

# CS480 Assignment 4

Nikhil Nanda

20727218

## Exercise 1:

$$\rightarrow s_j = \frac{\sum_{i=1}^n r_{ik} (x_{ij} - \mu_j)^2}{\sum_{i=1}^n r_{ik}} \\ = \frac{\sum_{i=1}^n r_{ik} x_{ij}^2}{\sum_{i=1}^n r_{ik}} - \mu_j^2$$

The above equation comes from the Expectation step of the EM algorithm, where the covariance matrices  $S_k$  are constrained to be diagonal.

On fixing the responsibility  $r_{ik}$ ,

$$\min_{\theta} \sum_{i=1}^n \sum_{k=1}^K r_{ik} \left[ -\log \pi_k + \frac{1}{2} \log |S_k| + \frac{1}{2} (x_i - \mu_k)^T S_k^{-1} (x_i - \mu_k) \right]$$

On taking derivative and setting it to zero , we get that

for each  $K$  , the  $j$ -th diagonal

$$s_j = \frac{\sum_{i=1}^n r_{ik}(x_{ij} - \mu_j)^2}{\sum_{i=1}^n r_{ik}}$$

→ python3 qvl.py

## Exercise 2:

$$\begin{aligned}
 \text{error}(f) &= P(Y \neq f(x)) \\
 &= 1 - P(Y = f(x)) \\
 &= 1 - \sum_{m=1}^C P(Y=m, f(x)=m) \\
 &= 1 - \sum_{m=1}^C E\left[P(Y=m, f(x)=m | X)\right] \\
 &= 1 - \sum_{m=1}^C E\left[\mathbb{1}_{f(x)=m} P(Y=m | X)\right] \\
 &= 1 - E\sum_{m=1}^C \mathbb{1}_{f(x)=m} P(Y=m | X)
 \end{aligned}$$

$\Rightarrow$  
$$\boxed{\text{error}(f) = E\left[1 - \sum_{m=1}^C \mathbb{1}_{f(x)=m} P(Y=m | X)\right]}$$

- Thus, the minimum classification

$$\text{error}(f) = E\left[1 - \max_{m=1, \dots, C} P(Y=m | X)\right]$$

- This means that Bayes rule (with ties breaking arbitrarily)

$$f^*(x) = \operatorname{argmax}_{m=1, \dots, C} P(Y=m | X)$$

achieves the minimum classification error

### Exercise 3:

→ Pseudocode :

$\text{thresh}_{:,0} \leftarrow$  best threshold values of all features

$\text{thresh}_{:,1} \leftarrow$  error for all features with threshold value  $\text{thresh}_{:,0}$

$x_{ij} \leftarrow$  value at i-th row and j-th column of  $X$

$n \leftarrow$  number of training examples

Let  $\text{thresh}$  be a 2D-array of size  $(n+1) \times 2$

For every feature  $j$

Sort the  $n \times x_{ij}$  values

Set  $\text{sum}$  be an array of size  $n+1$

$\text{sum}[0] \leftarrow 0$

$\text{sumSqr} \leftarrow 0$

for  $i \leftarrow 1$  to  $n$

$\text{sum}[i] \leftarrow y_i + \text{sum}[i-1]$

$\text{sumSqr} \leftarrow (y_i)^2 + \text{sumSqr}$

Let  $\text{error}$  be an array of size  $n+1$   
whose elements are initialized to 0

for  $i \leftarrow 0$  to  $n$

$\mu \leftarrow \text{sum}[i] / i$

$\gamma \leftarrow (\text{sum}[n] - \text{sum}[i]) / (n-i)$

$\text{error}[i] \leftarrow \text{sumSqr} - (2 * \mu * \text{sum}[i])$

$+ (i * (\mu)^2) + ((n-i) * (\gamma)^2)$

$- (2 * \gamma * (\text{sum}[n] - \text{sum}[i]))$

index  $\leftarrow$  index of min value of error

if index == 0

$$\text{thresh}[j][0] \leftarrow x_{1,j} - 0.1$$

$$\text{thresh}[j][1] \leftarrow \text{error}[index]$$

else

$$\text{thresh}[j][0] \leftarrow x_{\text{index}, j}$$

$$\text{thresh}[j][1] \leftarrow \text{error}[index]$$

choose feature  $j$  with minimum  $\text{thresh}[j][1]$  value

$\rightarrow$  Runtime:

$O(d * (n \log n + k * (n+1))) + (n+1)$ , where  
k is a constant

$\therefore \boxed{O(d n \log n)}$

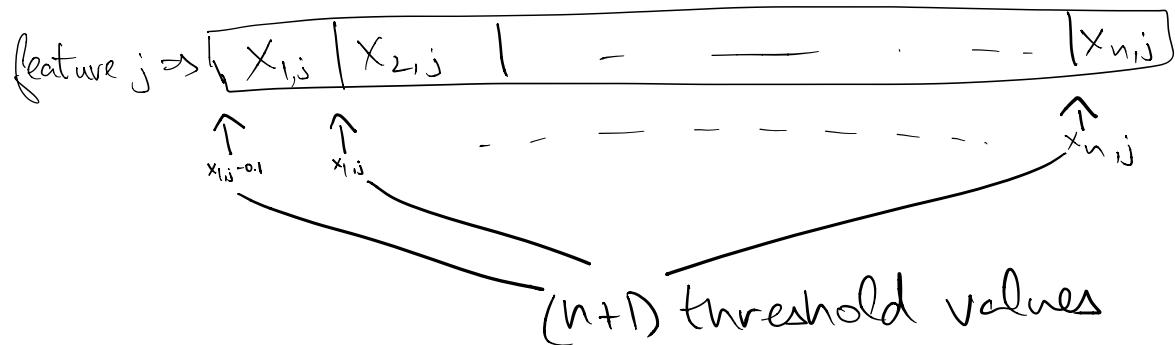
$\rightarrow$  Space Complexity:

$O(k \cdot n)$ , where k is a constant

$\Rightarrow \boxed{O(n)}$ , assuming quicksort  
is used for sorting

→ The above algorithm goes through every feature in order to decide which feature to split at the root of the regression tree. Each feature is considered independently and a best threshold value (pertaining to minimum error variance) is calculated for each feature.

The best threshold for each feature  
is the one with minimum error  $\oplus$   
variance , chosen from  $(n+1)$  threshold  
values .



The error for each  $(n+1)$  threshold values of a particular feature is computed in  $O(n)$  time by the following pre-computations

$$\text{sum}[i] \leftarrow \sum_{a=1}^i y_a$$

$$\text{sumsq} \leftarrow \sum_{i=1}^n (y_i)^2$$

$\mu, \nu$  are set to be the average of all  $y_i$  give the constraint on  $x_{ij}$

Thus,

$$\begin{aligned}
& \sum_{i: x_{ij} \leq t_j} (y_i - \mu)^2 + \sum_{i: x_{ij} > t_j} (y_i - \nu)^2 \\
= & \sum_{i: x_{ij} \leq t_j} y_i^2 - 2\mu \sum_{i: x_{ij} \leq t_j} y_i + \mu^2 \sum_{i: x_{ij} \leq t_j} 1 \\
& + \sum_{i: x_{ij} > t_j} y_i^2 - 2\nu \sum_{i: x_{ij} > t_j} y_i + \nu^2 \sum_{i: x_{ij} > t_j} 1 \\
= & \sum_{i=1}^n y_i^2 - 2\mu \sum_{i: x_{ij} \leq t_j} y_i + \mu^2 \sum_{i: x_{ij} \leq t_j} 1 \\
& - 2\nu \sum_{i: x_{ij} > t_j} y_i + \nu^2 \sum_{i: x_{ij} > t_j} 1 \\
= & \text{error for each of the different (int) threshold values for each feature}
\end{aligned}$$

Thus, the required result returned by the algorithm is the feature with the minimum error ~~or~~ variance as computed above along with the best threshold value (minimum threshold value) for that feature.

Exercise 4:

$$\min_{w \in \mathbb{R}^d} \max_{\substack{\max_{j, \|z_j\|_2 \leq \lambda} \\ \|w\|_2}} \| (x + z) w - y \|_2$$

$$\Rightarrow \min_{w \in \mathbb{R}^d} \max_{\substack{\max_{j, \|z_j\|_2 \leq \lambda} \\ \|w\|_2}} \left[ (x + z) w - y \right]^T u$$

$$\Rightarrow \min_{w \in \mathbb{R}^d} \max_{\substack{\max_{j, \|z_j\|_2 \leq \lambda} \\ \|w\|_2 \leq 1}} \left[ (x_w - y)^T u + (z_w)^T u \right]$$

$$\Rightarrow \min_{w \in \mathbb{R}^d} \max_{\substack{\max_{j, \|z_j\|_2 \leq \lambda} \\ \|w\|_2 \leq 1}} \left[ (x_w - y)^T u + \max_{j, \|z_j\|_2 \leq \lambda} (z_w)^T u \right]$$

$$\Rightarrow \min_{w \in \mathbb{R}^d} \left[ \|x_w - y\|_2 + \max_{j, \|z_j\|_2 \leq \lambda} \|z_w\|_2 \right]$$

$$\Rightarrow \min_{w \in \mathbb{R}^d} \left[ \|x_w - y\|_2 + \max_{j, \|z_j\|_2 \leq \lambda} \sqrt{\sum_j z_j w_j^2} \right]$$

$$\Rightarrow \min_{w \in \mathbb{R}^d} \|x_w - y\|_2 + \lambda \sqrt{\sum_j w_j^2}$$

$$\Rightarrow \boxed{\min_{w \in \mathbb{R}^d} \|x_w - y\|_2 + \lambda \|w\|_1}$$