

Exercise 2:

Algorithm:

sort y , temp-u=temp-v=temp-us=temp-vs=0

$$Y \leftarrow \sum_{i \in n} y_i$$

$$Y_s \leftarrow \sum_{i \in n} y_i^2$$

for $i = 1$ to d sort $X[:,j]$ for $j = 1$ to n

$$\text{temp-u} += y_i$$

$$\text{temp-us} += y_i^2$$

$$\text{temp-v} = Y - \text{temp-u}$$

$$\text{temp-vs} = Y_s - \text{temp-us}$$

$$u\text{-mean} = \frac{\text{temp-u}}{j}$$

$$v\text{-mean} = \frac{\text{temp-v}}{n-j}$$

$$\text{loss}(i) = \text{temp-us} - 2 \cdot u\text{-mean} \cdot \text{temp-u} + j \cdot (u\text{-mean})^2 + \\ \text{temp-vs} - 2 \cdot v\text{-mean} \cdot \text{temp-v} + (n-j) \cdot (v\text{-mean})^2$$

end loop

$$\text{min-loss} \\ \text{thresh}(j) = \arg \min_i \text{loss}(i)$$

end loop

$$\text{best-thresh} = \min_j \text{thresh}(j)$$

$$\text{best-feature} = \arg \min_j \text{thresh}(j)$$

return (best-thresh, best-feature)

Correctness:

To solve $\min_j \min_{t_j} \left[\min_{u \in R} \sum_{i: x_{ij} \leq t_j} (y_i - u)^2 + \min_{v \in R} \sum_{i: x_{ij} > t_j} (y_i - v)^2 \right] :$

if we split elements by a threshold for a specific feature, the value of other features don't matter, so we can consider each feature independently if we only need threshold for one feature.

Now for any feature selected, we only need to try using each element as threshold which has total of n .

if we specified a threshold, for any $\min_{u \in R} \sum_{i: x_{ij} \leq t_j} (y_i - u)^2$, the best value would be ~~mean(y_i)~~ $\text{mean}(y_i)$, it holds true for v as well since this will be variance of y_i and the value will increase if u is not the mean for v .

$$\begin{aligned} \therefore \min_{u \in R} \sum_{i: x_{ij} \leq t_j} (y_i - u)^2 + \min_{v \in R} \sum_{i: x_{ij} > t_j} (y_i - v)^2 &= \min_{u \in R} \sum_{i: x_{ij} \leq t_j} (y_i - \text{mean}(y_i))^2 + \min_{v \in R} \sum_{i: x_{ij} > t_j} (y_i - \text{mean}(y_i))^2 \\ &= \min_{u \in R} \sum_{i: x_{ij} \leq t_j} y_i^2 - 2y_i \text{mean}(y_i) + (\text{mean}(y_i))^2 + \\ &\quad \min_{v \in R} \sum_{i: x_{ij} > t_j} y_i^2 - 2y_i \text{mean}(y_i) + (\text{mean}(y_i))^2 \\ &= \min_{u \in R} \sum y_i^2 - 2 \cdot \text{mean}(y_i) \cdot \sum y_i + |y_i| \cdot (\text{mean}(y_i))^2 + \\ &\quad \min_{v \in R} \sum y_i^2 - 2 \text{mean}(y_i) \cdot \sum y_i + (y_i - \text{mean}(y_i))^2 \end{aligned}$$

\therefore The algorithm solves the equation.

run time: total d features, try n elements for each feature, for each try: sort takes $O(n \log n)$, calculations are linear, constant-time.

Final selection: n for each threshold, totally $d \cdot n$ then, then choose 1 feature from d

\therefore Runtime $T(n) = O(d \cdot (n \log n + n)) + O(nd) + O(n) = O(d \cdot n \log n)$