# Clustering (contd), and Dimensionality Reduction

CS771: Introduction to Machine Learning

# Soft Clustering

- If clusters overlap, doing a "soft" clustering is more desirable

- Instead of hard assignment to a cluster, report cluster membership <u>probabilities</u>

Hard Clustering (assuming K=5)

| 0 | 1 | 0 | 0 | 0 |
|---|---|---|---|---|

$z_n$ is a one-hot vector (only one $z_{nk}$ is 1)

Soft Clustering (assuming K=5)

| 0.04 | 0.65 | 0.15 | 0.10 | 0.06 |
|------|------|------|------|------|

$z_n$ is a probability vector ($\sum_{k=1}^{K} z_{nk} = 1$)

- Several methods exist for soft clustering such as latent variable models like Gaussian mixture models (will see later), and heuristics such as "soft" $K$-means

Hard K-means: each iteration

A heuristic to convert distances into the probability of $\boldsymbol{x_n}$ belonging to the $\boldsymbol{k^{th}}$ cluster

Soft/Fuzzy K-means: each iteration

$$\boldsymbol{z_n} = \mathrm{argmin}_k \|\boldsymbol{x_n} - \boldsymbol{\mu_k}\|^2 \quad \forall n$$

Only inputs with $z_{nk} = 1$ contribute to $\boldsymbol{\mu_k}$

$$z_{nk} = \frac{\exp(-\|\boldsymbol{x_n} - \boldsymbol{\mu_k}\|^2)}{\sum_{\ell=1}^{K} \exp(-\|\boldsymbol{x_n} - \boldsymbol{\mu_\ell}\|^2)} \quad \forall n, k$$
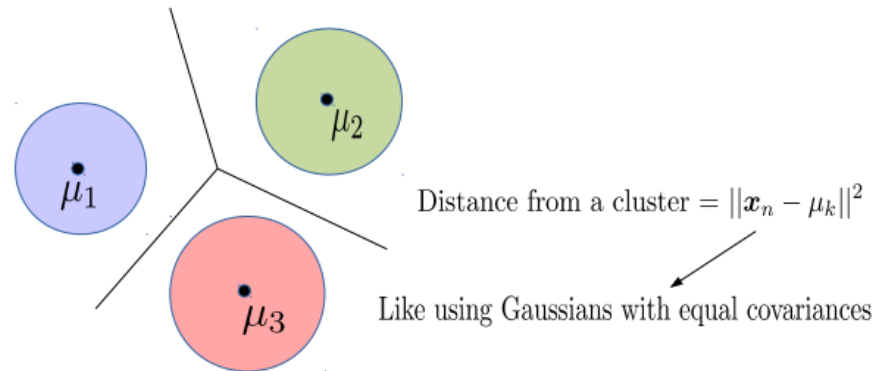
All $\boldsymbol{N}$ inputs contribute to $\boldsymbol{\mu_k}$ since $z_{nk} > 0 \ \forall k$

$$\boldsymbol{\mu_k} = \frac{1}{N_k} \sum_{n:z_n=k} \boldsymbol{x_n} = \frac{\sum_{n=1}^{N} z_{nk}\boldsymbol{x_n}}{\sum_{n=1}^{N} z_{nk}} \quad \forall k$$

$$\boldsymbol{\mu_k} = \frac{\sum_{n=1}^{N} z_{nk}\boldsymbol{x_n}}{\sum_{n=1}^{N} z_{nk}} \quad \forall k$$

# *K*-means: Some Other Limitations

- *K*-mean assumes that the decision boundary between any two clusters is linear
- Reason: The *K*-means loss function implies assumes equal-sized, spherical clusters
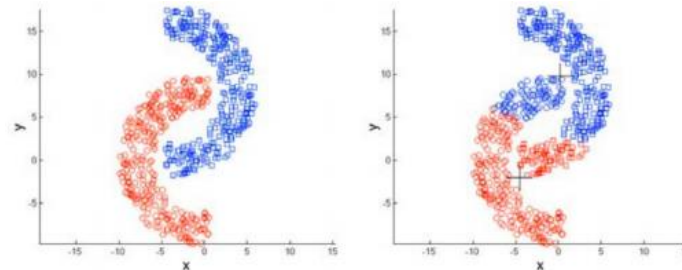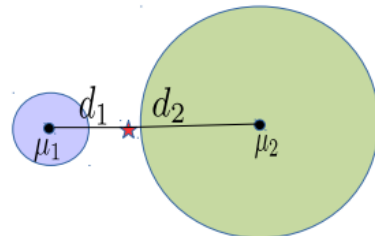


Distance from a cluster $= \|\boldsymbol{x}_n - \mu_k\|^2$

Like using Gaussians with equal covariances

Reason: Use of Euclidean distances

Some of these issues can be addressed using probabilistic models for clustering (like mixture models) or using kernels

- May do badly if clusters are not roughly equi-sized and convex-shaped

# Kernel *K*-means

Helps learn non-spherical clusters and nonlinear cluster boundaries

- Basic idea: Replace the Eucl. distances in $K$-means by the kernelized versions

Kernelized distance between input $\boldsymbol{x}_n$ and mean of cluster $k$

$$||\phi(\boldsymbol{x}_n) - \phi(\boldsymbol{\mu}_k)||^2 \quad = \quad ||\phi(\boldsymbol{x}_n)||^2 + ||\phi(\boldsymbol{\mu}_k)||^2 - 2\phi(\boldsymbol{x}_n)^\top \phi(\boldsymbol{\mu}_k)$$

- Here $k(.,.)$ denotes the kernel function and $\phi$ is its (implicit) feature map

- Note: $\phi(\mu_k)$ is the mean of $\phi$ mappings of the data points assigned to cluster $k$

<u>Not</u> the same as the $\phi$ mapping of the mean of the data points assigned to cluster $k$

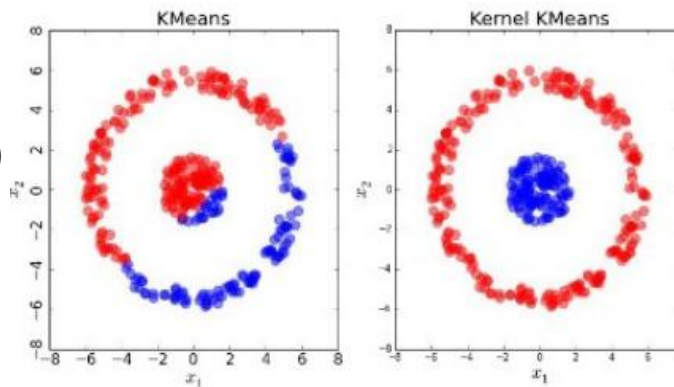$$\phi(\boldsymbol{\mu}_k) = \frac{1}{|\mathcal{C}_k|} \sum_{n:z_n=k} \phi(\boldsymbol{x}_n)$$

Can also used landmarks or kernel random features idea to get new features and run standard k-means on those

$$||\phi(\boldsymbol{\mu}_k)||^2 = \phi(\boldsymbol{\mu}_k)^\top \phi(\boldsymbol{\mu}_k)$$

$$= \frac{1}{|\mathcal{C}_k|^2} \sum_{n:z_n=k} \sum_{n:z_m=k} k(\boldsymbol{x}_n, \boldsymbol{x}_m)$$

$$\phi(\boldsymbol{x}_n)^\top \phi(\boldsymbol{\mu}_k) = \frac{1}{|\mathcal{C}_k|} \sum_{m:z_m=k} k(\boldsymbol{x}_n, \boldsymbol{x}_m)$$

Note: Apart from kernels, it is also possible to use other distance functions in $K$-means. Bregman Divergence* is such a family of distances (Euclidean and Mahalanobis are special cases)


KMeans


Kernel KMeans

*Clustering with Bregman Divergences (Banerjee et al, 2005)

# Overlapping Clustering

- Have seen hard clustering and soft clustering

- In hard clustering, $z_n$ is a one-hot vector

- In soft clustering, $z_n$ is a vector of probabilities

- Overlapping Clustering: A point can <u>simultaneously</u> belong to multiple clusters
    - This is different from soft-clustering
    - $z_n$ would be a binary vector, rather than a one hot or probability vector, e.g.,

$$z_n = [1\ 0\ 0\ 1\ 0]$$

- In general, more difficult than hard/soft clustering (for $N$ data points and $K$ clusters, the size of the space of possible solutions is not $K^N$ but $2^{NK}$ - exp in both $N$ and $K$)

- K-means has extensions* for doing overlapping clustering. There also exist latent variable models for doing overlapping clustering

*An extended version of the k-means method for overlapping clustering (Cleuziou, 2008); Non-exhaustive, Overlapping k-means (Whang et al, 2015)

# Clustering as Matrix Factorization

Also an unsupervised learning problem

- Clustering (hard, soft, overlapping) can also be posed as matrix factorization



$$N \underset{X}{\overset{D}{\boxed{X}}} \approx N \underset{Z}{\overset{K}{\boxed{Z}}} \; K \underset{\mu}{\overset{D}{\boxed{\mu}}}$$

- Minimize the distortion $\|X - Z\mu\|_F^2$ subject to suitable constraints on $Z$, e.g.,
  - $z_n$ is a one-hot vector
  - Entries of $z_n$ are non-negative and sum to 1
  - $z_n$ is a binary vector

# Hierarchical Clustering

Similarity between two clusters (or two set of points) is needed in HC algos (e.g., this can be average pairwise similarity between the inputs in the two clusters)

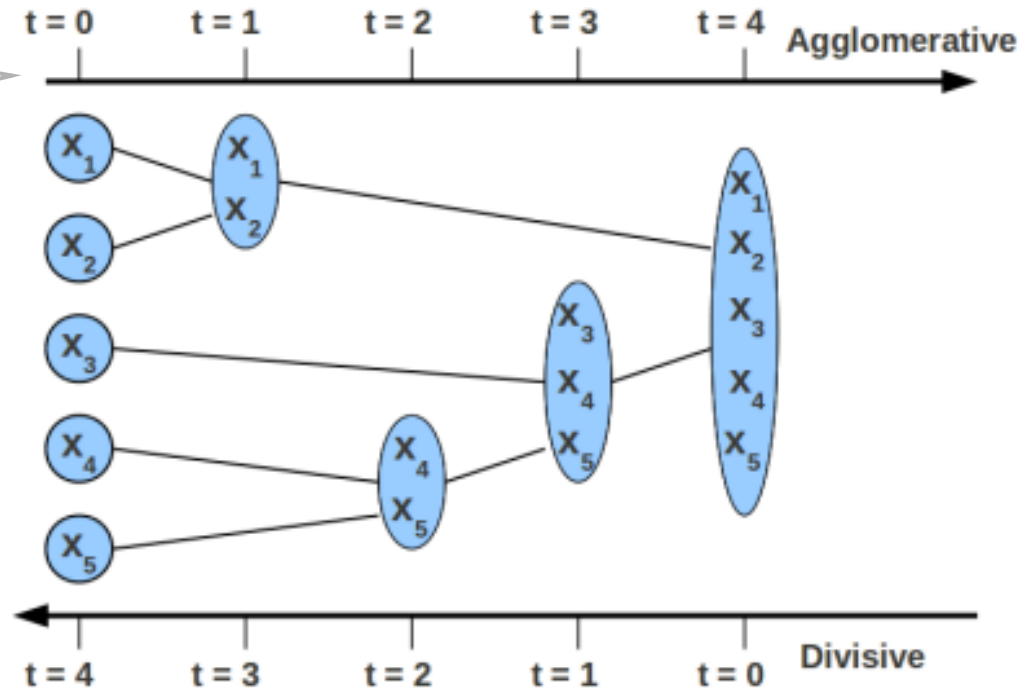▪ Can be done in two ways: Agglomerative or Divisive

Agglomerative: Start with each point being in a singleton cluster

At each step, greedily merge two most "similar" sub-clusters

Stop when there is a single cluster containing all the points

Learns a dendrogram-like structure with inputs at the leaf nodes. Can then choose how many clusters we want

Keep recursing until the desired number of clusters found

At each step, break a cluster into (at least) two smaller homogeneous sub-clusters

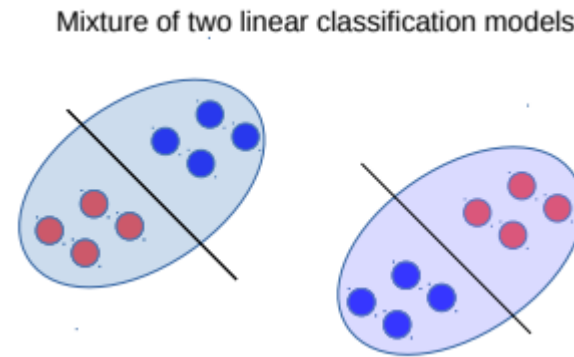Divisive: Start with all points being in a single cluster

Tricky because no labels (unlike Decision Trees)

t = 0   t = 1   t = 2   t = 3   t = 4   **Agglomerative**

$X_1$   $X_1$                           $X_1$
        $X_2$                           $X_2$
$X_2$
                                $X_3$   $X_3$
$X_3$                           $X_4$   $X_4$
                        $X_4$   $X_5$   $X_5$
$X_4$                   $X_5$
                $X_5$
$X_5$

t = 4   t = 3   t = 2   t = 1   t = 0   **Divisive**

▪ Agglomerative is more popular and simpler than divisive (the latter usually needs complicated heuristics to decide cluster splitting).
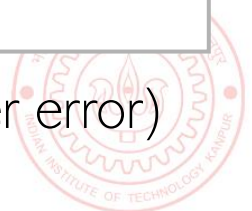
▪ Neither uses any loss function

CS771: Intro to ML

# Clustering can help supervised learning, too

- Often "difficult" sup. learning problems can be seen as mixture of simpler models

- Example: Nonlinear regression or nonlinear classification as mixture of linear models

Mixture of two linear regression models

Mixture of two linear classification models

- Don't know which point should be modeled by which linear model ⇒ Clustering

- Can therefore solve such problems as follows

  Such an approach is also an example of divide and conquer and is also known as "mixture of experts" (will see it more formally when we discuss latent variable models)

  - Initialize each linear model somehow (maybe randomly)
  - Cluster the data by assigning each point to its "closest" linear model (one that gives lower error)
  - (Re-)Learn a linear model for each cluster's data. Go to step 2 if not converged.
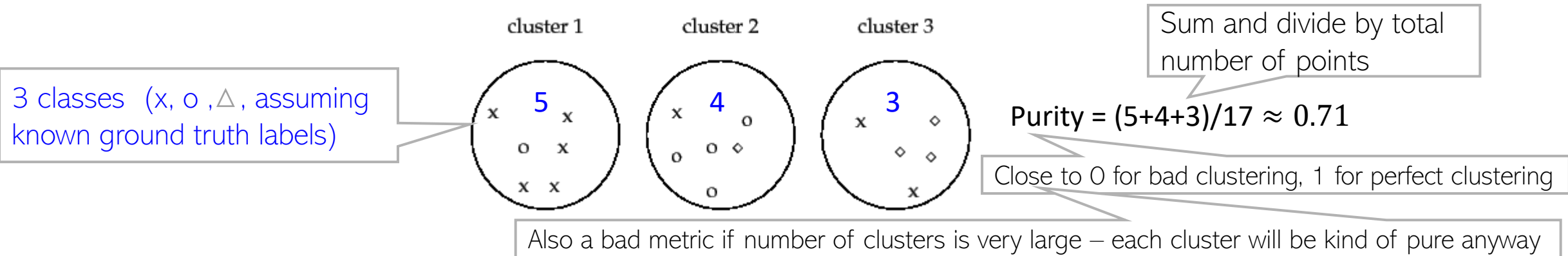
# Evaluating Clustering Algorithms

- Clustering algos are in general harder to evaluate since we rarely know the ground truth clustering (since clustering is unsupervised)

- If ground truth labels not available, use output of clustering for some other task
  - For example, use cluster assignment $z_n$ (hard or soft) as a new feature representation
  - Performance on some task using this new rep. is a measure of goodness of clustering

- If ground truth labels are available, can compare them with clustering based labels
  - Not straightforward to compute accuracy since the label identities may not be the same, e.g.,

    Ground truth = [1 1 1 0 0 0]     Clustering = [0 0 0 1 1 1]

    (Perfect clustering but zero "accuracy" if we just do a direct match)
  - There are various metrics that take into account the above fact
    - Purity, Rand Index, F-score, Normalized Mutual Information, distortion or loss on test data etc
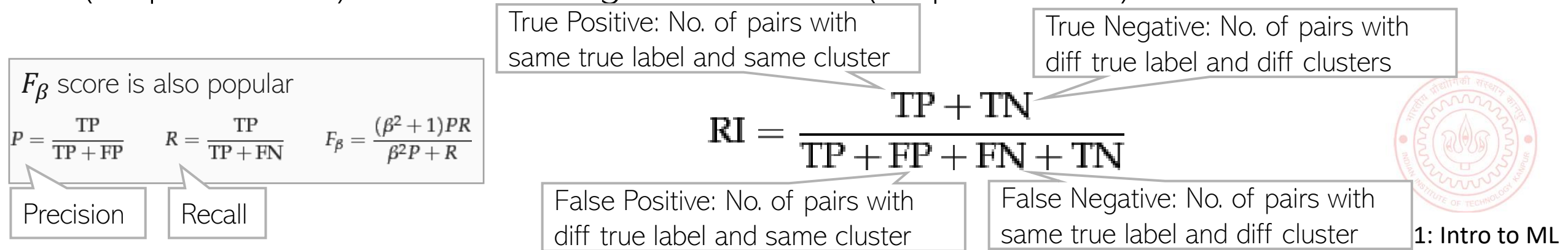
# Evaluating Clustering Algorithms

- Purity: Looks at how many points in each cluster belong to the majority class in that cluster

cluster 1    cluster 2    cluster 3

3 classes (x, o ,△, assuming known ground truth labels)

5 | 4 | 3

Sum and divide by total number of points

Purity = (5+4+3)/17 ≈ 0.71

Close to 0 for bad clustering, 1 for perfect clustering

Also a bad metric if number of clusters is very large – each cluster will be kind of pure anyway

- Rand Index (RI): Can also look at what fractions of pairs of points with same (resp. different) label are assigned to same (resp. different) cluster
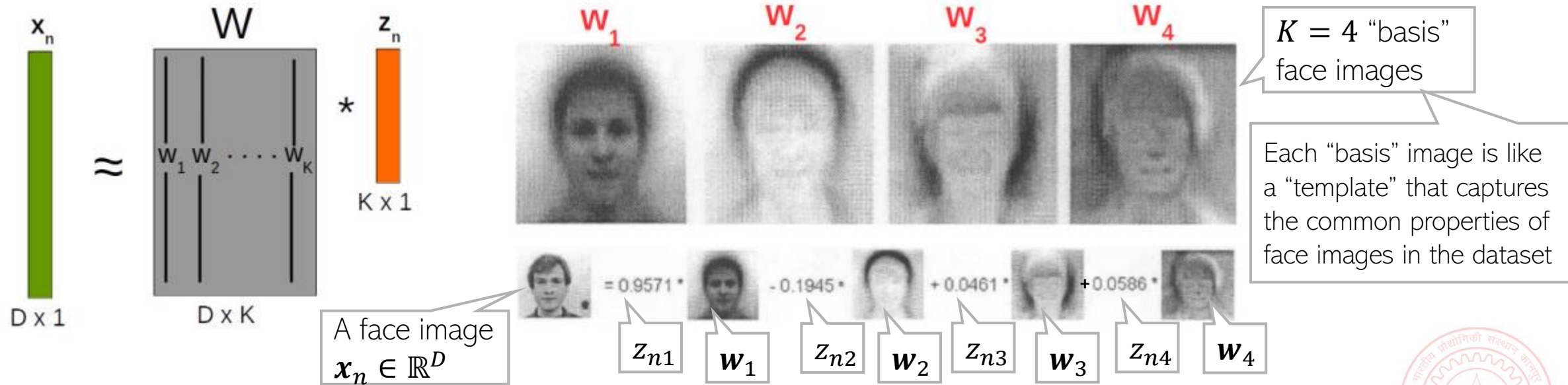
True Positive: No. of pairs with same true label and same cluster

True Negative: No. of pairs with diff true label and diff clusters

$F_\beta$ score is also popular

$$P = \frac{TP}{TP+FP} \qquad R = \frac{TP}{TP+FN} \qquad F_\beta = \frac{(\beta^2+1)PR}{\beta^2 P + R}$$

Precision    Recall

$$RI = \frac{TP+TN}{TP+FP+FN+TN}$$

False Positive: No. of pairs with diff true label and same cluster

False Negative: No. of pairs with same true label and diff cluster

1: Intro to ML

# Dimensionality Reduction: A motivative example

- Consider a linear model of the form

$$x_n \approx \widetilde{x}_n = Wz_n = \sum_{k=1}^{K} z_{nk}w_k$$

$w_k$ is the $k$-th column of $W$

- Above means that each $x_n$ is appox a linear comb of $K$ vectors $w_1, w_2, \ldots, w_K$



$K = 4$ "basis" face images

Each "basis" image is like a "template" that captures the common properties of face images in the dataset

A face image $x_n \in \mathbb{R}^D$

$= 0.9571 *$  $- 0.1945 *$  $+ 0.0461 *$  $+ 0.0586 *$

$z_{n1}$  $w_1$  $z_{n2}$  $w_2$  $z_{n3}$  $w_3$  $z_{n4}$  $w_4$
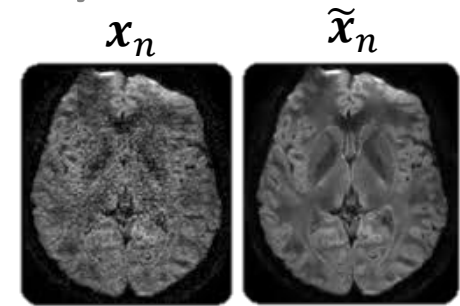
- In this example, $z_n \in \mathbb{R}^K$ ($K = 4$) is a low-dim feature rep. for each image $x_n \in \mathbb{R}^D$

# Dimensionality Reduction: More formally..

$x_n$  $\tilde{x}_n$

- Goal: Reduce the dimensionality of each input $x_n \in \mathbb{R}^D$

$z_n \in \mathbb{R}^K$ $(K \ll D)$ is a compressed version of $x_n$

$$z_n = f(x_n)$$

- Also want to be able to (approximately) reconstruct $x_n$ from $z_n$

Often $\tilde{x}_n$ is a "cleaned" version of $x_n$ (the loss in information is often the noise/redundant information in $x_n$)

$$\tilde{x}_n = g(z_n) = g(f(x_n)) \approx x_n$$

- Sometimes $f$ is called "encoder" and $g$ is called "decoder". Can be linear/nonlinear

- These functions are learned by minimizing the distortion/reconstruction error of inputs

$$\mathcal{L} = \sum_{n=1}^{N} \|x_n - \tilde{x}_n\|^2 = \sum_{n=1}^{N} \|x_n - g(f(x_n))\|^2$$

Image source: Manjon et al (2013): Diffusion Weighted Image Denoising using overcomplete Local PCA