

# IAML Past Paper 2015/2016

## I. K - Nearest Neighbors

Apply KNN to classify tweets to politics/not politics

### a.i.Tweet representation

Represent tweets as a 'BAG OF WORDS'

create vector, each column = words rows = documents

attributes: unique words

remove stop words or assign TF-IDF weights.

	F	B	C	D	E
Doc 1	1	0	6	9	10
Doc 2	2	5	7	0	3

$$\text{dist}(\text{Doc1}, \text{Doc2}) = \sqrt{(1 - 2)^2 + (0 - 5)^2 + \dots + (10 - 3)^2}$$

## ii. KNN Algorithm

measure distance with all dogs

e.g.  $K = 3$

dist (x,

dst cx,

did't Cx,

dist  $C_x$ ,

$\text{dist}(x, \text{Doc}1) = 1$   
 $\text{dist}(x, \text{Doc}2) = 3$   
 $\text{dist}(x, \text{Doc}3) = 2$   
 $\text{dist}(x, \text{Doc}4) = 5$

} nearest neighbors, check  
the NNS labels

### iii. Decide K

Cross-validation, run on multiple  $K$ s, check which one's best

#### iv. Distance function

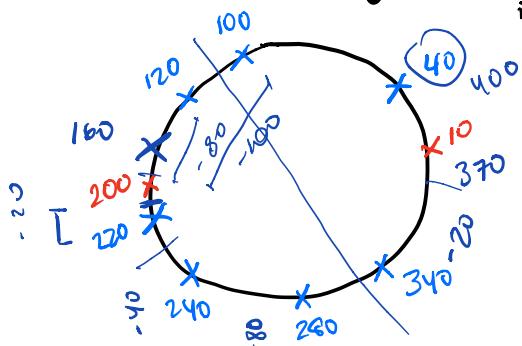
$$d(D_1, D_2) = \sqrt{\sum_{i=1}^n (x_i - y_i)^2}$$

#### v. Using K-D Trees?

No, because K-D trees are for low-dimensional real valued data and we have discrete data.

## 2. K-means clustering

- i. Describe how k-Means can be used
    - Place centroids
    - repeat until convergence
      - each pt.  $x_i$ :
        - find nearest centroid
        - assign point  $i$  to cluster  $j$
        - $j$  -- new centroid = mean of all pts.  $x_i$
    - stop when  $j$  does not change



iii. K-means

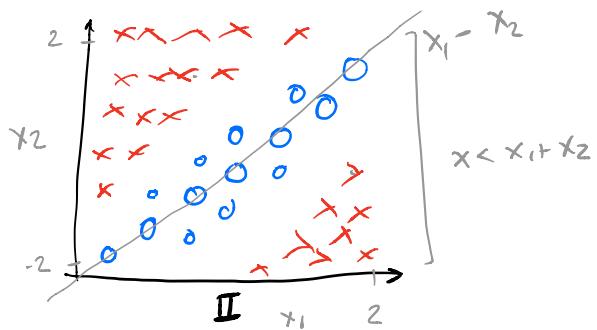
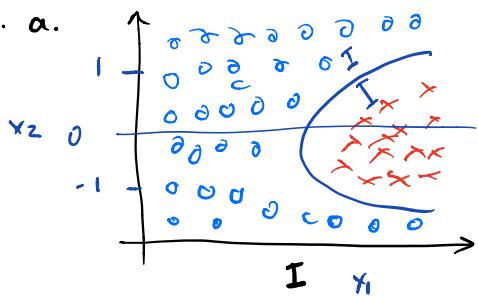
$$\text{centre} = 10: \frac{40, 340, 100}{3} = \cancel{400} 40 \quad \frac{340 + 400 + 460}{3} = 400 \approx 40$$

$$\text{centre} = 200 \quad \frac{120 + 220 + 240 + 280}{4} = 215$$

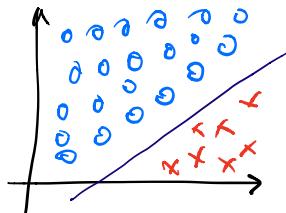
iii. How would you use K-means for hierarchy

- get clusters from each data
  - repeat K-means on each cluster
- WHAT IS THE DIFFERENCE?

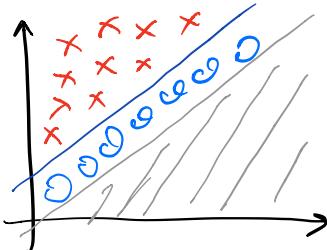
2. a.



i. For I  $\rightarrow$  square  $x_2$ , and then draw a line



ii. For II  $\rightarrow (x_1, x_2) \rightarrow (x_2, x_1)$  for all  $x < x_1 - x_2$



b. PCA  $\rightarrow x_1 = [1 2 3 4]$

$$x_2 = [5 6 7 8]$$

$$x_3 = [8 0 1 2]$$

$$x_4 = [1 6 9 8]$$

$$x_5 = [1 7 8 9]$$

- i. Describe how PCA can be used to transform the matrix  $X$  into a lower-dimensional equivalent  $Z$ .

1. Center the points  $\rightarrow$  subtract column means.
2. compute covariance matrix

$$\text{cov}(h_j, v_i) = \frac{1}{n} \sum_{i=1}^n h_j v_i$$

3. eigenvectors and eigenvalues

4. pick  $m < d$  eigenvectors w/ highest eigenvalues
5. project data points to those eigenvectors

$$b_j (x - \bar{x})^T \leq j \text{ for } j \dots m$$

ii. eigenvectors  $\rightarrow e_1 = [0.00 \ 0.35 \ 0.44 \ -0.82]$   $\lambda_1 = 0.02$   
 $e_2 = [0.40 \ -0.76 \ 0.64 \ 0.02]$   $\lambda_2 = 0.20$   
 $e_3 = [0.89 \ 0.26 \ 0.25 \ 0.26]$   $\lambda_3 = 6.44$   
 $e_4 = [-0.44 \ 0.47 \ 0.57 \ 0.50]$   $\lambda_4 = 38.83$

coordinates of  $x_i$  on  $e_1$  and  $e_2$

?

$$\begin{matrix} 1 & 2 & 3 & 4 \end{matrix} \begin{bmatrix} 0.00 & 0.40 \\ 0.35 & -0.76 \\ 0.44 & 0.64 \\ -0.82 & 0.02 \end{bmatrix} = \begin{bmatrix} x, y \end{bmatrix} \rightarrow [-1.26 \ 0.88]$$

$1 \times 4$        $4 \times 2$

$$x = 2 \times 0.35 + 3 \times 0.44 + 4 \times -0.82 \approx -1.26$$

$$y = 0.4 + 2 \times -0.76 + 3 \times 0.64 + 4 \times 0.02 \approx 0.88$$

- iii. How to select dimensions in PCA?

sort eigenvectors such that  $\lambda_1 > \lambda_2 > \dots > \lambda_d$ .

pick first  $m$  eigenvectors which explain 95% of the total variance

i.e. eigenvalue  $\lambda_i$  = variance along  $e_i$

$$\text{cum. var} \quad 38.83, 6.44, 0.20, 0.02 \quad \text{total} = 45.49$$

$$\underbrace{38.83, 6.44}_{2 \text{ dimensions}}, 0.20, 0.02 \quad 95\% = 43.21$$

### 3. Generalization and over-fitting

#### Linear Regression

Goal = predict a real valued output  $y$  for instance vector  $\underline{x}$   
via a weight vector  $\underline{w}$ .

Given = training set  $(\underline{x}_i, y_i)$

Objective Function: Mean Squared Error

- i. Difference between generalization error and training error

\* gen error  $\rightarrow$  error we will make on future data

+ training error  $\rightarrow E_{\text{train}} = \frac{1}{n} \sum_{i=1}^n \text{error}(f_p(\underline{x}_i), y_i)$

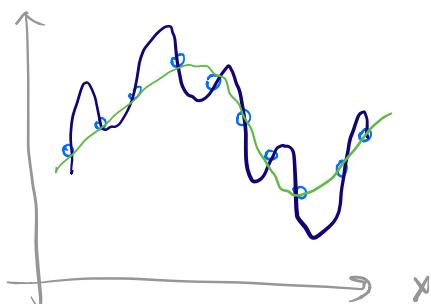
$$E_{\text{gen}} = \int \text{error}(f_p(\underline{x}), y) p(y, \underline{x}) d\underline{x}$$

- ii. Equation for generalization error

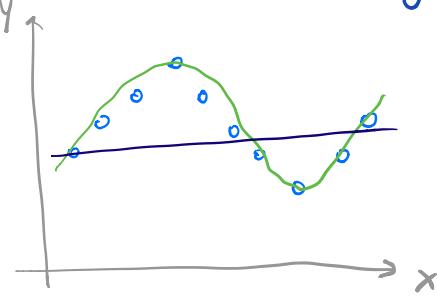
$$E_{\text{gen}} = \underbrace{\int \text{error}(f_p(\underline{x}), y) \underbrace{p(y, \underline{x}) d\underline{x}}}_{\substack{\text{error as before} \\ \text{how often we see such } \underline{x} \text{ and } y}}$$

- iii. Overfitting and underfitting in this case

makes more mistakes on training data but fewer on unseen data.  
can find another predictor  $f'$  w/ smaller  $E_{\text{tr}}$  &  $E_{\text{gen}}$



overfitting



underfitting

#### b. Evaluation and ROC Curves

?

$$\text{TPR} = \frac{\# \text{correct} (+)}{\# \text{correct} (+ \& -)} \quad \text{FPR} = \frac{\# (-) \text{incorrect}}{\# (-) \text{total}} \quad \text{ACC} = \frac{\# \text{correct}}{\# \text{total}}$$

- i. simple strategy w/ high accuracy

$$\text{ACC} = \frac{\text{TP} + \text{TN}}{\text{TP} + \text{TN} + \text{FP} + \text{FN}}$$

go with the majority of the training data labels  
(the one w/ highest prior). TPR high, FPR = 0

# IAML 2016 with solutions

## i. K-Nearest Neighbors

### i. 'Tweets' representation

- use a 'Bag of words' representation
- the tweets will be represented as binary (or integer-valued) vectors over the vocabulary, where the vocab = all words in training data
- there will be a separate numeric attr. for every word in vocab.
  - if the word is absent in a tweet = 0
  - present in a tweet = 1
  - \* or frequency if using integer-value

### ii. Describe KNN algorithm

- computes a distance acc. to a 'distance function', bet. new tweet and each training tweet, and assigns that distance to the training tweet.
- Then, it takes K trn. instances which have the smallest distance and counts whether more of them are labelled as 'politics' or 'not politics'
- tweet label = most frequent label
- if both labels equally frequent, do tie-breaker
  - random coin toss
  - most frequent label in dataset
  - use 1NN.

### iii. What distance function would you use?

Euclidean distance.

$$D(x, y) = \sqrt{\sum_{v \in V} (x_v - y_v)^2}$$

for tweets  $x$  and  $y$ , where  $V$  is the vocabulary.

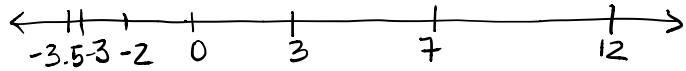
### iv. Selecting the optimal value of $K$ .

- held-out validation dataset: pick the value of  $K$  that gives the lowest classification error over the validation dataset.

### v. Can we use K-D trees to speed up this KNN?

- no
- dimensionality of tweets is too high (size of vocab)
- Inverted Indexes are a more appropriate technique

## 2. Hierarchical Agglomerative Clustering



- i. Describe how agglomerative clustering can be used to arrange an arbitrary set of datapoints.

1. Start by assigning each instance to its own cluster
2. Iteratively merge pairs of clusters into a representative cluster until we are left with a single cluster containing all the instances.
3. This single cluster is the root of the hierarchy (called a dendrogram), the two clusters that are merged into it are its children, and so on recursively.

\* The algorithm runs for  $n-1$  iterations where  $n$  is the no. of datapoints.

At each iteration we merge the two clusters that have the smallest distance between them, according to some distance metric.

- ii. Single link - defines the distance between two closest elements of the two clusters  $C_1$  and  $C_2$  as the dist. between two closest elements of the two clusters.  
How about complete link and average link?

Complete link - distance between two clusters is the distance between the farthest elements in each cluster,  
i.e.  $\max_{x_i \in C_1, x_j \in C_2} D(x_i, x_j)$   
e.g. Euclidean distance

Average link - mean of all pairwise distances between pairs of points where one point is on one cluster and the other point in the other.

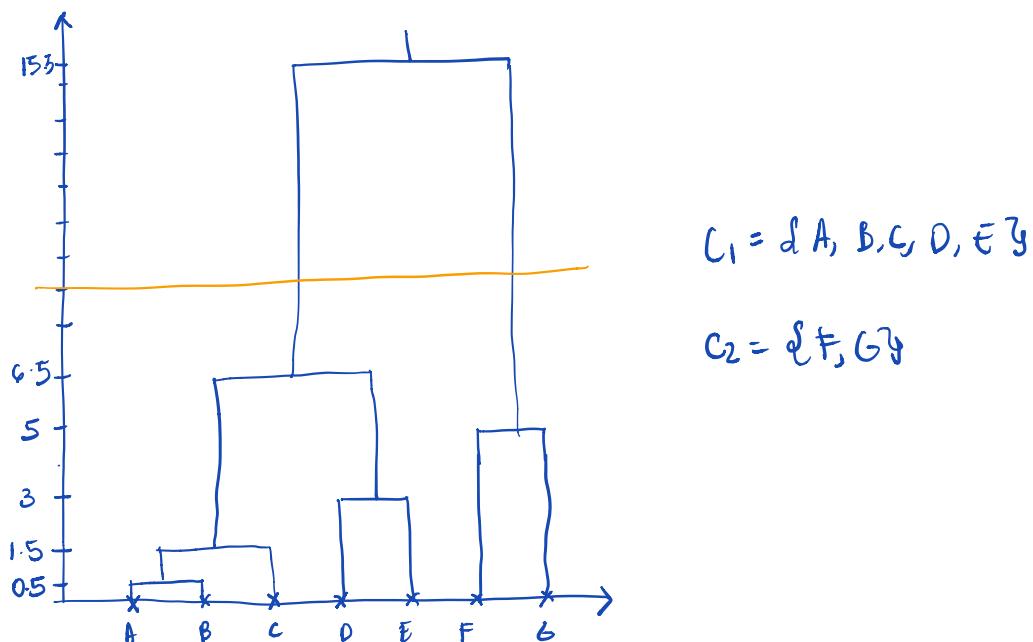
$$\frac{1}{n_1 n_2} \sum_{x_i \in C_1, x_j \in C_2} D(x_i, x_j)$$

iii. Run single link on dataset and complemlink.

single link



complete link



- \* Single link has a tendency to cluster instances to long "chains", adding elements one-at-a-time to a large cluster.
- \* Complete-link  $\rightarrow$  "spherical" clusters, where the diameter of each cluster (dist. bet. 2 dissimilar points in the cluster) is comparable among all clusters.
  - $\rightarrow$  more balanced: join all singletons to many small clusters, then proceeding to merge them to larger ones.