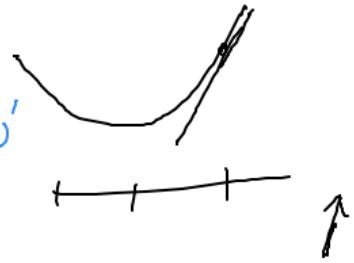


Ridge regression:

minimise for some given $\lambda \geq 0$

$$L = \sum_{i=1}^N \left(y_i - \theta_0 - \sum_{j=1}^n \theta_j x_{ij} \right)^2 + \lambda \sum_{j=1}^n \theta_j^2$$

Suppose:
 $x_{ij} > 0$
 we prohibit too small value
 i.e. $\theta_0 + \sum \theta_j x_{ij} < y_i$



to find the minimiser, i.e., $\theta_0, \dots, \theta_n$ for which L is the smallest, one can (also) use gradient descent

$$\frac{\partial L}{\partial \theta_k} = \sum_{i=1}^N 2 \left(y_i - \theta_0 - \sum_{j=1}^n \theta_j x_{ij} \right) x_{ik} + 2\lambda \theta_k \quad k=1, \dots, n$$

$$\tilde{\theta}_k = \theta_k - c \cdot \frac{\partial L}{\partial \theta_k}$$

\uparrow adjusted θ_k \uparrow learning constant

should be repeated many times

-||-
 In the neural networks world, this corresponds to a batch size = size of the training data

- One can also use regularisation for neural networks

$$L = \sum_j (t_j - y_j)^2 + \lambda \sum W^2$$

loss function

t_j target value

y_j actual value

W : weights different than these encoding bias

$$2\lambda W_{ij}^{(k)}$$

Decision trees

1) Divide the predictor space - i.e., the set of possible values (x_1, \dots, x_p) into J distinct disjoint regions R_1, \dots, R_J



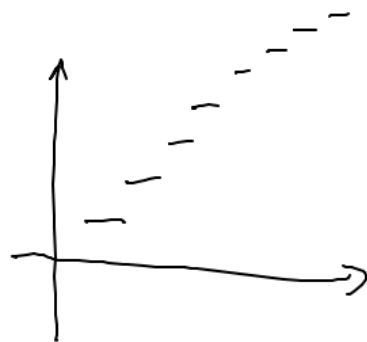
2) Prediction: for every observation that falls into R_j , we make the same

prediction: regression: mean/median of the training obs. that fell into R_j

classification: category that occurs most often among training obs. that fell into R_j
(majority vote)

Model function:

$$f(x) = \sum_{j=1}^J c_j \mathbb{1}_{R_j}(x)$$



For growing the tree:

One uses a greedy algorithm:

- start from the root \rightarrow regions = {whole space}
- when we have a tree, ^{with regions $\{R_1, \dots, R_k\}$} , for a region $R = R_k$ consider all regions of the form:

$$R^{(1)}(j, s) = R_k \cap \{(x_1, \dots, x_p) : x_j < s\}$$

$$R^{(2)}(j, s) = R_k \cap \{(x_1, \dots, x_p) : x_j \geq s\}$$

and choose one that minimizes some error

regression trees:
$$\sum_{i: x_i \in R^{(1)}(j, s)} (y_i - \hat{y}_{R^{(1)}})^2 + \sum_{i: x_i \in R^{(2)}(j, s)} (y_i - \hat{y}_{R^{(2)}})^2$$

classification:
$$\hat{p}_{mk} = \frac{\# \text{ samples in } R_m \text{ of class } k}{\# \text{ samples in } R_m}$$

$$\hat{p}_{m1} + \hat{p}_{m2} = 1$$

minimize:

$$E = 1 - \max_k \hat{p}_{mk} \quad \text{error of classification}$$

not a good one

$$G = \sum_{k=1}^K \hat{p}_{mk} (1 - \hat{p}_{mk}) = 1 - \sum_{k=1}^K (\hat{p}_{mk})^2 \quad \text{Gini}$$

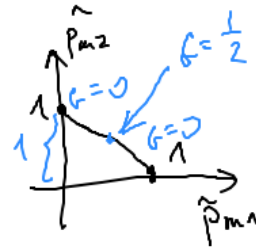
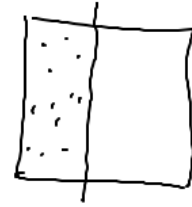
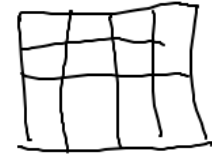
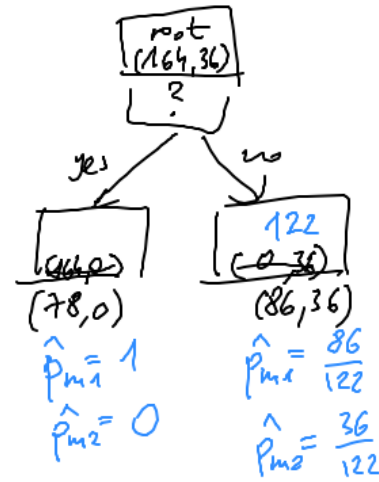
-Gini

$$D = - \sum_{k=1}^K \hat{p}_{mk} \log \hat{p}_{mk}$$

-Entropy

We stop when all leaves have a number of samples $<$ threshold.

weighted area should be minimized for all leaves \rightarrow for one node



Second step: pruning

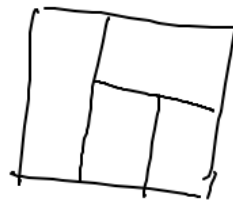
cost of complexity:

$$L \geq 0$$

Minimise

$$\sum_{m=1}^M \sum_{i: x_i \in R_m} (y_i - \hat{y}_{R_m})^2 + \underbrace{2|T|}_{\text{the number of leaves}}$$

prediction error (on training data) \leftarrow regression



- for each node, we check for which L the cost of complexity would decrease by pruning the subtree of that node
- we prune the subtree of the node with the smallest $L \rightarrow$ obtain a new tree
- repeat for the new tree



number of leaves will be smaller by 2

the error of prediction will increase by ϵ

the cost will change by

$$\epsilon - 2L \leq 0 \quad L \geq \frac{\epsilon}{2}$$

sequence of trees!

