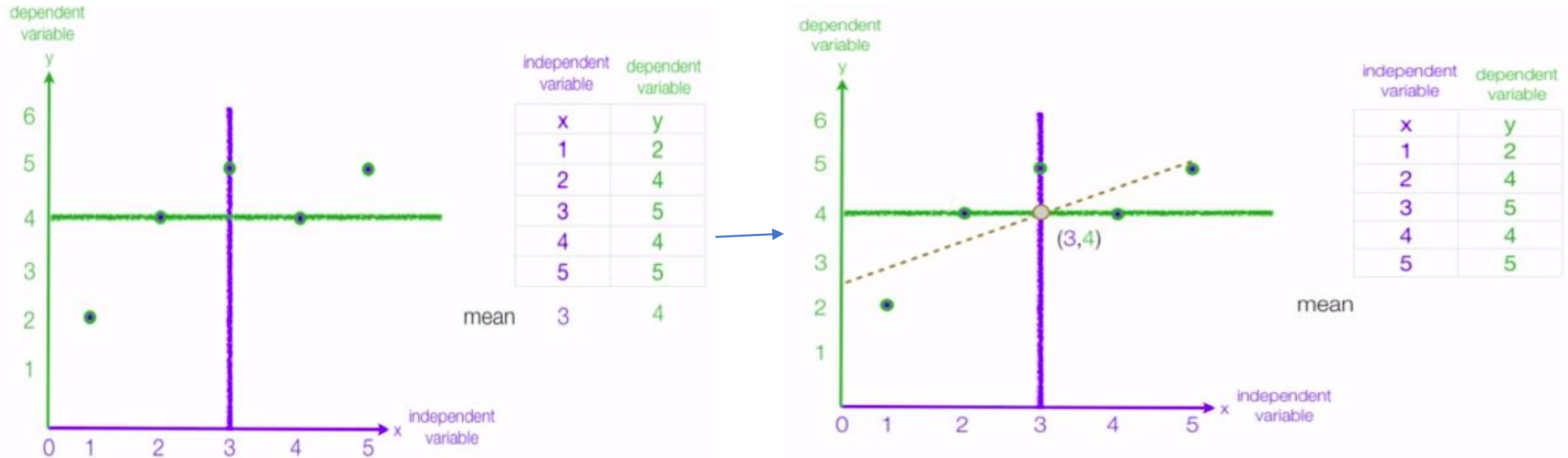
# Linear regression- intro4

- $\hat{y}$  is the estimated value (dependent variable or output)
- $x$  is the independent value (input, we can control)
- $b_0$  is the y-intercept
- $b_1$  is the slop of the regression line



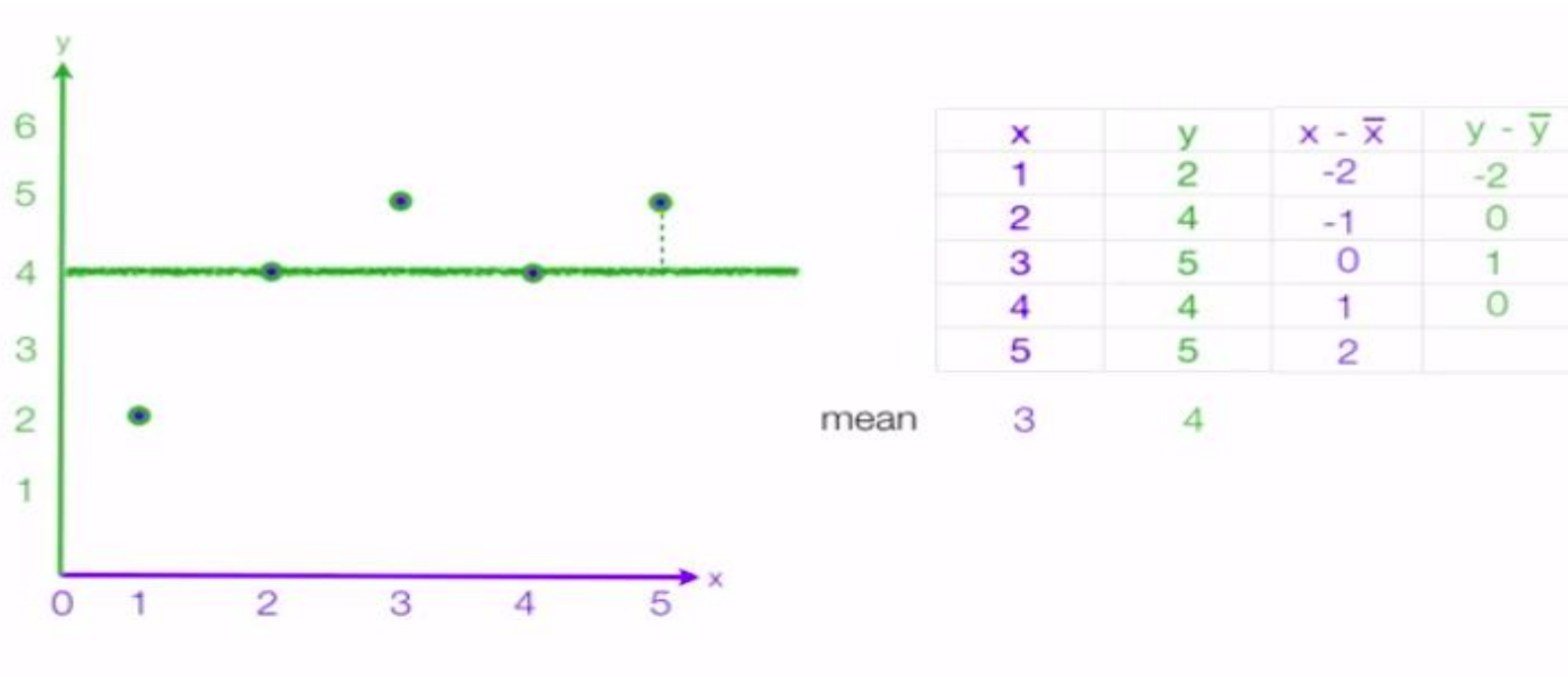
*Image source: David Longstreet, Professor of the Universe, MyBookSucks, (Accessed, August 2018)*

# Regression using Least Square method



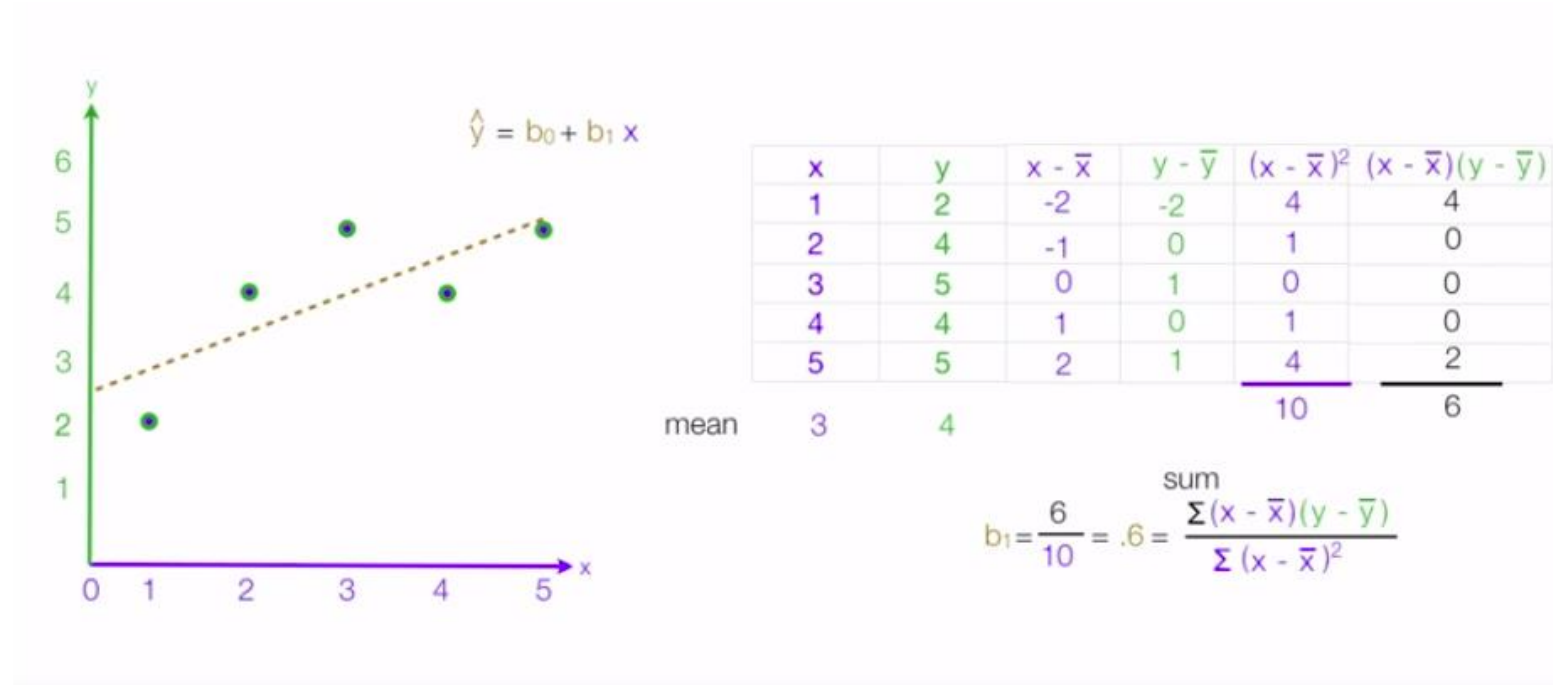
- Find  $\bar{x}$  and  $\bar{y}$ , then plot and draw the line pass through the point
- All the possible regression line has to go through the point

# Regression using Least Square method (cont.)



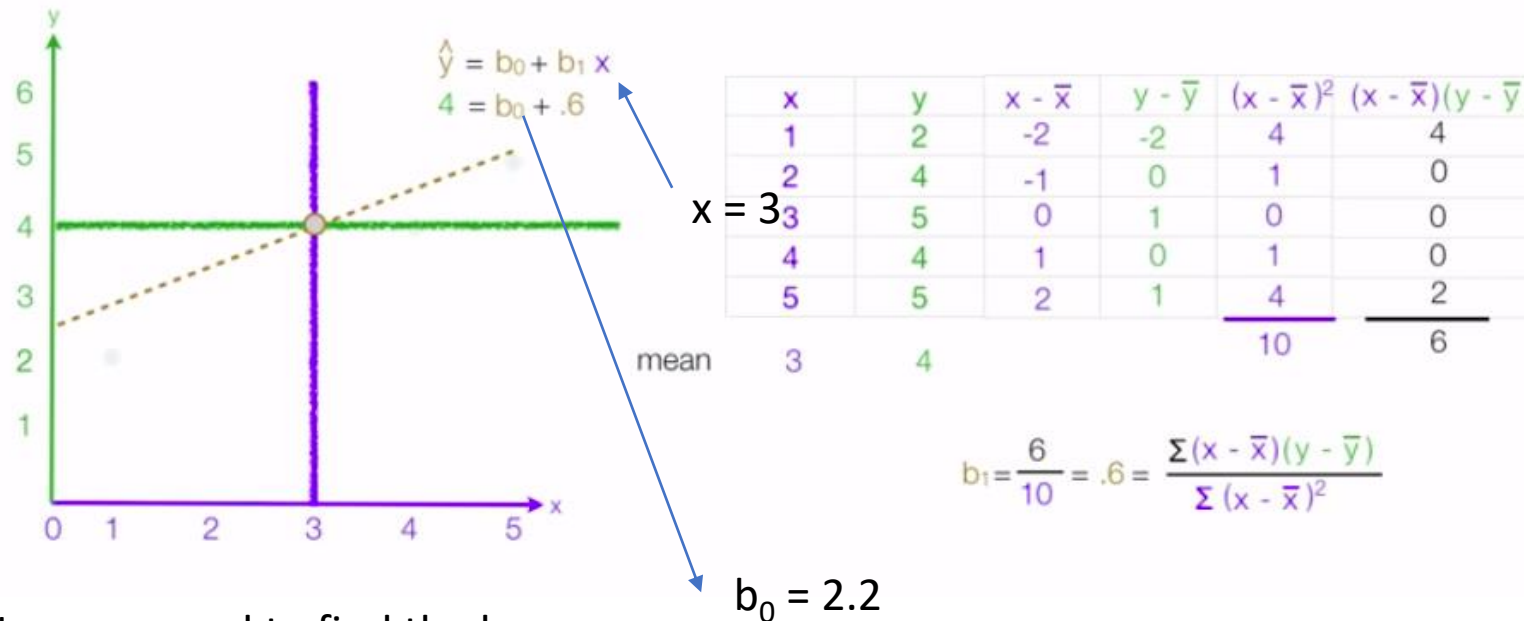
- Find the distance of observation points of x and y from the mean

# Regression using Least Square method (cont.)



- Now, we want to find the slope of the line by compute the square of  $(x - \bar{x})$  and  $(x - \bar{x})(y - \bar{y})$ ,
- Hence, we get  $b_1 = 0.6$

# Regression using Least Square method (cont.)

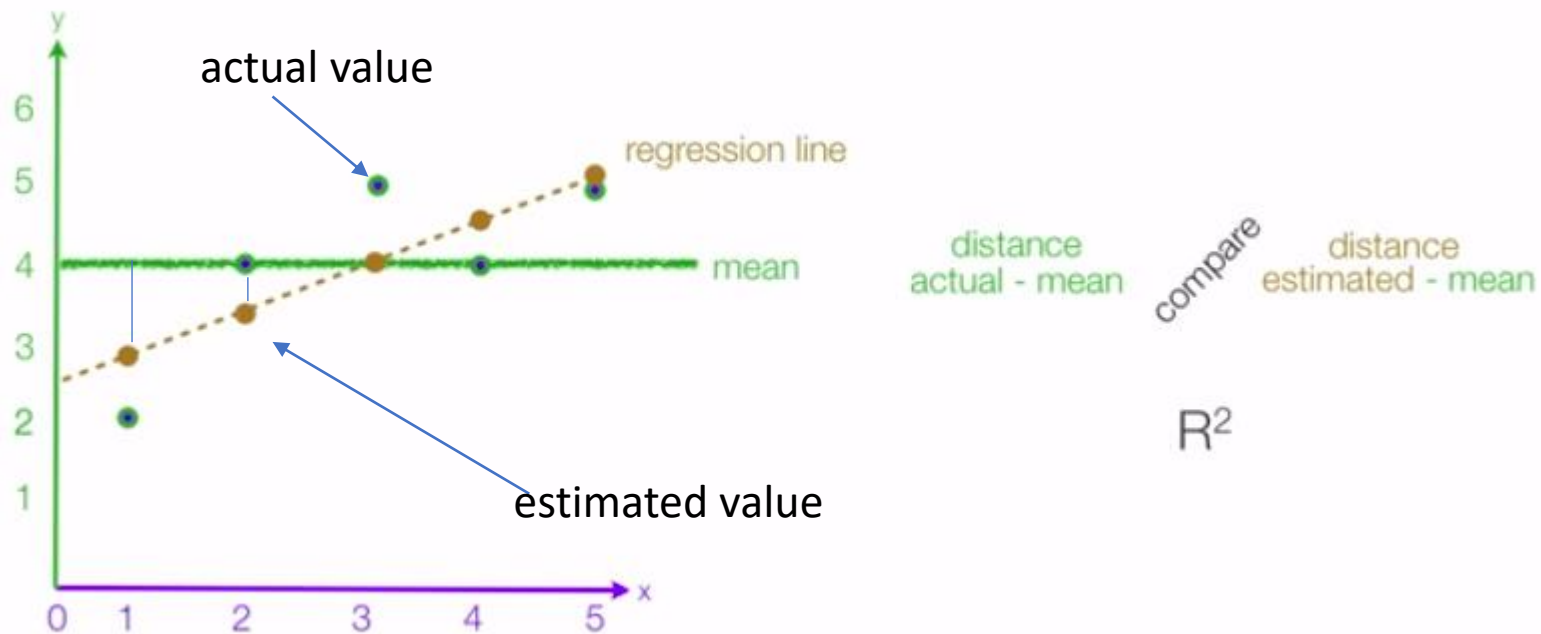


Now, we need to find the  $b_0$

- We know that the regression line will always cross the  $p(3,4)$ , so finding y-intercept is simple.
- Hence, we have  $b_0 = 2.2$  and  $b_1 = 0.6$
- $\hat{y} = 2.2 + 0.6x$

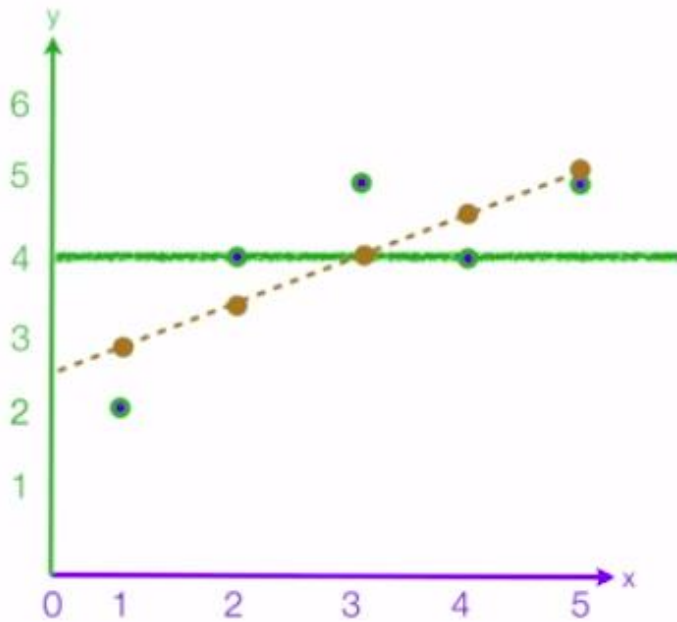
Image source: David Longstreet, Professor of the Universe, MyBookSucks, (Accessed, August 2018)

# Evaluate regression model with R square



- $R^2$  = (R square)
- $R^2$  tells us how well our regression line can estimate predict value.

# Evaluate regression model with R square



x	y	$y - \bar{y}$	$(y - \bar{y})^2$	$\hat{y}$	$\hat{y} - \bar{y}$	$(\hat{y} - \bar{y})^2$
1	2	-2	4	2.8	-1.2	1.44
2	4	0	0	3.4	-.6	.36
3	5	1	1	4	0	0
4	4	0	0	4.6	.6	.36
5	5	1	1	5.2	1.2	1.44
mean		4	6			3.6

$$R^2 = \frac{3.6}{6} = .6$$

$R^2$  ranges [0,1], the more closer 1 the better our regression line



# Steps

- Find  $x^2$  and  $xy$
- Find slope
  - $\frac{\sum(x-\bar{x})(y-\bar{y})}{\sum(x-\bar{x})^2}$
- Get the y-intercept
- Create the equation

# Multicollinearity

- two or more **independent variables** in a multiple regression model are highly correlated with each other.
- Lead to overfitting

# Multicollinearity

- How to check for Multicollinearity

- Use Variance Inflation Factor (VIF)

```
from statsmodels.stats.outliers_influence import  
variance_inflation_factor  
vif = pd.DataFrame()  
vif["features"] = house_selected.columns  
vif["vif_Factor"] =  
[variance_inflation_factor(house_selected.values, i)  
for i in range(house_selected.shape[1])]  
vif
```

	features	vif_value
0	OverallQual	50.558204
1	YearBuilt	7015.226148
2	TotalBsmntSF	23.974354
3	1stFlrSF	700.064200
4	2ndFlrSF	141.592671
5	GrLivArea	1141.861313
6	PoolArea	1.046374
7	MoSold	6.529697
8	YrSold	6542.476374

# Multicollinearity

- How to remedy?
  - Collect more data (Very expensive, but most effective)
  - Feature selection, transformation, PCA
  - Try ridge regression

# Ridge regression

- Aka (L2 regularization), an alpha value to be tuned.
- Present regularization term which is not included in the linear regression
- Use when you have **large number of feature** and **multicollinearity** is happened
- **Try with OLS first as based line and move to more complex like ridge**

# Ridge regression

- Loss function = OLS + alpha \* summation (squared coefficient values)

$$\text{Ridge} = \frac{1}{N} \sum_{i=1}^N (y_i - (mx_i + z))^2 + \lambda \sum_{i=1}^p (mx_i + z)^2$$

# Lasso regression

- Least Absolute Shrinkage Selector operator
- Aka. L1 regularization
- For the penalty term, take abs instead of square
- Can be used a feature selection

# Discussion and Class activities



# Sklearn for linear regressssion

```
1  import numpy as np
2  from sklearn import datasets
3  from sklearn.model_selection import train_test_split
4  from sklearn.linear_model import LinearRegression
5  import matplotlib.pyplot as plt
6
7  # Load the diabetes dataset
8  diabetes = datasets.load_diabetes()
9
10 # Use only one feature for simplicity
11 X = diabetes.data[:, np.newaxis, 2]
12
13 # Split the data into training/testing sets
14 X_train, X_test, y_train, y_test = train_test_split(X, diabetes.target, test_size=0.2, random_state=42)
15
16 # Create linear regression object
17 regr = LinearRegression()
18
19 # Train the model using the training set
20 regr.fit(X_train, y_train)
21
22 # Make predictions using the testing set
23 y_pred = regr.predict(X_test)
```

# Sklearn for Ridge and how to tune

```
2 ridge_regr = Ridge(alpha=1.0)
3
4 # Train the model using the training set
5 ridge_regr.fit(X_train, y_train)
6
7
8 # Split the data into training/testing sets
9 X_train, X_test, y_train, y_test = train_test_split(X, diabetes.target, test_size=0.2, random_state=42)
10
11 # Define a range of alpha values
12 alphas = np.logspace(-4, 4, 100)
13
14 # Create a RidgeCV regressor object
15 ridge_regr = RidgeCV(alphas=alphas, store_cv_values=True)
16
17 # Train the model using the training set
18 ridge_regr.fit(X_train, y_train)
19
20 # Get the best alpha value
21 best_alpha = ridge_regr.alpha_
22 print(f"Best alpha value: {best_alpha}")
23
24 # Make predictions using the testing set
25 y_pred_ridge = ridge_regr.predict(X_test)
26
27 # Print the coefficients
28 print('Coefficients: \n', ridge_regr.coef_)
29
30 # The mean squared error
31 print('Mean squared error: %.2f' % mean_squared_error(y_test, y_pred_ridge))
32
33 # The coefficient of determination: 1 is perfect prediction
34
```

# Linear regression (by hand)

Number of Chimpanzees		Percent Successful Hunts	↕
0	1	30	
1	2	45	
2	3	51	
3	4	57	
4	5	60	
5	6	65	
6	7	70	
7	8	71	



**CMU**  
CHIANG MAI  
MARATHON

**Run** the legacy  
60 years of CMU



**NOV 26, 2023**  
MARATHON - HALF MARATHON  
MINI MARATHON - FUN RUN - VR RACE

เข้าสู่เว็บไซต์



**CMU**  
CHIANG MAI  
MARATHON

**Run** the legacy  
60 years of CMU



**NOV 26, 2023**  
MARATHON - HALF MARATHON  
MINI MARATHON - FUN RUN - VR RACE

**RUNNING SHIRT  
VR RACE**



Front



Back

ติดต่อสอบถามได้ที่  CMU Marathon



# KNN: Now's it your turn

- Let  $k = 3$
- normalize the data in range.
- Plot a figure to support your result

Height (cms)	Weight (kgs)	Size
158	58	M
158	59	M
158	63	M
160	59	M
160	60	M
163	60	M
163	61	M
160	64	L
163	64	L
165	61	L
165	62	L
165	65	L
168	62	L
168	63	L
168	66	L
170	63	L
170	64	L
170	68	?

# Term project

- Workshop :
- Please fill out the term project roaster under MSTeam

# References

- Mahdavinejad, M. S., Rezvan, M., Barekatain, M., Adibi, P., Barnaghi, P., & Sheth, A. P. (2018). Machine learning for internet of things data analysis: A survey. *Digital Communications and Networks*, 4(3), 161–175. <https://doi.org/10.1016/j.dcan.2017.10.002>
- (*Data Science (The MIT Press Essential Knowledge Series)*): Kelleher, John D., Tierney, Brendan: 0000262535432: Amazon.Com: Books, n.d.)
- Andreas C.Muller and Sarah Guido. 2017. Introduction to machine learning with pyhton
- Busse, C. D. (1978). Do Chimpanzees Hunt Cooperatively? *The American Naturalist*, 112(986), 767–770. <https://doi.org/10.1086/283318>