

Machine Learning (Supervised Learning)



Hello!

I am Eslam Ahmed

I am a software engineer.

You can find me at jeksogsa@gmail.com



Hello!

I am Eman Ehab

I am a ML research engineer.

You can find me at
emanehab.ieee@gmail.com





Machine Learning



Supervised Learning

- Regression
 - Simple Linear Regression
 - Multiple Linear Regression
 - Polynomial Regression
 - Evaluating Model Performance
- Classification
 - Logistic Regression
 - K-Nearest Neighbors (KNN)
 - Naive Bayes
 - SVM
 - Decision Trees
 - Ensemble Methods
 - What is Bagging & Boosting
 - Random Forests
 - XGBoost
 - Evaluating Model Performance



Unsupervised Learning

- Clustering
 - KMeans
 - Hierarchical Clustering
 - Density Based Clustering - DBSCAN
- Association rule mining
 - Apriori
- Dimension Reduction
 - PCA
 - LDA
- Evaluating Model Performance



Model Selection & Evaluation

- Cross Validation
- Hyperparameter Tuning
 - Grid Search
 - Randomized Search



Recommendation Systems



Machine Learning



Supervised Learning

- Regression
 - Simple Linear Regression
 - Multiple Linear Regression
 - Polynomial Regression
 - Evaluating Model Performance
- Classification
 - Logistic Regression
 - K-Nearest Neighbors (KNN)
 - Naive Bayes
 - SVM
 - Decision Trees
 - Ensemble Methods
 - What is Bagging & Boosting
 - Random Forests
 - XGBoost
 - Evaluating Model Performance



Unsupervised Learning

- Clustering
 - KMeans
 - Hierarchical Clustering
 - Density Based Clustering - DBSCAN
- Association rule mining
 - Apriori
- Dimension Reduction
 - PCA
 - LDA
- Evaluating Model Performance



Model Selection & Evaluation

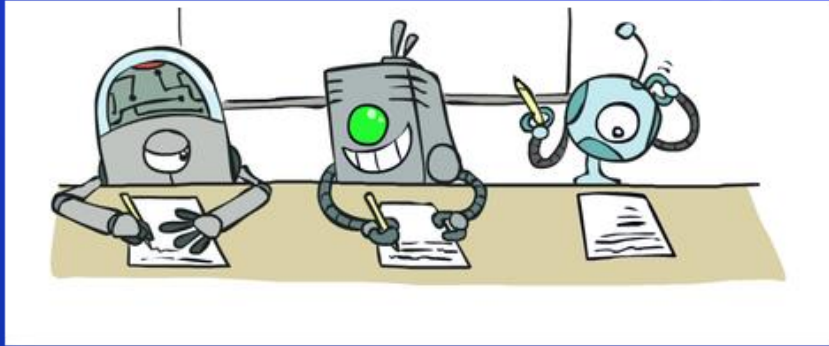
- Cross Validation
- Hyperparameter Tuning
 - Grid Search
 - Randomized Search



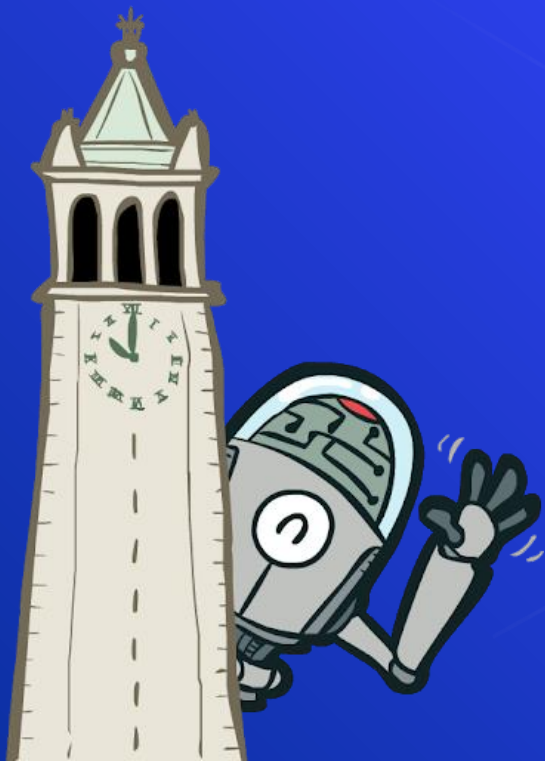
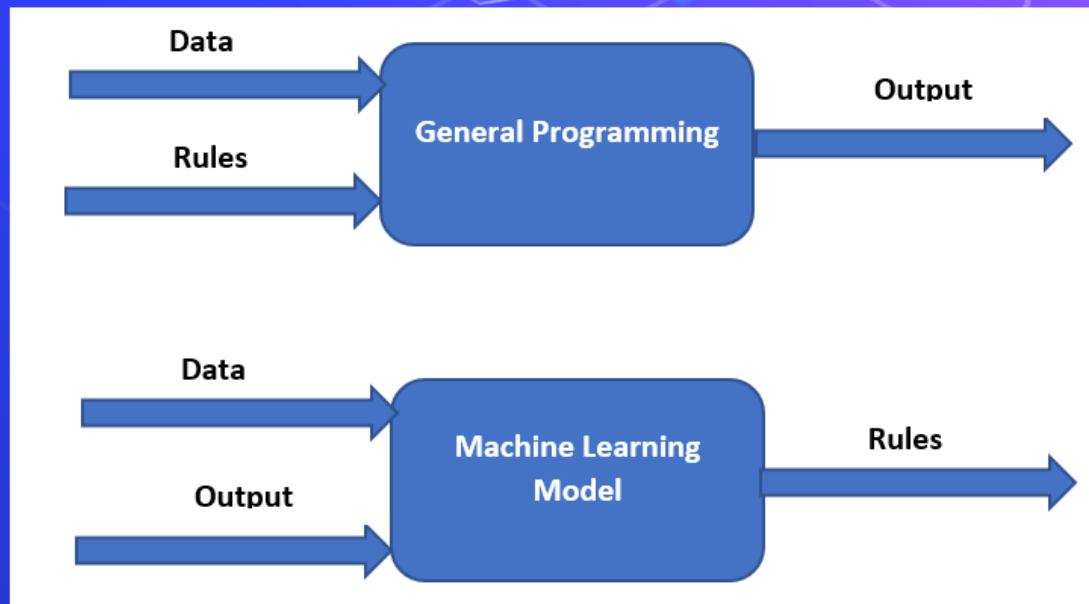
Recommendation Systems

Machine Learning

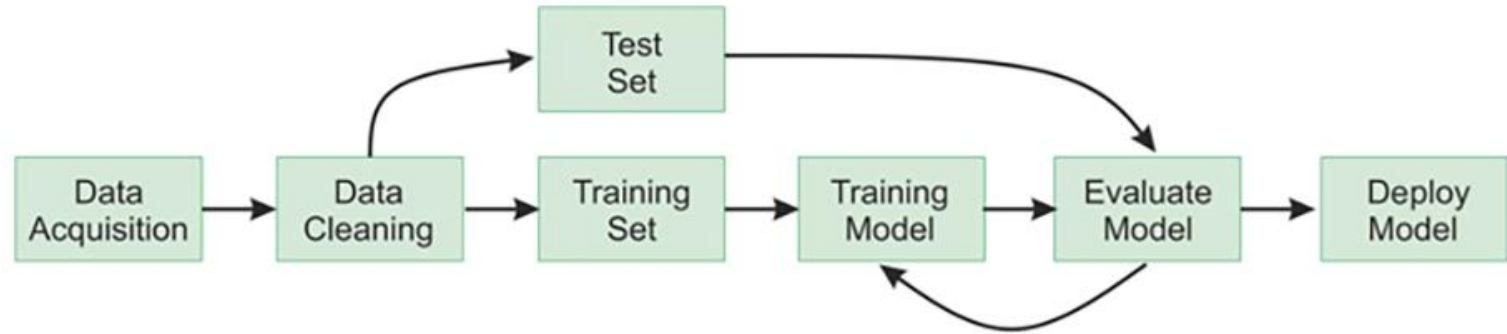
Machine learning involves computers discovering how they can perform tasks without being explicitly programmed to do so. It involves computers learning from data provided so that they carry out certain tasks.



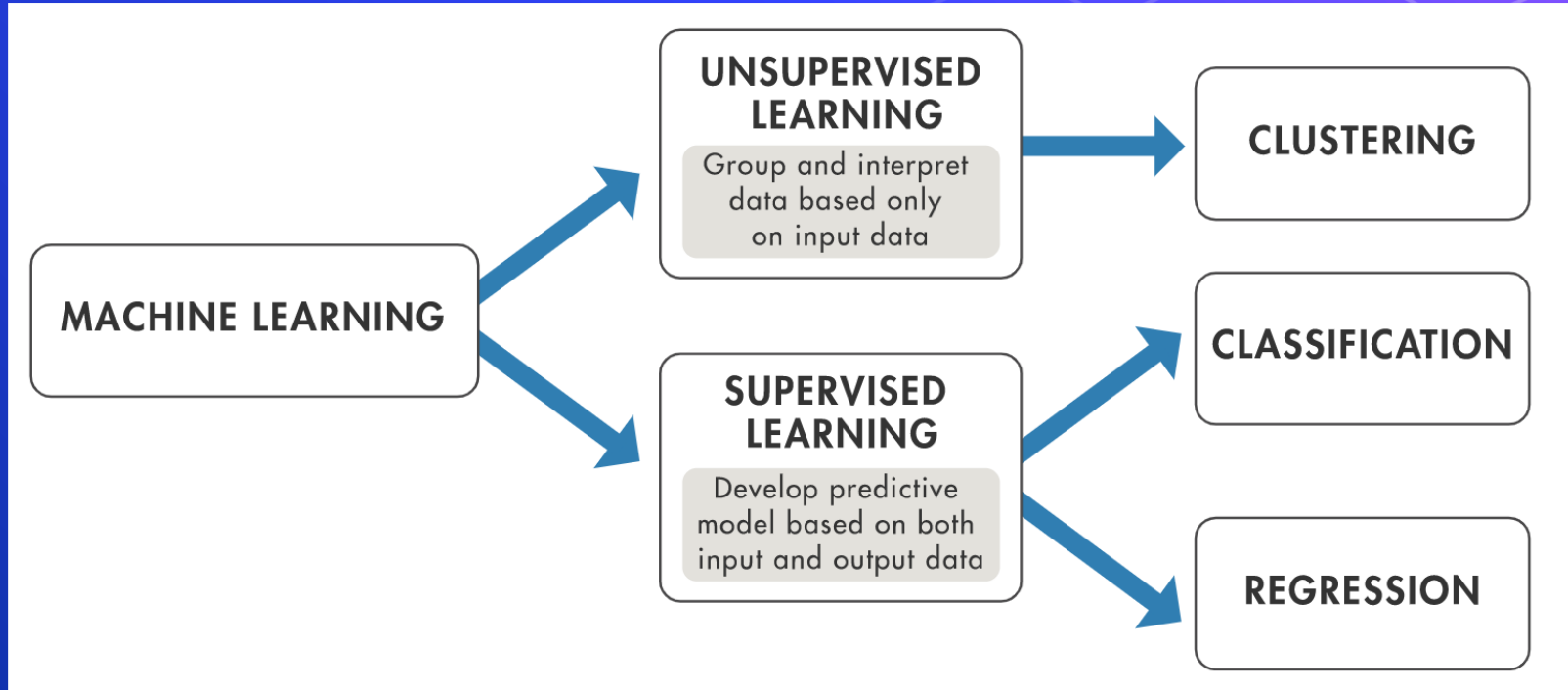
Machine Learning



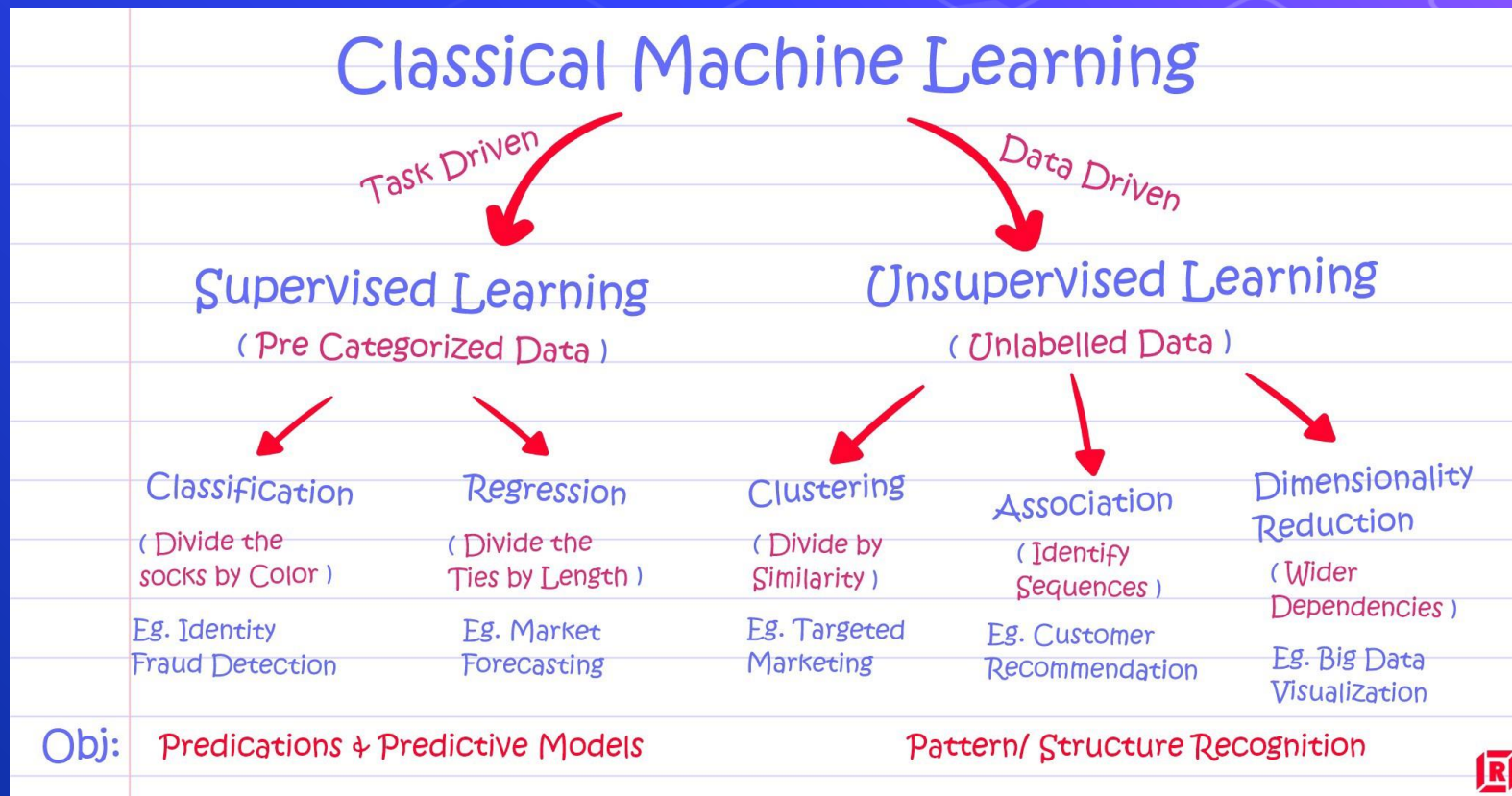
Machine Learning



Machine Learning



Machine Learning



Machine Learning

- ⬡ Supervised learning : Task Driven (Classification, Regression)
- ⬡ Unsupervised learning : Data Driven (Clustering)
- ⬡ Reinforcement learning :
 - Close to human learning.
 - Algorithm learns a policy of how to act in a given environment.
 - Every action has some impact in the environment, and the environment provides rewards that guides the learning algorithm.

Machine Learning

Supervised Learning

- Regression

- Simple Linear Regression
- Multiple Linear Regression
- Polynomial Regression
- Evaluating Model Performance

- Classification

- Logistic Regression
- K-Nearest Neighbors (KNN)
- Naive Bayes
- SVM
- Decision Trees
- Ensemble Methods
 - What is Bagging & Boosting
 - Random Forests
 - XGBoost
- Evaluating Model Performance

Unsupervised Learning

- Clustering

- KMeans
- Hierarchical Clustering
- Density Based Clustering - DBSCAN

- Association rule mining

- Apriori

- Dimension Reduction

- PCA
- LDA

- Evaluating Model Performance

Model Selection & Evaluation

- Cross Validation
- Hyperparameter Tuning
 - Grid Search
 - Randomized Search

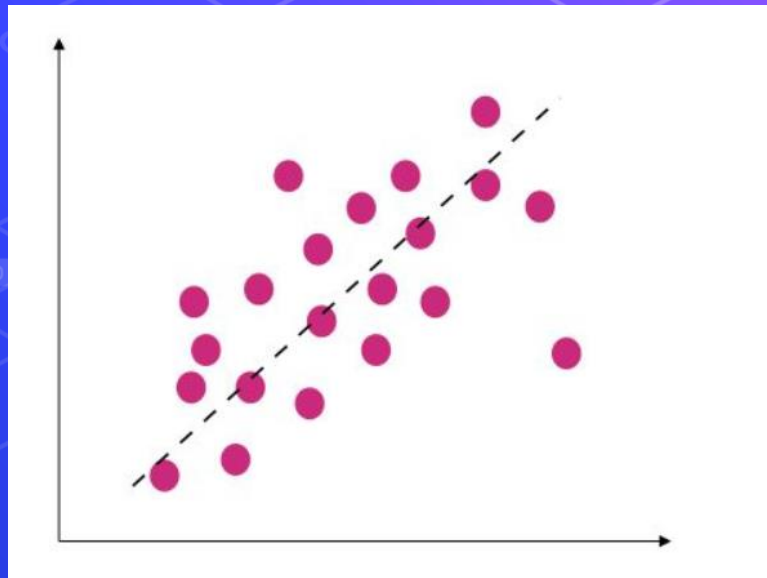
Recommendation Systems

Regression

Regression quantifies the relationship between one or more predictor variable(s) and one outcome variable.

For example,
it can be used to quantify the relative impacts of age, gender,
and diet (the predictors) on height (the output or dependent).

- House prices based on size, locations,...
- Salary prediction based on experience
- Stock Market prediction
- Weather prediction
- ...



Machine Learning

Supervised Learning

Regression

- Simple Linear Regression
- Multiple Linear Regression
- Polynomial Regression
- Evaluating Model Performance

Classification

- Logistic Regression
- K-Nearest Neighbors (KNN)
- Naive Bayes
- SVM
- Decision Trees
- Ensemble Methods
 - What is Bagging & Boosting
 - Random Forests
 - XGBoost
- Evaluating Model Performance

Unsupervised Learning

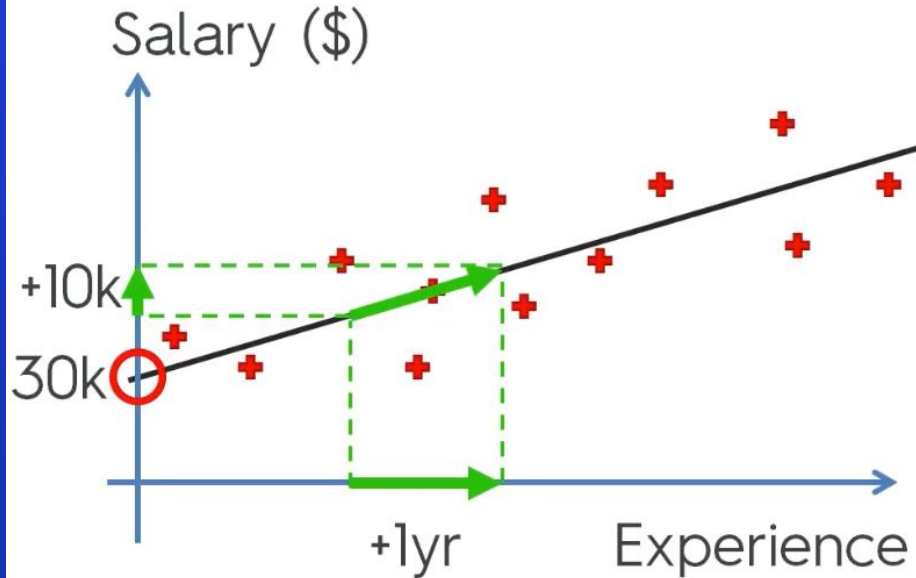
- Clustering
 - KMeans
 - Hierarchical Clustering
 - Density Based Clustering - DBSCAN
- Association rule mining
 - Apriori
- Dimension Reduction
 - PCA
 - LDA
- Evaluating Model Performance

Model Selection & Evaluation

- Cross Validation
- Hyperparameter Tuning
 - Grid Search
 - Randomized Search

Recommendation Systems

Simple Linear Regression (Equation of a straight line)



$$y = b_0 + b_1 * x$$

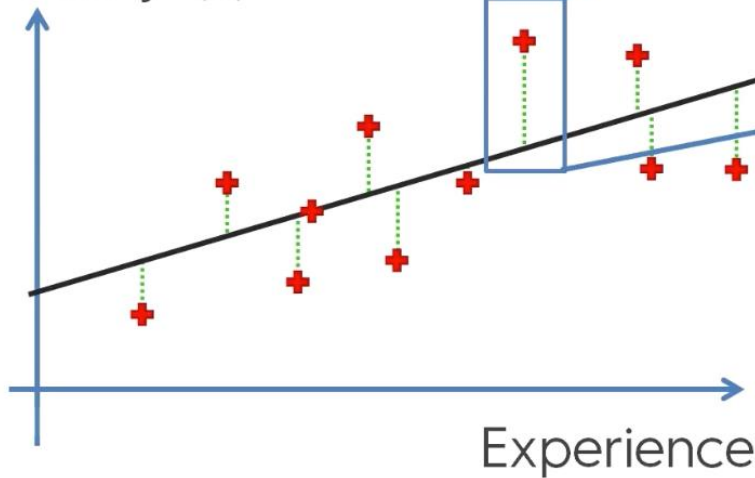
↓

$$\text{Salary} = b_0 + b_1 * \text{Experience}$$

Simple Linear Regression (Calculate Cost Function)

Simple Linear Regression:

Salary (\$)



Hypothesis: $h_{\theta}(x) = \theta_0 + \theta_1 x$

Parameters: θ_0, θ_1

Cost Function: $J(\theta_0, \theta_1) = \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2$

Goal: minimize $J(\theta_0, \theta_1)$

Simple Linear Regression (Example)

Theta0 = 5 , theta 1 = 2

Equation $h(x) = 5 + 2x$

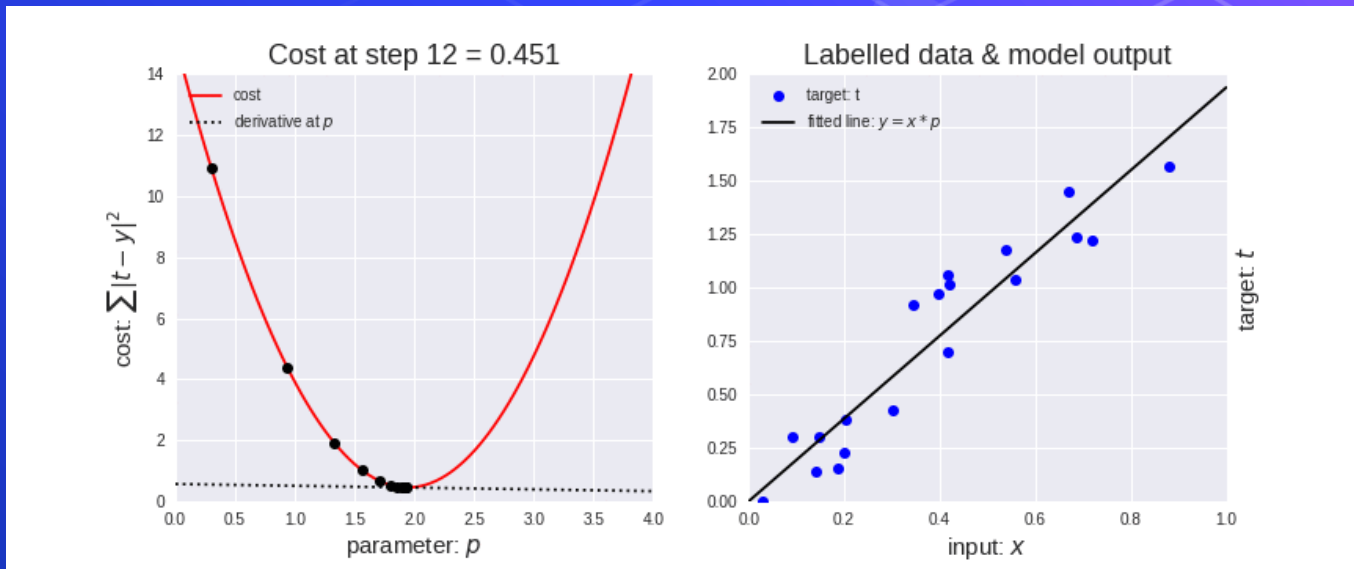
X	Y	h(x)	h(x) - y	(h(x) - y) ²
1	7	7	0	0
2	8	9	1	1
2	7	9	2	4
3	9	11	2	4
4	11	13	2	4
5	10	15	5	25
5	12	15	3	9

$$J(\theta_0, \theta_1) = \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2$$

$$J = 1 / 14 (0+1+4+4+4+25+9)$$

$$J = 47/14 = 3.3$$

Simple Linear Regression (Minimize cost using Gradient Descent)



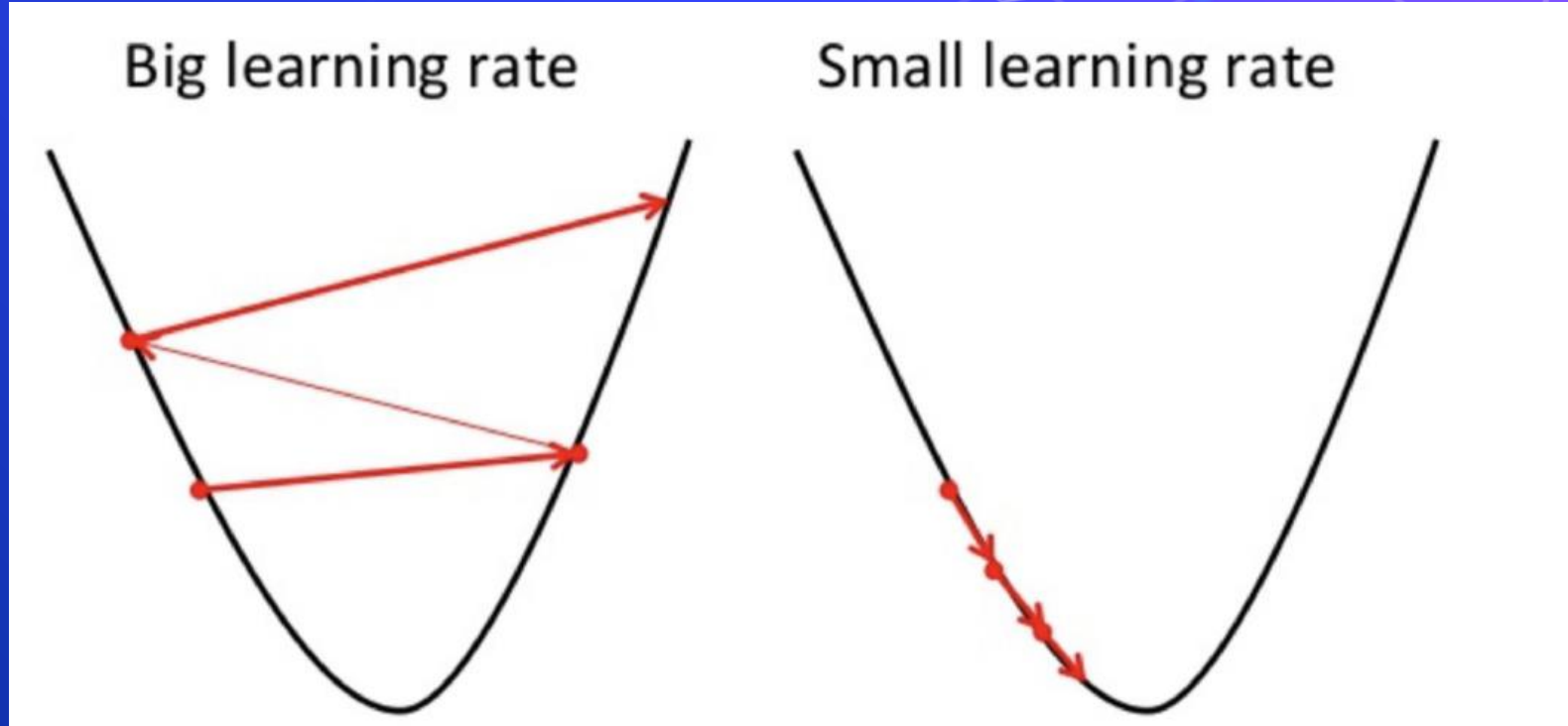
repeat until convergence:

$$\theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta_0, \theta_1)$$

$$\theta_0 := \theta_0 - \alpha \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x_i) - y_i)$$

$$\theta_1 := \theta_1 - \alpha \frac{1}{m} \sum_{i=1}^m ((h_{\theta}(x_i) - y_i) x_i)$$

Simple Linear Regression (Alpha constant)



Simple Linear Regression (Code)

```
1 from sklearn.linear_model import LinearRegression
2
3 # make model object
4 model = LinearRegression()
5
6 # train
7 model.fit(x_train, y_train)
8
9 # test
10 model.predict(x_test)
11
12 # calculate R2 score on training data
13 model.score(x_train, y_train)
14
15 # calculate R2 score on testing data
16 model.score(x_test, y_test)
17
18 # get model parameters
19 model.coef_
20 model.intercept_
```



Machine Learning

Supervised Learning

Regression

- Simple Linear Regression
- Multiple Linear Regression
- Polynomial Regression
- Evaluating Model Performance

Classification

- Logistic Regression
- K-Nearest Neighbors (KNN)
- Naive Bayes
- SVM
- Decision Trees
- Ensemble Methods
 - What is Bagging & Boosting
 - Random Forests
 - XGBoost
- Evaluating Model Performance

Unsupervised Learning

Clustering

- KMeans
- Hierarchical Clustering
- Density Based Clustering - DBSCAN

Association rule mining

- Apriori

Dimension Reduction

- PCA
- LDA

Evaluating Model Performance

Model Selection & Evaluation

- Cross Validation
- Hyperparameter Tuning
 - Grid Search
 - Randomized Search

Recommendation Systems

Multiple Linear Regression (Equation)

Simple
Linear
Regression

$$y = b_0 + b_1 * x_1$$

Multiple
Linear
Regression

Dependent variable (DV)

Independent variables (IVs)

$$y = b_0 + b_1 * x_1 + b_2 * x_2 + \dots + b_n * x_n$$

Constant

Coefficients

Multiple Linear Regression (Update theta values)

repeat until convergence: {

$$\theta_0 := \theta_0 - \alpha \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) \cdot x_0^{(i)}$$

$$\theta_1 := \theta_1 - \alpha \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) \cdot x_1^{(i)}$$

$$\theta_2 := \theta_2 - \alpha \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) \cdot x_2^{(i)}$$

...

}

Multiple Linear Regression (Code)

```
1 from sklearn.linear_model import LinearRegression
2
3 # make model object
4 model = LinearRegression()
5
6 # train
7 model.fit(x_train, y_train)
8
9 # test
10 model.predict(x_test)
11
12 # calculate R2 score on training data
13 model.score(x_train, y_train)
14
15 # calculate R2 score on testing data
16 model.score(x_test, y_test)
17
18 # get model parameters
19 model.coef_
20 model.intercept_
```



Machine Learning

Supervised Learning

Regression

- Simple Linear Regression
- Multiple Linear Regression
- **Polynomial Regression**
- Evaluating Model Performance

Classification

- Logistic Regression
- K-Nearest Neighbors (KNN)
- Naive Bayes
- SVM
- Decision Trees
- Ensemble Methods
 - What is Bagging & Boosting
 - Random Forests
 - XGBoost
- Evaluating Model Performance

Unsupervised Learning

Clustering

- KMeans
- Hierarchical Clustering
- Density Based Clustering - DBSCAN

Association rule mining

- Apriori

Dimension Reduction

- PCA
- LDA

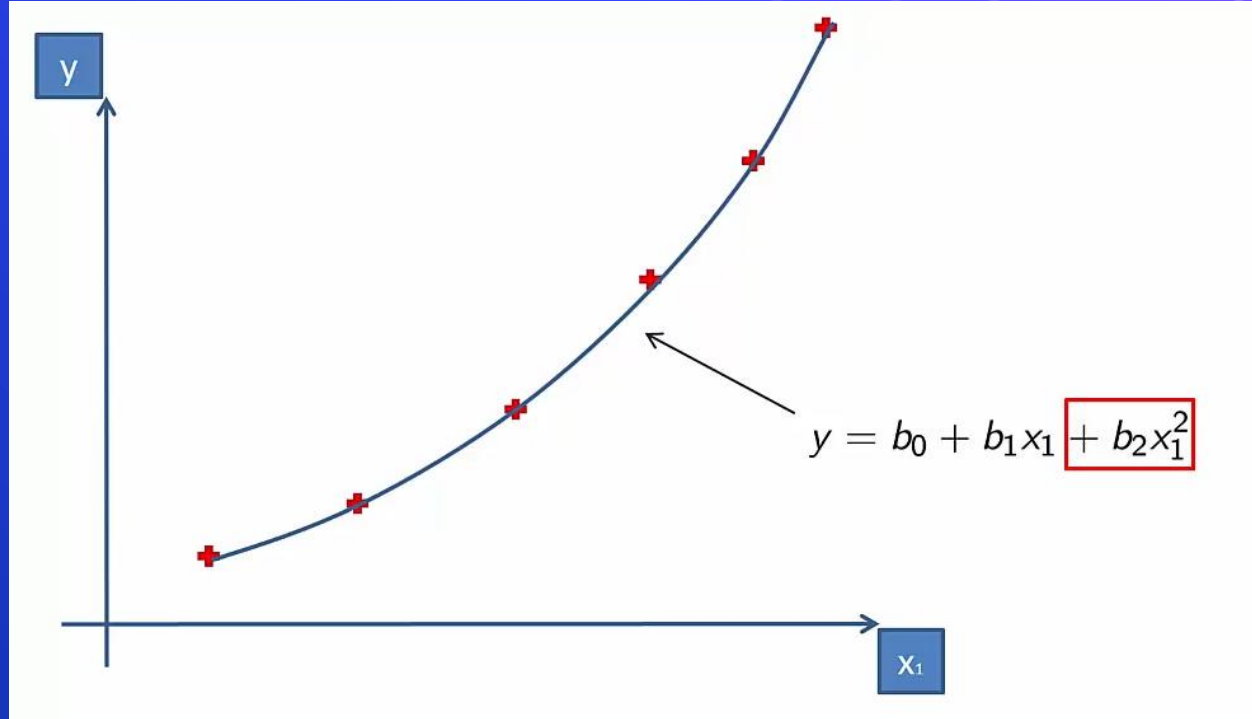
Evaluating Model Performance

Model Selection & Evaluation

- Cross Validation
- Hyperparameter Tuning
 - Grid Search
 - Randomized Search

Recommendation Systems

Polynomial Regression (Poly equation)



Polynomial Regression (Poly equation)

Polynomial
Linear
Regression

$$y = b_0 + b_1x_1 + b_2x_1^2 + \dots + b_nx_1^n$$

Polynomial Regression (Code)

```
1 from sklearn.linear_model import LinearRegression
2 from sklearn.preprocessing import PolynomialFeatures
3
4 # create poly features for a feature
5 poly = PolynomialFeatures(degree=2)
6 x_train_poly = poly.fit_transform(x_train)
7 x_test_poly = poly.fit_transform(x_test)
8
9 # make model object
10 model = LinearRegression()
11
12 # train
13 model.fit(x_train_poly, y_train)
14
15 # test
16 model.predict(x_test_poly)
17
18 # calculate R2 score on training data
19 model.score(x_train_poly, y_train)
20
21 # calculate R2 score on testing data
22 model.score(x_test_poly, y_test)
23
24 # get model parameters
25 model.coef_
26 model.intercept_
```





Machine Learning



Supervised Learning

- **Regression**

- Simple Linear Regression
- Multiple Linear Regression
- Polynomial Regression
- **Evaluating Model Performance**

- **Classification**

- Logistic Regression
- K-Nearest Neighbors (KNN)
- Naive Bayes
- SVM
- Decision Trees
- Ensemble Methods
 - What is Bagging & Boosting
 - Random Forests
 - XGBoost
- Evaluating Model Performance



Unsupervised Learning

- **Clustering**

- KMeans
- Hierarchical Clustering
- Density Based Clustering - DBSCAN

- **Association rule mining**

- Apriori

- **Dimension Reduction**

- PCA
- LDA

- **Evaluating Model Performance**



Model Selection & Evaluation

- Cross Validation
- Hyperparameter Tuning
 - Grid Search
 - Randomized Search



Recommendation Systems

Evaluating Model Performance (R²)

Hrs Studied (X)	Marks (Y)
0	40
2	52
3	53
4	55
4	56
5	72
6	71
6	88
7	56
7	74
8	89
9	67
9	89
5.38	66.31
Mean	

$$R^2 = \text{SSR} / \text{SST}$$

$$= 1844.12 / 3028.77$$

$$= 0.60886$$

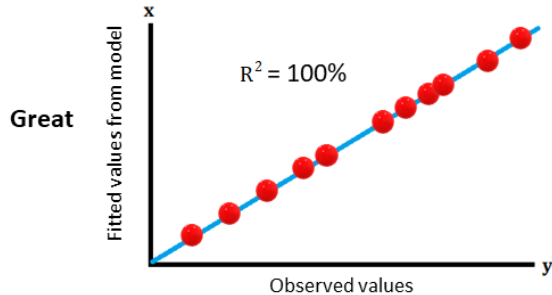
Higher the value → Variation in Y is explained by variation in X.

$$\hat{y} \rightarrow y \quad \text{SSR} \rightarrow \text{SST} \quad R^2 \approx 1$$

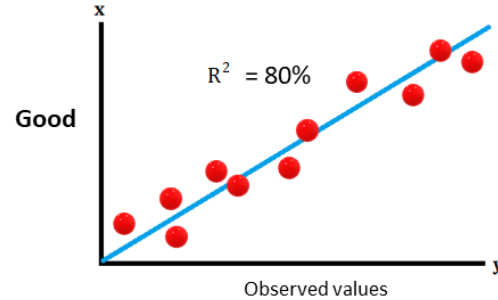
$(Y - \bar{Y})^2$	$(\hat{Y} - \bar{Y})^2$
692.22	600.74
204.78	237.47
177.16	117.94
127.92	39.82
106.30	39.82
32.38	3.10
22.00	7.78
470.46	7.78
106.30	53.88
59.14	53.88
514.84	141.37
0.48	270.27
514.84	270.27
3028.77	1844.12
SST	SSR

Evaluating Model Performance (R²)

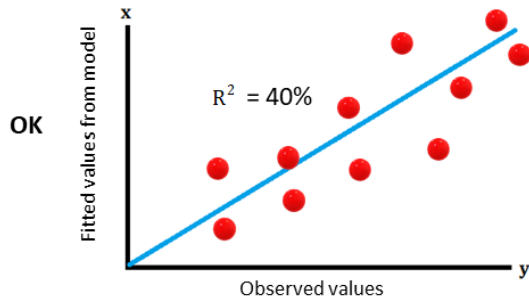
Comparison of R-Squared for Different Linear Models (Same Data Set)



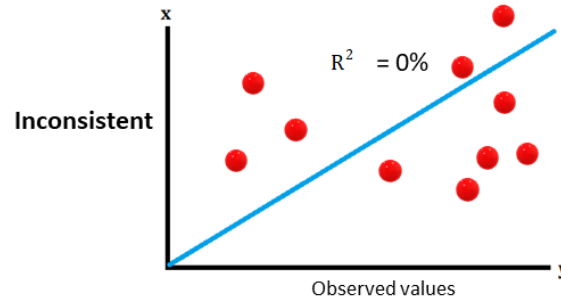
Fitted = observed: Model explains all variance



Model explains bulk of variance



Model explains 40% of variance, so is reasonable.



Model fails to explain any variance



Machine Learning



Supervised Learning

- Regression
 - Simple Linear Regression
 - Multiple Linear Regression
 - Polynomial Regression
 - Evaluating Model Performance
- **Classification**
 - Logistic Regression
 - K-Nearest Neighbors (KNN)
 - Naive Bayes
 - SVM
 - Decision Trees
 - Ensemble Methods
 - What is Bagging & Boosting
 - Random Forests
 - XGBoost
 - Evaluating Model Performance



Unsupervised Learning

- Clustering
 - KMeans
 - Hierarchical Clustering
 - Density Based Clustering - DBSCAN
- Association rule mining
 - Apriori
- Dimension Reduction
 - PCA
 - LDA
- Evaluating Model Performance



Model Selection & Evaluation

- Cross Validation
- Hyperparameter Tuning
 - Grid Search
 - Randomized Search



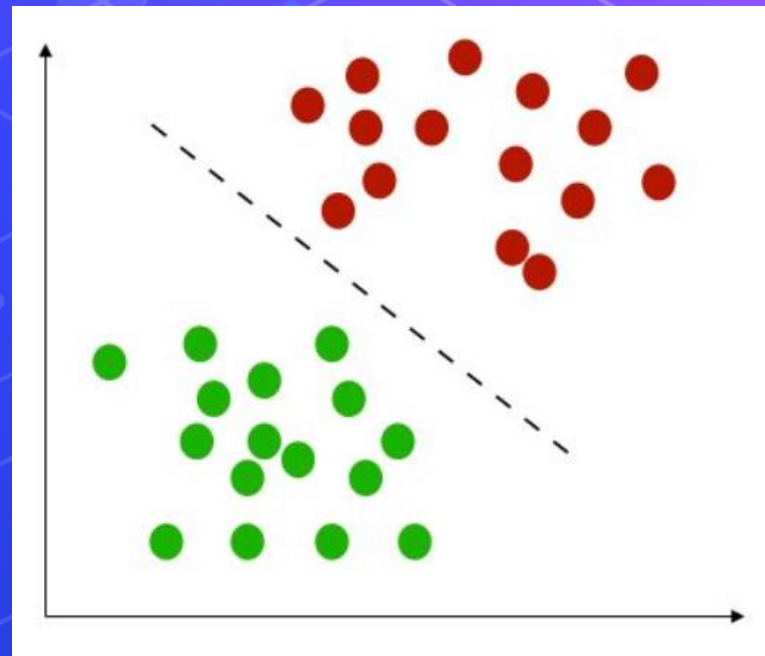
Recommendation Systems

Classification

Classification specifies the class to which data elements belong to and is best used when the output has finite and discrete values. It predicts a class for an input variable as well.

For example,
it can be used to classify the animal in an image, it's either a cat or a dog.

- ⬡ Email spam detection
- ⬡ Face classification
- ⬡ Patient has cancer or not
- ⬡ Fraud detection
- ⬡ ...





Machine Learning



Supervised Learning

- Regression
 - Simple Linear Regression
 - Multiple Linear Regression
 - Polynomial Regression
 - Evaluating Model Performance
- **Classification**
 - Logistic Regression
 - K-Nearest Neighbors (KNN)
 - Naive Bayes
 - SVM
 - Decision Trees
 - Ensemble Methods
 - What is Bagging & Boosting
 - Random Forests
 - XGBoost
 - Evaluating Model Performance



Unsupervised Learning

- Clustering
 - KMeans
 - Hierarchical Clustering
 - Density Based Clustering - DBSCAN
- Association rule mining
 - Apriori
- Dimension Reduction
 - PCA
 - LDA
- Evaluating Model Performance



Model Selection & Evaluation

- Cross Validation
- Hyperparameter Tuning
 - Grid Search
 - Randomized Search



Recommendation Systems

Logistic Regression

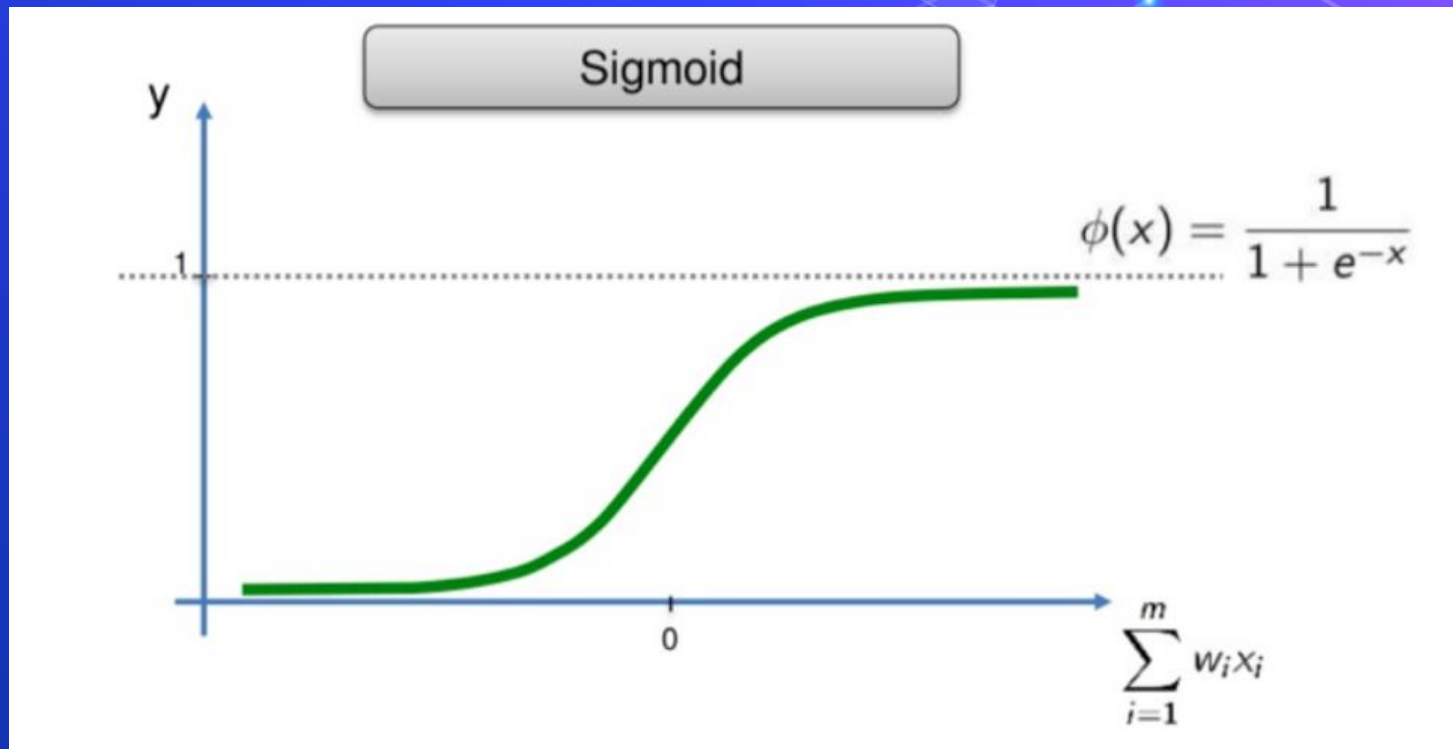
The Logistic Function

$$y = \frac{1}{1 + e^{-f(x_1, x_2, \dots, x_n)}} \in (0, 1)$$

where

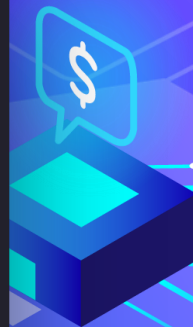
$$f(x_1, x_2, \dots, x_n) = a_0 + a_1x_1 + \dots + a_nx_n \in (-\infty, +\infty)$$

Logistic Regression



Logistic Regression

```
1 from sklearn.linear_model import LogisticRegression
2 from sklearn.metrics import classification_report
3 from sklearn.metrics import accuracy_score
4 from sklearn.metrics import confusion_matrix
5
6
7 model = LogisticRegression()
8 model.fit(x_train, y_train)
9 y_pred = model.predict(x_test)
10
11 print(confusion_matrix(y_test, y_pred))
12 print(accuracy_score(y_test, y_pred))
13 print(classification_report(y_test, y_pred))
```





Machine Learning



Supervised Learning

- Regression
 - Simple Linear Regression
 - Multiple Linear Regression
 - Polynomial Regression
 - Evaluating Model Performance
- **Classification**
 - Logistic Regression
 - **K-Nearest Neighbors (KNN)**
 - Naive Bayes
 - SVM
 - Decision Trees
 - Ensemble Methods
 - What is Bagging & Boosting
 - Random Forests
 - XGBoost
 - Evaluating Model Performance



Unsupervised Learning

- Clustering
 - KMeans
 - Hierarchical Clustering
 - Density Based Clustering - DBSCAN
- Association rule mining
 - Apriori
- Dimension Reduction
 - PCA
 - LDA
- Evaluating Model Performance



Model Selection & Evaluation

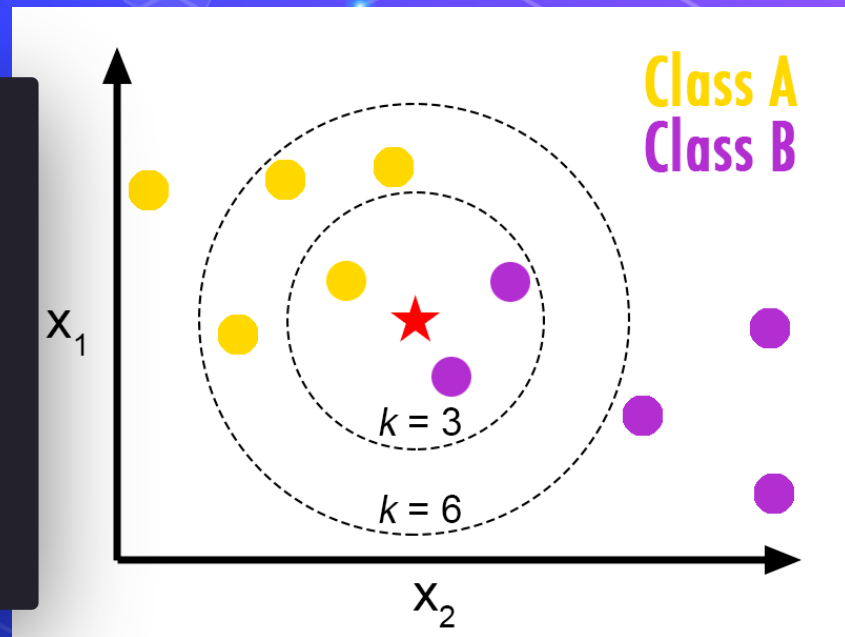
- Cross Validation
- Hyperparameter Tuning
 - Grid Search
 - Randomized Search



Recommendation Systems

K-Nearest Neighbors (KNN)

```
1 from sklearn.neighbors import KNeighborsClassifier
2 from sklearn.metrics import classification_report
3 from sklearn.metrics import accuracy_score
4 from sklearn.metrics import confusion_matrix
5
6
7 model = KNeighborsClassifier()
8 model.fit(x_train, y_train)
9 y_pred = model.predict(x_test)
10
11 print(confusion_matrix(y_test, y_pred))
12 print(accuracy_score(y_test, y_pred))
13 print(classification_report(y_test, y_pred))
```





Machine Learning



Supervised Learning

- Regression
 - Simple Linear Regression
 - Multiple Linear Regression
 - Polynomial Regression
 - Evaluating Model Performance
- **Classification**
 - Logistic Regression
 - K-Nearest Neighbors (KNN)
 - **Naive Bayes**
 - SVM
 - Decision Trees
 - Ensemble Methods
 - What is Bagging & Boosting
 - Random Forests
 - XGBoost
 - Evaluating Model Performance



Unsupervised Learning

- Clustering
 - KMeans
 - Hierarchical Clustering
 - Density Based Clustering - DBSCAN
- Association rule mining
 - Apriori
- Dimension Reduction
 - PCA
 - LDA
- Evaluating Model Performance



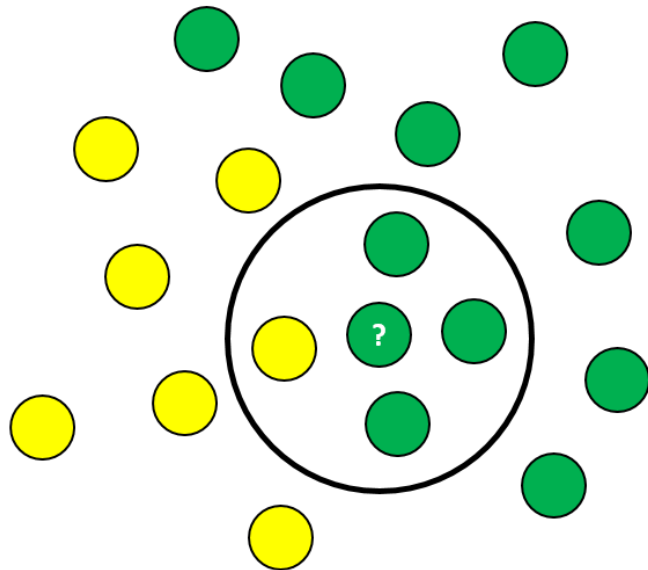
Model Selection & Evaluation

- Cross Validation
- Hyperparameter Tuning
 - Grid Search
 - Randomized Search



Recommendation Systems

Naive Bayes



posterior probability

posterior probability

$$P(\text{yellow}) = \frac{7}{17}$$

$$P(\text{green}) = \frac{10}{17}$$

$$P'(? \mid \text{green}) = \frac{3}{10}$$

$$P'(? \mid \text{yellow}) = \frac{1}{7}$$

prior probabilities

number of samples in a given class
divided by the total number of samples

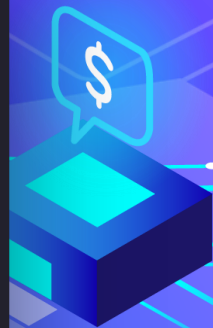
we consider just the vicinity
of the new sample we wan to classify

$$P''(? \text{ is green}) = P(\text{green}) * P'(? \mid \text{green}) = \frac{10}{17} * \frac{3}{10} = \frac{30}{170}$$

$$P''(? \text{ is yellow}) = P(\text{yellow}) * P'(? \mid \text{yellow}) = \frac{7}{17} * \frac{1}{7} = \frac{7}{119}$$

Naive Bayes

```
1 from sklearn.naive_bayes import MultinomialNB
2 from sklearn.metrics import classification_report
3 from sklearn.metrics import accuracy_score
4 from sklearn.metrics import confusion_matrix
5
6
7 model = MultinomialNB()
8 model.fit(x_train, y_train)
9 y_pred = model.predict(x_test)
10
11 print(confusion_matrix(y_test, y_pred))
12 print(accuracy_score(y_test, y_pred))
13 print(classification_report(y_test, y_pred))
```



Questions ?!



Thanks!

>_ Live long and prosper

