

Machine Learning

Week 8: Unsupervised Learning (Cont'd) Support Vector Machines

Maheesan Niranjana

School of Electronics and Computer Science
University of Southampton

Autumn Semester 2015/16

Gaussian Mixture Model and K-Means Clustering

$$p(\mathbf{x}) = \sum_{k=1}^K \pi_j \mathcal{N}(\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \quad \pi_j \geq 0, \quad \sum_{j=1}^K \pi_j = 1.$$

```
Input:  $\mathbf{X} = \{\mathbf{x}_n^t\}_{n=1}^N, K$   
Output:  $\mathbf{C}, \text{Idx}$   
initialize:  $\mathbf{C} = \{\mathbf{c}_j^t\}_{j=1}^K$   
  
repeat  
  . assign  $n^{\text{th}}$  sample to nearest  $\mathbf{c}_j$   
  .  $\text{Idx}(n) = \min_j \|\mathbf{x}_n - \mathbf{c}_j\|^2$   
  
  . recompute  $\mathbf{c}_j = \frac{1}{N_j} \sum_{n=j} \mathbf{x}_n$   
until no change in  $\mathbf{c}_1, \mathbf{c}_2, \dots, \mathbf{c}_K$   
  
return  $\mathbf{C}, \text{Idx}$ 
```

Objective Function for Clustering

Setting up an error function and minimizing it

$$J_e = \sum_{i=1}^K \sum_{\mathbf{x} \in \mathcal{D}_i} \|\mathbf{x} - \mathbf{m}_i\|^2$$
$$\mathbf{m}_i = \frac{1}{n_i} \sum_{\mathbf{x} \in \mathcal{D}_i} \mathbf{x}$$

Which is also the same as (in terms of scatter)

$$J_e = \frac{1}{2} \sum_{i=1}^K n_i \bar{s}_i$$
$$\bar{s}_i = \frac{1}{n_i^2} \sum_{\mathbf{x} \in \mathcal{D}_i} \sum_{\mathbf{y} \in \mathcal{D}_i} \|\mathbf{x} - \mathbf{y}\|^2$$

Homework: Show this i.e. sum of average distance to cluster means and sum of within cluster scatter are the same.

Iterative Optimization

Note: Discrete optimization

$$J_e = \sum_{i=1}^K J_i$$
$$= \sum_{i=1}^K \sum_{\mathbf{x} \in \mathcal{D}_i} \|\mathbf{x} - \mathbf{m}_i\|^2$$

- Mean of each cluster: $\mathbf{m}_i = \sum_{\mathbf{x} \in \mathcal{D}_i} \mathbf{x}$
- Move sample (data) $\hat{\mathbf{x}}$ from cluster \mathcal{D}_i to \mathcal{D}_j ; Say new J_j is J_j^* and new \mathbf{m}_j is \mathbf{m}_j^*

$$\mathbf{m}_j^* = \mathbf{m}_j + \frac{1}{n_j + 1} (\hat{\mathbf{x}} - \mathbf{m}_j)$$

$$J_j^* = \sum_{\mathbf{x} \in \mathcal{D}_j} \|\mathbf{x} - \mathbf{m}_j^*\|^2 + \|\hat{\mathbf{x}} - \mathbf{m}_j^*\|^2$$
$$= \left(\sum_{\mathbf{x} \in \mathcal{D}_j} \left\| \mathbf{x} - \mathbf{m}_j - \frac{1}{n_j + 1} (\hat{\mathbf{x}} - \mathbf{m}_j) \right\|^2 \right) = \left\| \frac{n_j}{n_j + 1} (\hat{\mathbf{x}} - \mathbf{m}_j) \right\|^2$$
$$= J_j + \frac{n_j}{n_j + 1} \|\hat{\mathbf{x}} - \mathbf{m}_j\|^2$$

Iterative optimization (cont'd)

...similarly, \mathbf{m}_i changes to

$$\mathbf{m}_i^* = \mathbf{m}_i - \frac{1}{n_i - 1}(\hat{\mathbf{x}} - \mathbf{m}_i)$$

$$J_i = J_i - \frac{n_i}{n_i - 1} \|\hat{\mathbf{x}} - \mathbf{m}_i\|^2$$

So, if

$$\frac{n_i}{n_i - 1} \|\hat{\mathbf{x}} - \mathbf{m}_i\|^2 > \frac{n_j}{n_j + 1} \|\hat{\mathbf{x}} - \mathbf{m}_j\|^2$$

then it is advantageous to move $\hat{\mathbf{x}}$ from \mathcal{D}_i to \mathcal{D}_j

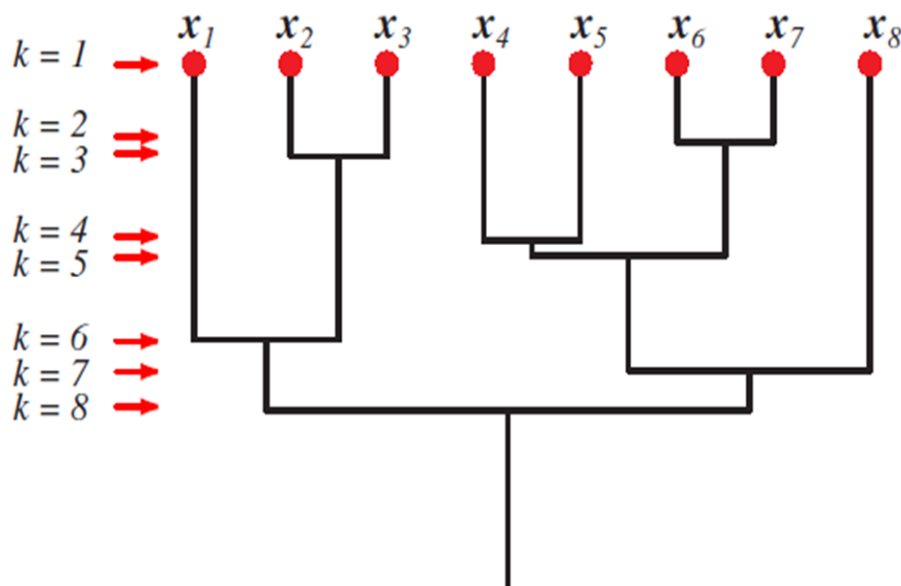
Algorithm:

- Select a data point at random
- Move it to cluster for which $\frac{n_j}{n_j + 1} \|\hat{\mathbf{x}} - \mathbf{m}_j\|^2$ is minimum.
- Recalculate means \mathbf{m}_i , $i=1, \dots, K$

This will be a sequential version of K -means algorithm; *i.e.* update at each data, rather than wait till we classify all data.

Hierarchical Clustering

Dendrogram



Homework: Use MATLAB to draw a dendrogram for the Boston Housing data

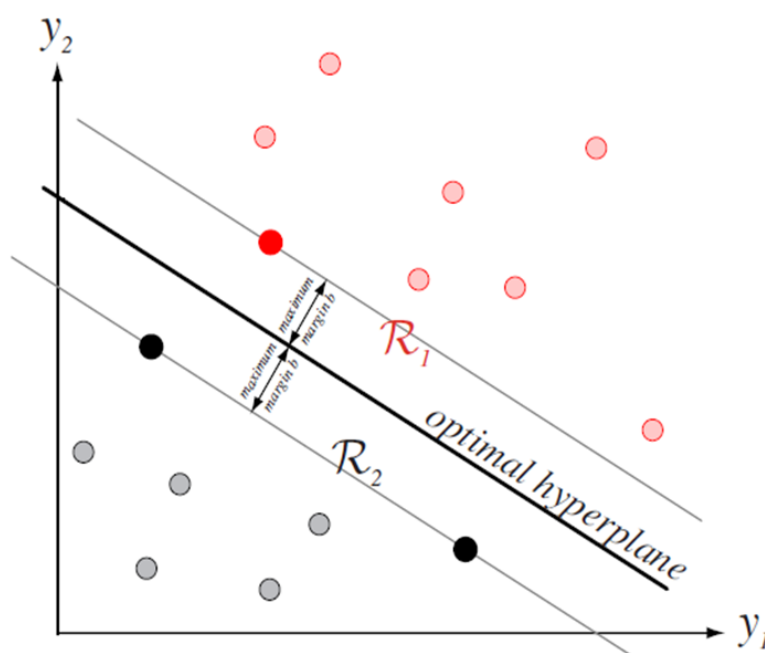
Agglomerative Hierarchical Clustering

- Initialize: $\hat{K} = n$ (No. clusters = No. data)
- Repeat (until $\hat{K} = K$)
 - find nearest clusters \mathcal{D}_i and \mathcal{D}_j
 - merge \mathcal{D}_i and \mathcal{D}_j
 - $\hat{c} \leftarrow \hat{c} - 1$

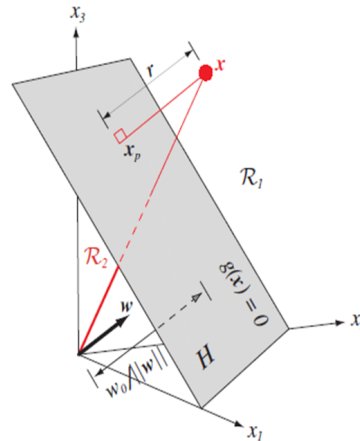
Defining nearest clusters

$$D_{\min}(\mathcal{D}_i, \mathcal{D}_j) = \min_{\mathbf{x} \in \mathcal{D}_i, \mathbf{y} \in \mathcal{D}_j} \|\mathbf{x} - \mathbf{y}\|^2$$
$$D_{\text{avg}}(\mathcal{D}_i, \mathcal{D}_j) = \frac{1}{n_i n_j} \sum_{\mathbf{x} \in \mathcal{D}_i} \sum_{\mathbf{y} \in \mathcal{D}_j} \|\mathbf{x} - \mathbf{y}\|^2$$

Now for something completely different!



Margin



(b in formula is w_0 in figure)

- Hyperplane: $\mathbf{w}^t \mathbf{x} + b = 0$ See Lab 2

vfill

- Data:

$$\mathcal{D} = \{\mathbf{x}_n, y_n\}_{n=1}^N, \mathbf{x}_n \in \mathcal{R}^d, y_n \in \{-1, +1\}$$

vfill

- Learning problem:

$$y_n [\mathbf{w}^t \mathbf{x}_n + b] \geq 1, n = 1, \dots, N$$

Margin

- Distance from data \mathbf{x}_n to a hyperplane (\mathbf{w}, b) :

$$d(\mathbf{w}, b, \mathbf{x}_n) = \frac{|\mathbf{w}^t \mathbf{x}_n + b|}{\|\mathbf{w}\|}$$

- The margin – distance between data closest to the hyperplane on either side

$$\begin{aligned} \rho(\mathbf{w}, b) &= \min_{\mathbf{x}_n: y_n = -1} d(\mathbf{w}, b, \mathbf{x}_n) + \min_{\mathbf{x}_n: y_n = +1} d(\mathbf{w}, b, \mathbf{x}_n) \\ &= \min_{\mathbf{x}_n: y_n = -1} \frac{|\mathbf{w}^t \mathbf{x}_n + b|}{\|\mathbf{w}\|} + \min_{\mathbf{x}_n: y_n = +1} \frac{|\mathbf{w}^t \mathbf{x}_n + b|}{\|\mathbf{w}\|} \\ &= \frac{1}{\|\mathbf{w}\|} \left(\min_{\mathbf{x}_n: y_n = -1} |\mathbf{w}^t \mathbf{x}_n + b| + \min_{\mathbf{x}_n: y_n = +1} |\mathbf{w}^t \mathbf{x}_n + b| \right) \\ &= \frac{2}{\|\mathbf{w}\|} \end{aligned}$$

$$\mathcal{L}(\mathbf{w}, b, \alpha) = \frac{1}{2} \|\mathbf{w}\|^2 - \sum_{n=1}^N \alpha_n (y_n [\mathbf{w}^t \mathbf{x}_n + b] - 1), \alpha_n \geq 0$$

- Setting $\frac{\partial \mathcal{L}}{\partial b}$ to zero, gives $\sum_{n=1}^N \alpha_n y_n = 0$
- Setting $\frac{\partial \mathcal{L}}{\partial \mathbf{w}}$ to zero, gives $\mathbf{w} = \sum_{n=1}^N \alpha_n y_n \mathbf{x}_n$
- Note: the unknown weights are computed as a weighted sum of the training examples; do you see a similarity to the perceptron algorithm?
- Substitute to get the dual problem

$$\max_{\alpha} -\frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \alpha_i \alpha_j y_i y_j \mathbf{x}_i^t \mathbf{x}_j + \sum_{k=1}^N \alpha_k$$

$$\min_{\alpha} \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \alpha_i \alpha_j y_i y_j \mathbf{x}_i^t \mathbf{x}_j - \sum_{k=1}^N \alpha_k$$

subject to $\alpha_n \geq 0$ and $\sum_{n=1}^N \alpha_n y_n = 0$

- Quadratic programming

MATLAB> help quadprog

$$\min_{\mathbf{x}} \frac{1}{2} \mathbf{x}^t \mathbf{H} \mathbf{x} + \mathbf{f}^t \mathbf{x}$$

Subject to

$$\begin{aligned} \mathbf{A} \mathbf{x} &\leq \mathbf{b} \\ \mathbf{A}_{\text{eq}} \mathbf{x} &= \mathbf{b}_{\text{eq}} \\ lb &\leq \mathbf{x} \leq ub \end{aligned}$$

MATLAB> `x = quadprog(H, f, A, b, Aeq, beq, lb, ub);`

Calculating the Bias Term

- Constraints $\alpha_n \geq 0$; Parameters $\mathbf{w} = \sum_{n=1}^N y_n \alpha_n \mathbf{x}_n$
- Non-zero α_n 's correspond to Support Vectors
- For any of these support vectors (\mathbf{x}_s): $y_s[\mathbf{w}^t \mathbf{x}_s + b] = 1$; we can compute the bias term b from this.

$$y_s \left[\sum_{m \in S} \alpha_m y_m \mathbf{x}_m^t \mathbf{x}_s + b \right] = 1$$

$$y_s^2 \left(\sum_{m \in S} \alpha_m y_m \mathbf{x}_m^t \mathbf{x}_s + b \right) = y_s$$

$$\text{Note : } y_s^2 = 1; \text{ Hence } b = y_s - \sum_{m \in S} \alpha_m y_m \mathbf{x}_m^t \mathbf{x}_s$$

- In practice, instead of using any one support vector, use we average:

$$b = \frac{1}{N_s} \sum_{s \in S} \left(y_s - \sum_{m \in S} \alpha_m y_m \mathbf{x}_m^t \mathbf{x}_s \right)$$