# Exam II, Fall 2009
# Foundation of Data Mining Course
# CSEE Department, University of Maryland Baltimore County

*Note: Closed book exam.*
*Time: 1 hour and 15 minutes*

Sandipan dey

# 1. Orthogonal Transformations: [4+3+3+10+10=30]

a. What kind of basis functions are used in Fourier transformation? How are they different from the ones used in case of PCA?

Fourier transforms can be represented by

$$F_w = \frac{1}{\sqrt{n}} \sum_{t=0}^{n-1} f_t e^{-\frac{2\pi i}{n} wt} \quad \text{(unitary version)}$$

The basis function is: $e^{-\frac{2\pi i}{n} wt}$ and the Fourier coefficients are expressed as a linear combination of the basis function (orthogonal)

In case of PCA, our problem is to find

$$\arg\max (\text{var}(w^T x)),$$
subject to $ww^T = 1$, the

solution of the optimization problem becomes $wCw^T$, where $W$ represents set of orthonormal eigenvectors of the covariance matrix $C$.

Here, the set of eigenvectors $W$ represent the orthonormal basis functions, since any data point can be represented as linear combination of the set of eigenvectors $\sum_{i=1}^{n} \lambda_i W_i$, with $W_i$ being linearly independent set of vectors (Basis of $R^n$)

b. Name one major difference between the Fourier and Wavelet transformations.

Fourier transform only expresses (by means of Fourier coefficients) a frequency domain representation of the data, but it does not represent the time domain variations. It is not good for capturing spikes or sharp changes in data. It does not take care of localizations.

Wavelet transform represents both time and frequency variations simultaneously. By means of Mother wavelet's as a generating function and small building blocks the entire signal can be generated. Also, it performs better when capturing sharp changes in data

Wavelet's capture forest and trees in the data.
Also, takes care of the localizations of data, by having

• Mother ~~wint~~ wavelets as building blocks and ~~scaling~~ using scaling and
translation to shrink/expand and shift the blocks respective

c. If G denotes the estimator of a population parameter H then write down the expression for the bias of the estimator.

*Bias of an estimator is given by the following:*

$$B_\theta(G) = E_\theta[G - \theta] = E_\theta[G] - E_\theta[\theta] \quad \text{what is } \theta?$$

$$= E_\theta[G] - \theta \quad (\theta \text{ being a constant})$$

-1

*But* $E_\theta(G) = \theta \Rightarrow B_\theta(G) = 0$

$$\Rightarrow G \text{ is an unbiased estimator of } \theta.$$

d. Let $\bar{x}, j \in \{0,1\}^l$, $f_1 : \{0,1\}^l \to R$, and $f_2 : \{0,1\}^l \to R$ where $R$ denotes the real numbers. Consider the following transformations:

$$f_1(\bar{x}) = \sum_j w_{1,j} \Psi_j(\bar{x}); \qquad f_2(\bar{x}) = \sum_j w_{2,j} \Psi_j(\bar{x})$$

where $w_{1,j}$-s and $w_{2,j}$-s are real valued coefficients. $\Psi_j(x)$-s define a set of orthonormal functions. In other words

$$\sum_x \Psi_j(\bar{x})\Psi_i(\bar{x}) = 0 \text{ when } j \neq i \text{ and } \sum_x \Psi_j(\bar{x})\Psi_i(\bar{x}) = 1 \text{ when } i = j.$$

Prove that

$$\sum_x f_1(\bar{x})f_2(\bar{x}) = \sum_j w_{1,j} w_{2,j}$$

$$\text{L.H.S.} = \sum_{\bar{x}} f_1(\bar{x})\, f_2(\bar{x})$$

$$= \sum_{\bar{x}} \sum_{\hat{\jmath}} w_{1,\hat{\jmath}}\, \varphi_{\hat{\jmath}}(\bar{x}) \sum_{\dot{\jmath}} w_{2,\dot{\jmath}}\, \varphi_{\dot{\jmath}}(\bar{x})$$

$$= \sum_{\bar{x}} \sum_{i,\hat{\jmath}} w_{1,\hat{\jmath}}\, \varphi_{\hat{\jmath}}(\bar{x})\, w_{2,i}\, \varphi_{i}(\bar{x})$$

$$= \sum_{\bar{x}} \sum_{i,\hat{\jmath}} w_{1,\hat{\jmath}}\, w_{2,i}\, \varphi_{\hat{\jmath}}(\bar{x})\, \varphi_{i}(\bar{x})$$

$$= \sum_{i,\hat{\jmath}} w_{1,\hat{\jmath}}\, w_{2,i} \sum_{\bar{x}} \varphi_{\hat{\jmath}}(\bar{x})\, \varphi_{i}(\bar{x}) \qquad \left(\begin{array}{l}\text{since } w_{1,\hat{\jmath}} \text{ and} \\ w_{2,i} \text{ are} \quad \text{scalars, not dependent} \\ \text{on } \bar{x}\end{array}\right)$$

$$= \sum_{i,\hat{\jmath}} w_{1,\hat{\jmath}}\, w_{2,i} \left( \sum_{\substack{\bar{x} \\ i \neq \hat{\jmath}}} \varphi_{\hat{\jmath}}(\bar{x})\, \varphi_{i}(\bar{x}) + \sum_{\substack{\bar{x} \\ i = \hat{\jmath}}} \varphi_{\hat{\jmath}}(\bar{x})\, \varphi_{i}(\bar{x}) \right)$$

$$= \sum_{i,\hat{\jmath}} w_{1,\hat{\jmath}}\, w_{2,\hat{\jmath}} \left( 0 + 1 \right) \qquad \left(\begin{array}{l}\text{By orthonormal} \\ \text{condition}\end{array}\right)$$

$$\left(\begin{array}{l}\text{Since non-zero only} \\ \text{for } \hat{\jmath} = i\end{array}\right)$$

$$= \sum_{i,\hat{\jmath}} w_{1,\hat{\jmath}}\, w_{2,\hat{\jmath}}$$

$$= \text{R.H.S.} \qquad (\text{Proved})$$

(10)

e. Consider the following basis function: $\Psi_j(\bar{x}) = (-1)^{\bar{j}\cdot\bar{x}}$ for question 1 (d). Write down the value of all the four coefficients $(w_{00}, w_{01}, w_{10}, w_{11})$ for the following two bit function:

f(00)=0; f(01)=1; f(10)=1; f(11)=2

$\Psi_{\bar{j}}(\bar{x}) = (-1)^{\bar{j}\cdot\bar{x}}$ and we know that

$w_{\bar{j}} = \frac{1}{2^n} \sum_{\bar{x}} f(\bar{x})\, \varphi_{\bar{j}}(\bar{x}),$ & $f(\bar{x}) = \begin{bmatrix} f(0,0) \\ f(0,1) \\ f(1,0) \\ f(1,1) \end{bmatrix} = \begin{bmatrix} 0 \\ 1 \\ 1 \\ 2 \end{bmatrix}$

here $n=2$.

$\therefore W_{00} = \frac{1}{2^2} \sum_{\bar{x}} f(\bar{x})\, (-1)^{\overline{00}\cdot\bar{x}}$, let $\bar{x} = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}$

$= \frac{1}{4} \sum_{\bar{x}} f(\bar{x})\cdot(-1)^0$   $\left(\because [0.0].\begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = 0\right)$

$= \frac{1}{4} \sum_{\bar{x}} f(\bar{x}) = \frac{1}{4}(0+1+1+2) = 1$

$W_{01} = \frac{1}{4} \sum_{\bar{x}} f(\bar{x})\cdot(-1)^{\overline{01}\cdot\bar{x}}$   $\left(\because [0\;1].\begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = x_2\right)$

$= \frac{1}{4}\left(f(00)\cdot(-1)^{[01][00]^T} + f(01)\cdot(-1)^{[01][01]^T} + f(10)\cdot(-1)^{[01][10]^T} + f(11)\cdot(-1)^{[01][11]^T}\right)$

$= \frac{1}{4}\left(0 + 1\cdot(-1)^1 + 1\cdot(-1)^0 + 2\cdot(-1)^1\right)$

$= \frac{1}{4}(-1+1-2) = -\frac{1}{2}$

$W_{10} = \frac{1}{4}\left(f(00)\cdot(-1)^{[10][00]^T} + f(01)\cdot(-1)^{[10][01]^T} + f(10)\cdot(-1)^{[10][10]^T} + f(11)\cdot(-1)^{[10][11]^T}\right)$

$= \frac{1}{4}\left(0 + 1\cdot(-1)^0 + 1\cdot(-1)^1 + 2\cdot(-1)^1\right)$

$= \frac{1}{4}(1-1-2) = -\frac{1}{2}$

$W_{11} = \frac{1}{4}\left(0 + 1\cdot(-1)^1 + 1\cdot(-1)^1 + 2\cdot(-1)^2\right)$

$= \frac{1}{4}(0-1-1+2) = 0$

10

1. Choose $\bar{c}_i$, $\forall i = 1, 2, .., K$

2. $j \leftarrow \arg \min_{i \in \{1, 2, .., K\}} d(x, \bar{c}_i)$, $\forall x \in D$,

   assign $x$ to $C_j$

   ( $d$ is distance metric, e.g. Euclidian $d$ )

3. $\bar{c}_i \leftarrow \dfrac{1}{|C_i|} \sum_{x_i \in C_i} x_i$     (compute changed cluster centers)

4. Go to step 2, if $(\exists i) \mid c_i' \neq c_i$

---

$$f(\bar{x}) = \sum_{j} w_j \, \varphi_j(\bar{x}).$$

$$w_j = \frac{1}{2^n} \sum_{\bar{x}} f(\bar{x}) \, \varphi_j(\bar{x})$$

$$\varphi_j(\bar{x}) = (-1)^{\bar{j} \cdot \bar{x}}$$

$$f: X^n \to R$$

$$w_j \in R$$

## 2. Clustering: [4+3+3+3+2=15]

**a.** How does k-means clustering work? Write down the main steps.

<u>Steps</u>

K-Means Clustering (Input & parameter 'K': #clusters)

1. Randomly select $k$ centroids for the $k$ clusters from the data points, select $\bar{C}_1, \bar{C}_2, ..., \bar{C}_K$ (cluster centers)

2. For each data point find the distance (e.g. using Euclidian distance metric) from each of the clusters and find the minimum distance, place the point in the cluster with minimum distance, ~~data for data pt~~ $\forall x \in D$,

   $j \leftarrow \arg\min\limits_{i \in \{1...k\}} d(x, \bar{C}_i)$, $j$ will be the cluster index ~~which to which~~ which $x$ is to be assigned   (Maximize intra-cluster similarity)

3. Recompute the mean for each of the clusters and assign new centroids for each cluster, $(\bar{C}_1', \bar{C}_2', ..., \bar{C}_K')$

   $\left( \bar{C}_i' \leftarrow \frac{1}{m_i} \sum\limits_{x_i \in C_i} x_i, \right)$

4. Go to step 1 until no change, i.e., convergence. Stop if $\bar{C}_i = \bar{C}_i'$, $\forall i \in \{1, 2, ..., K\}$ and output clusters

**b.** Identify two major problems of the k-Means Clustering.

<u>Problems of k-means clustering</u>

① ~~8~~ Since the centroids are ~~mean~~ calculated by the statistic 'mean', since mean is very sensitive to extreme points (outliers), the ~~initial choice~~ method is <u>not very robust to outliers</u>. The initial choice of cluster-centers must be done very carefully for this reason.

② ~~Since approaches~~ Can't <u>detect arbitrary-shaped clusters</u>, clusters tend to be ~~near~~ of nearly symmetric shape. Also, dependent on the ~~order of points~~ ordering of the data points. Also, <u>number of clusters are needed to be specified initially</u>.

c. What is the time complexity of the k-Means Clustering?

Complexity of the k-means clustering

$$= O(Imnk)$$

Where $I$ = number of iterations for convergence

$m$ = 4 $\rightarrow$ data points

$n$ = 4 $\rightarrow$ dimensions

$k$ = 4 $\rightarrow$ clusters

Step1: $O(k)$ mee
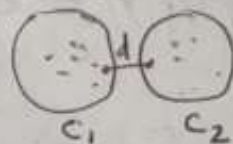Step 2: $O(mnk) \times I$ times
Step 3:

$$\overline{O(mnk\ I + k)}$$
$$= O(mnkI)$$

d. What is Single-link Hierarchical clustering? In bottom-up Hierarchical clustering ~~is a also~~ technique, while merging the lower level or smaller clusters down the tree to higher level larger clusters in the tree, if the inter-cluster distances are computed like the following:

$$d = \left\{ \min_{\substack{t \in C_i \\ t' \in C_j}} \|t - t'\| \right.$$

using some norm (e.g. $L^2$), it's called single link hierarchical clustering

$C_1$   $C_2$

e. How is the divisive hierarchical clustering different from agglomerative hierarchical clustering?

As the name suggests, agglomerative hierarchical clustering starts with a bottom-up approach, i.e., it starts with every data point as a ~~single~~ separate cluster and goes on merging the clusters hierarchically to ~~minimize intra cluster~~ maximize intra-cluster similarity.

So to the contrary, divisive cluster follows a top-down approach, starting from the entire data set and then partitioning it in hierarchical manner.