

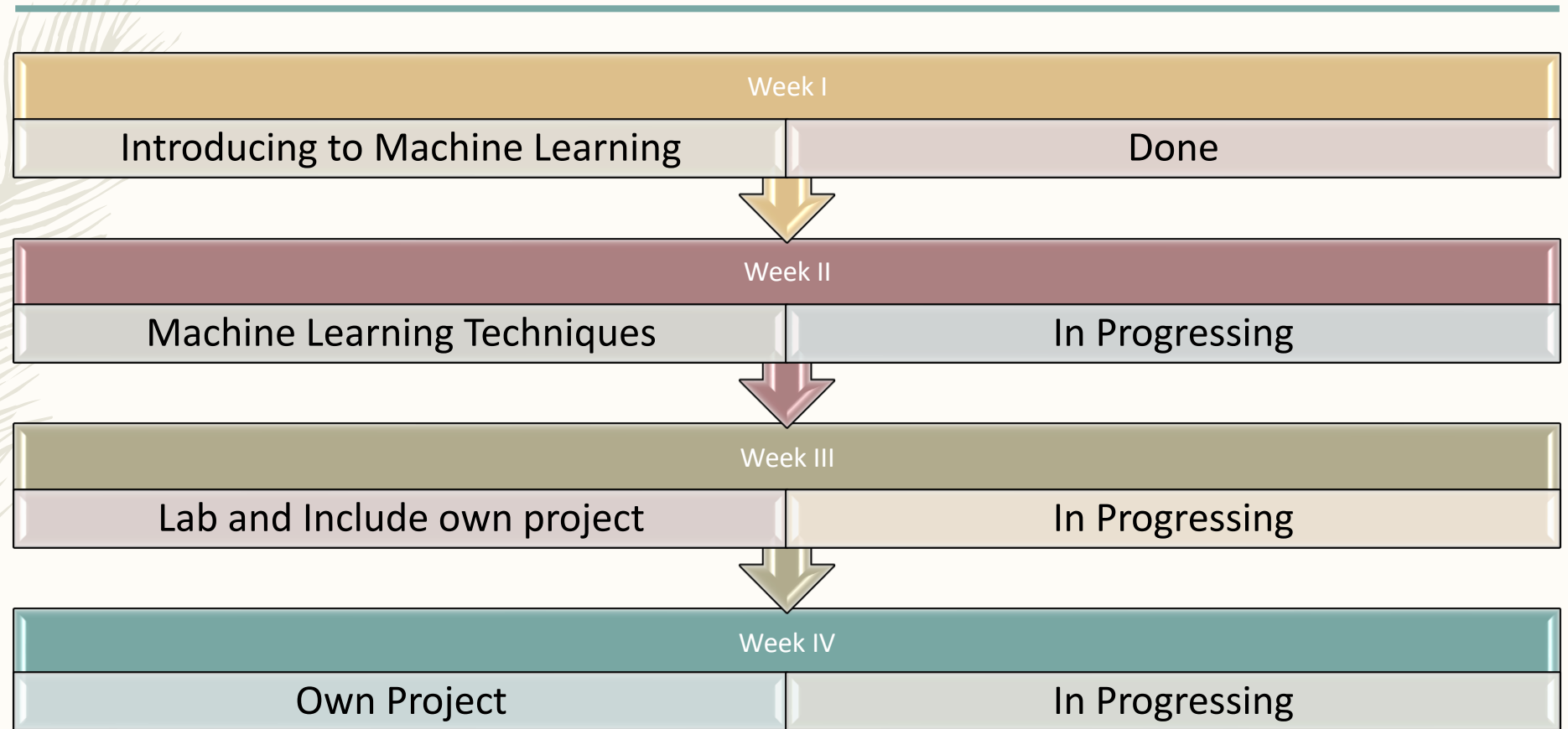
Machine Learning Techniques

Week II

Lecturer : Mr. Hel Chanthan

Student : Chek Nita

Planning





Outline

- Machine learning techniques
 - Regression
 - Simple Linear regression
 - Multiple Linear regression
 - Classification
 - K-Nearest Neighbor
 - Decision Tree
 - Logistic Linear
 - Clustering
 - Recommender Systems



Machine Learning Techniques

- Machine Learning techniques :

- Regression/Estimation

- *Predicting continuous values*

- Classification

- *Prediction the item class/ category of a case*

- Clustering

- *Finding the structure of data / summarization*

- Associations

- *Finding item that co-occur*

- Anomaly detection

- *Discovering abnormal and unusual case*

- Sequence mining

- *Predict next events*

- Dimension Reduction

- *Reducing the size of data*

- Recommendation systems

- *Recommending items*

Regression

– What is Regression ?

	X: Independent variable			Y: Dependent variable
	ENGINE SIZE	CYLINDERS	FUEL CONSUMPTION_COMB	CO2 EMISSIONS
0	2.0	4	8.5	196
1	2.4	4	9.6	221
2	1.5	4	5.9	136
3	3.5	6	11.1	255
4	3.5	6	10.6	244
5	3.5	6	10.0	230
6	3.5	6	10.1	232
7	3.7	6	11.1	255
8	3.7	6	11.6	267
9	2.4	4	9.2	?

Continuous Values

Regression is the process of predicting a continuous value



Regression

- Type of Regression models :
 - Simple Linear Regression
 - *Simple Linear Regression*
 - *Simple Non-Linear Regression*
 - Multiple Linear Regression
 - *Multiple Linear Regression*
 - *Multiple Non linear Regression*

Applications of regression

- Sales forecasting
- Satisfaction analysis
- Price estimation
- Employment income

Simple Linear Regression

	ENGINE SIZE	CYLINDERS	FUEL CONSUMPTION_COMB	CO2EMISSIONS
0	2.0	4	8.5	196
1	2.4	4	9.6	221
2	1.5	4	5.9	136
3	3.5	6	11.1	255
4	3.5	6	10.6	244
5	3.5	6	10.0	230
6	3.5	6	10.1	232
7	3.7	6	11.1	255
8	3.7	6	11.6	267
9	2.4	4	9.2	?

Simple Linear Regression

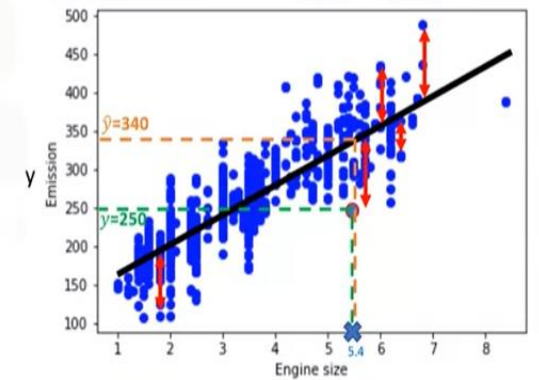
How to find the best fit?

$x_1 = 5.4$ independent variable
 $y = 250$ actual Co2 emission of x_1

$\hat{y} = \theta_0 + \theta_1 x_1$
 $\hat{y} = 340$ the predicted emission of x_1

Error = $y - \hat{y}$
= $250 - 340$
= -90

$$MSE = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2$$



Simple Linear Regression

Estimating the parameters

	ENGINE SIZE	CYLINDERS	FUEL CONSUMPTION_COMB	CO2 EMISSIONS
0	2.0	4	8.5	196
1	2.4	4	9.6	221
2	1.5	4	5.9	136
3	3.5	6	11.1	255
4	3.5	6	10.6	244
5	3.5	6	10.0	230
6	3.5	6	10.1	232
7	3.7	6	11.1	255
8	3.7	6	11.6	267

$$\hat{y} = \theta_0 + \theta_1 x_1$$

$$\theta_1 = \frac{\sum_{i=1}^s (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^s (x_i - \bar{x})^2}$$

$$\bar{x} = (2.0 + 2.4 + 1.5 + \dots) / 9 = 3.03$$

$$\bar{y} = (196 + 221 + 136 + \dots) / 9 = 226.22$$

$$\theta_1 = \frac{(2.0 - 3.03)(196 - 226.22) + (2.4 - 3.03)(221 - 226.22) + \dots}{(2.0 - 3.03)^2 + (2.4 - 3.03)^2 + \dots}$$

$$\theta_1 = 39$$

$$\theta_0 = \bar{y} - \theta_1 \bar{x}$$

$$\theta_0 = 226.22 - 39 * 3.03$$

$$\theta_0 = 125.74$$

Multiple Linear Regression

Predicting continuous values with multiple linear regression

$$\text{Co2 Em} = \theta_0 + \theta_1 \text{Engine size} + \theta_2 \text{Cylinders} + \dots$$

$$\hat{y} = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \dots + \theta_n x_n$$

$$\hat{y} = \theta^T X$$

$$\theta^T = [\theta_0, \theta_1, \theta_2, \dots] \quad X = \begin{bmatrix} 1 \\ x_1 \\ x_2 \\ \dots \end{bmatrix}$$

X: Independent variable

Y: Dependent variable

	ENGINE SIZE	CYLINDERS	FUEL CONSUMPTION_COMB	CO2 EMISSIONS
0	2.0	4	8.5	196
1	2.4	4	9.6	221
2	1.5	4	5.9	136
3	3.5	6	11.1	255
4	3.5	6	10.6	244
5	3.5	6	10.0	230
6	3.5	6	10.1	232
7	3.7	6	11.1	255
8	3.7	6	11.6	267
9	2.4	4	9.2	?



Multiple Linear Regression

Estimating multiple linear regression parameters

- How to estimate θ ?
 - Ordinary Least Squares
 - Linear algebra operations
 - Takes a long time for large datasets (10K+ rows)
 - An optimization algorithm
 - Gradient Descent
 - Proper approach if you have a very large dataset

Non-Linear Regression

[Share](#)

What is non-linear regression?

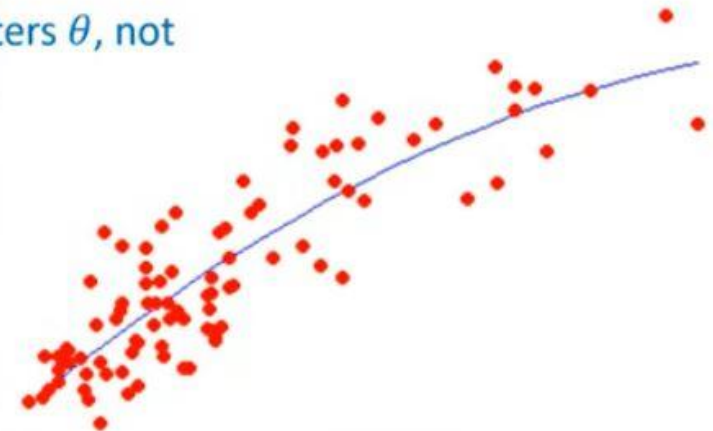
- To model non-linear relationship between the dependent variable and a set of independent variables
- \hat{y} must be a non-linear function of the parameters θ , not necessarily the features x

$$\hat{y} = \theta_0 + \theta_2^2 x$$

$$\hat{y} = \theta_0 + \theta_1 \theta_2^x$$

$$\hat{y} = \log(\theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3)$$

$$\hat{y} = \frac{\theta_0}{1 + \theta_1^{(x - \theta_2)}}$$





Classification

What is classification?

- A supervised learning approach
- Categorizing some unknown items into a discrete set of categories or “classes”
- The target attribute is a categorical variable

Classification

Normal

age	ed	employ	address	income	debtinc	creddebt	othdebt	default
41	3	17	12	176	9.3	11.359	5.009	1
27	1	10	6	31	17.3	1.362	4.001	0
40	1	15	14	55	5.5	0.856	2.169	0
41	1	15	14	120	2.9	2.659	0.821	0
24	2	2	0	28	17.3	1.787	3.057	1
41	2	5	5	25	10.2	0.393	2.157	0
39	1	20	9	67	30.6	3.834	16.668	0
43	1	12	11	38	3.6	0.129	1.239	0
24	1	3	4	19	24.4	1.358	3.278	1
36	1	0	13	25	19.7	2.778	2.147	0

Categorical Variable

age	ed	employ	address	income	debtinc	creddebt	othdebt	default
37	2	16	10	130	9.3	10.23	3.21	0

Multiple

Age	Sex	BP	Cholesterol	Na	K	Drug
23	F	HIGH	HIGH	0.793	0.031	drugY
47	M	LOW	HIGH	0.739	0.056	drugC
47	M	LOW	HIGH	0.697	0.069	drugC
28	F	NORMAL	HIGH	0.564	0.072	drugX
61	F	LOW	HIGH	0.559	0.031	drugY
22	F	NORMAL	HIGH	0.677	0.079	drugX
49	F	NORMAL	HIGH	0.79	0.049	drugY
41	M	LOW	HIGH	0.767	0.069	drugC
60	M	NORMAL	HIGH	0.777	0.051	drugY
43	M	LOW	NORMAL	0.526	0.027	drugY

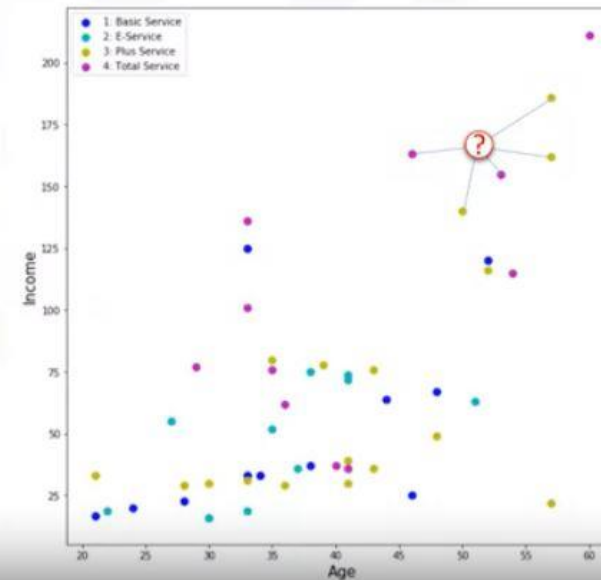
Categorical Variable

Age	Sex	BP	Cholesterol	Na	K	Drug
36	F	LOW	HIGH	0.697	0.069	

K-Nearest Neighbor

What is K-Nearest Neighbor (or KNN)?

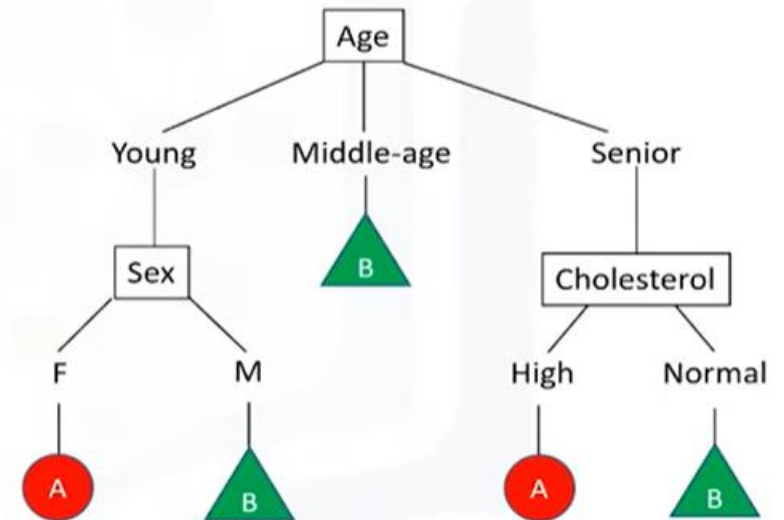
- A method for **classifying** cases based on their similarity to other cases
- Cases that are near each other are said to be “**neighbors**”
- Based on **similar cases with same class labels** are near each other



Decision Tree

Decision tree learning algorithm

1. Choose an attribute from your dataset.
2. Calculate the significance of attribute in splitting of data.





Logistic Regression

Logistic regression applications

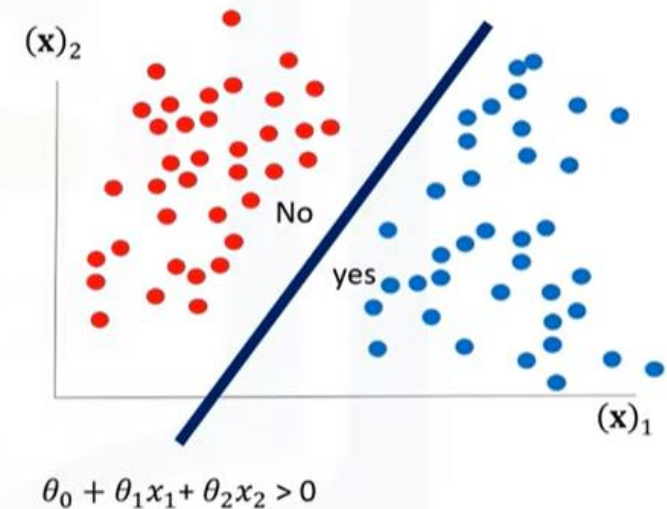
- Predicting the probability of a person having a heart attack
- Predicting the mortality in injured patients
- Predicting a customer's propensity to purchase a product or halt a subscription
- Predicting the probability of failure of a given process or product
- Predicting the likelihood of a homeowner defaulting on a mortgage

Logistic Regression

When is logistic regression suitable?

Share

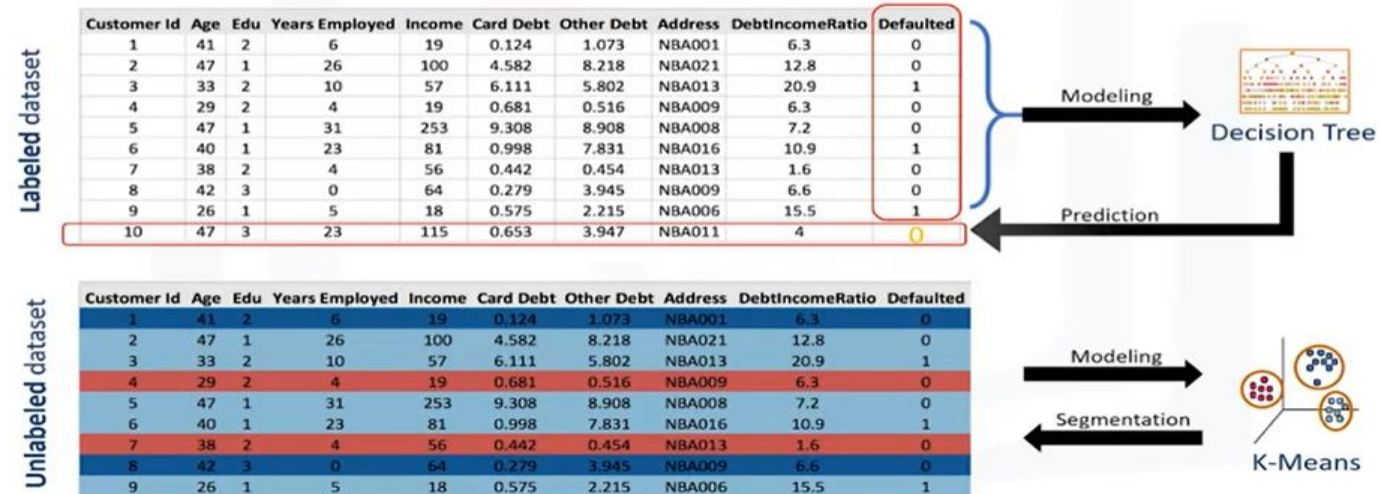
- If your data is binary
 - 0/1, YES/NO, True/False
- If you need probabilistic results
- When you need a linear decision boundary
- If you need to understand the impact of a feature



Clustering

- Clustering is similar to regression and classification but mostly use for unlabeled dataset.

Clustering Vs. classification





Clustering Applications

- **PUBLICATION:**

- Auto-categorizing news based on their content
- Recommending similar news articles

- **MEDICINE:**

- Characterizing patient behavior

- **BIOLOGY:**

- Clustering genetic markers to identify family ties

- **RETAIL/MARKETING:**

- Identifying buying patterns of customers
- Recommending new books or movies to new customers

- **BANKING:**

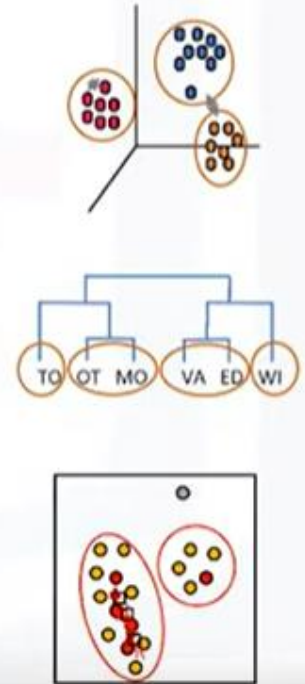
- Fraud detection in credit card use
- Identifying clusters of customers (e.g., loyal)

- **INSURANCE:**

- Fraud detection in claims analysis
- Insurance risk of customers

Clustering Algorism

- Partitioned-based Clustering
 - Relatively efficient
 - E.g. k-Means, k-Median, Fuzzy c-Means
- Hierarchical Clustering
 - Produces trees of clusters
 - E.g. Agglomerative, Divisive
- Density-based Clustering ★
 - Produces arbitrary shaped clusters
 - E.g. DBSCAN



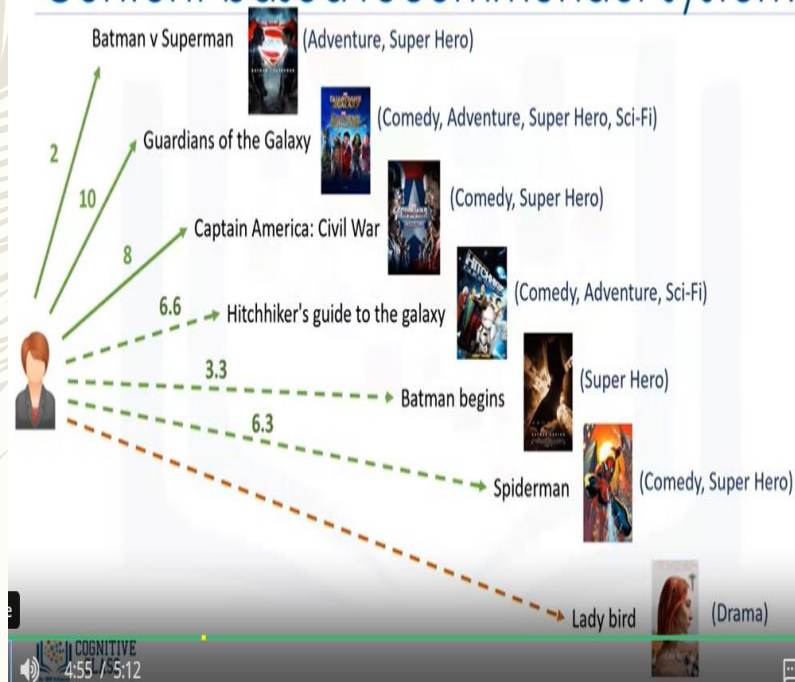


Recommender Systems

- Recommender systems capture the pattern of peoples' behavior and use it to predict what else they might want or like.
- Type of recommender systems :
 - Content-Based recommendation
 - Collaborative Filtering

Content-based vs Collaboration

Content-based recommender systems



Collaborative filtering

