# Distance-based Classification for Product Recommendation

Unit 6 Machine Learning

# Product Recommendation

When we search a product, retail websites suggest similar products



Image source: Flipkart tech blog

# Product Recommendation

When we search a product, retail websites suggest similar products

How does it know which products are "similar" to the original product?

What is even meant by "similarity of products"?

# Product Recommendation

When we search a product, retail websites suggest similar products

How does it know which products are "similar" to the original product?

What is even meant by "similarity of products"?

Similarity notion 1: Based on customer statistics

Similarity notion 2: Based on features of products

# Recommendation based on statistics

Two products which are often bought together may be considered "similar"



Image source: Flipkart

# Recommendation based on features

- Shirt and tie often bought together, but are they similar?

- What if the user is looking for choices among shirts?

- Product Type Matching The recommender system should show other shirts!

Recommended shirts should be "similar" to the shirt being viewed

Similarity defined in terms of attributes

Colour, size, brand, perhaps design and texture ………

Query          Positive          In-class negative          Out-of-class negative

Image source: Google Images

# Data Representation

- Represent each product by a vector of attributes or "features"
- Training data:  $(X\_i, Y\_i)$ where X: feature vector, Y: label
- Test data: $(X_{test},$ ?)  [Label to be predicted]

# Data Representation

- Represent each product by a vector of attributes or "features"
- Training data:  (X_i, Y_i) where X: feature vector, Y: label
- Test data: ($X_{test,}$ ?)  [Label to be predicted]

- Label: can be binary, discrete or continuous
- Binary label: product recommendation (similar or not?)
- Discrete label: rating by one user (1star, 2star, 3star, 4star or 5star)?
- Continuous label: mean rating by many users (1,5)

# A simple idea for classification

- Training data: $(X_i, Y_i)$ where X: feature vector, Y: label
- Test data: $(X_{test}, ?)$
- X: representation of data-points (feature)


- Idea: things that "look" similar, are usually the same type!
- If $X_{test} = X_i$, then probably $Y_{test} = Y_i$ !

# But is it a good idea?

- Training data: $(X_i, Y_i)$ where X: feature vector, Y: label
- Test data: $(X_{test}, ?)$
- X: representation of data-points (feature)

- Idea: things that "look" similar, are usually the same type!
- If $X_{test} = X_i$, then probably $Y_{test} = Y_i$ !

- Is X a good enough representation?
- In reality, $X_{test} = X_i$ will rarely happen (especially if X is continuous!)

# But is it a good idea?

- Training data: $(X_i, Y_i)$ where X: feature vector, Y: label
- Test data: $(X_{test}, ?)$
- X: representation of data-points (feature)

- Idea: things that "look" similar, are usually the same type!
- If $X_{test} = X_i$, then probably $Y_{test} = Y_i$ !

- Is X a good enough representation?   - Let's assume it is
- In reality, $X_{test} = X_i$ will rarely happen (especially if X is continuous!)

# Distance between feature vectors

- If X is continuous-valued, $X_{test}$ = X_i will almost surely never happen!
- But, $X_{test}$ ~ X_i is possible!

- $X_{test}$ ~ X_i : Euclidean Distance between the two points is very less
- $|| X_{test} - X\_i ||_2$ is very low!

- $||a - b||_2 = \sqrt{\sum(a\_i - b\_i)2}$ (also called the $l_2$-norm of a-b)

# Nearest-Neighbor Classification

- Training: N labelled examples $(X_i, Y_i)$ where i: 1 to N
- Function learnt: the training set itself!
- Testing: X_test
- $Y_{pred} = Y_n$, where $n = \arg\min_i ||X_{test} - X_i||_2$

# Nearest-Neighbor Classification

- Training: N labelled examples $(X_i, Y_i)$ where i: 1 to N
- Function learnt: the training set itself!
- Testing: $X_{test}$
- $Y_{pred} = Y_n$, where $n = \arg\min_i ||X_{test} - X_i||_2$

- Compute the Euclidean distance between the test datapoint and each of the N training datapoints
- Choose that training point for which this distance is minimum (Nearest-Neighbor)
- Use its label as the predicted label of the test point!

# Nearest-Neighbor Classification

- Training: N labelled examples $(X_i, Y_i)$ where $i$: 1 to N

- Function learnt: the training set itself!

- Testing: $X_{test}$

- $Y_{pred} = Y_n$, where $n = \arg\min_i ||X_{test} - X_i||_2$

- Not very robust due to outliers!

- Too much computation and storage required!

# K-nearest Neighbors Classification

- Problem 1: The nearest neighbour of the test point may be outlier!

- Outliers are rare
- K nearest neighbors: unlikely to contain many outliers

- Solution:
1) Sort the training points according to distance from test point
2) Choose the first K training points
3) Predicted label = most frequent label among them!

# K-nearest Neighbors Regression

- Problem 1: The nearest neighbour of the test point may be outlier!

- Outliers are rare
- K nearest neighbors: unlikely to contain many outliers

- Solution:
1) Sort the training points according to distance from test point
2) Choose the first K training points
3) Predicted label = mean of their labels!

# Nearest Mean Classification

- Problem 2: Too much computation (N) and storage (N*(D+1)) required!

- One solution: keep only one representative from each class
- How to choose the representative?
- Mean of feature vectors all data-points in that class!

- Mean for class k: $\mu_k = \sum 1(Y_i=k) * X_i \, 1(Y_i=k) \, / \, \sum 1(Y_i=k)$
- Compare test point $X_{test}$ with each $\mu_k$ and choose label of the closest!
- $Y_{pred} = \text{argmin}_k \, ||X_{test} - \mu_k||_2$

# Nearest Mean Classification

- Problem 2: Too much computation (N) and storage (N*(D+1)) required!

- One solution: keep only one representative from each class

- How to choose the representative?

- Mean of feature vectors all data-points in that class!

- Mean for class k: $\mu_k = \sum 1(Y_i=k) * X_i\ 1(Y_i=k) / \sum 1(Y_i=k)$

- Compare test point $X_{test}$ with each $\mu_k$ and choose label of the closest!

- $Y_{pred} = \text{argmin}_k ||X_{test} - \mu_k||_2$  [K comparisons instead of N]