

(12-1)

Useful "non-parametric" learning methods

Non-parametric really means that we don't pick a fixed-size parameterization of the hypothesis class in advance. We'll look at two really useful, easy-to-implement, fast-to-run methods for supervised learning.

Decision trees

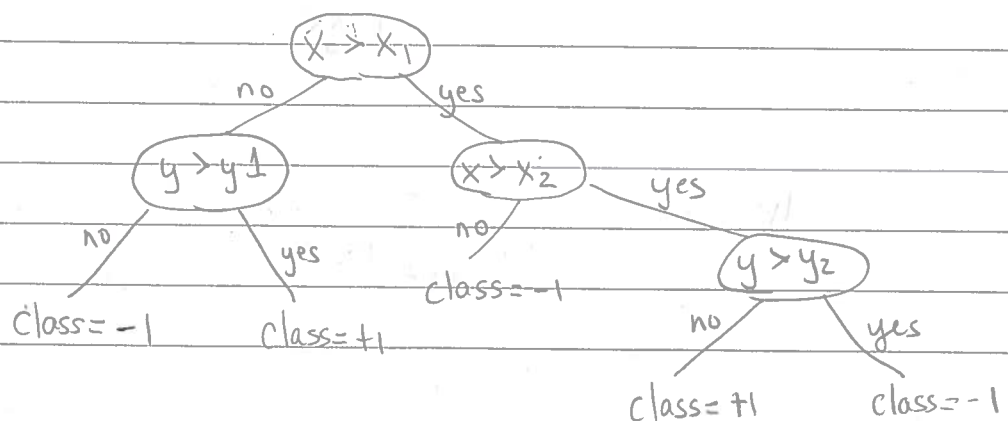
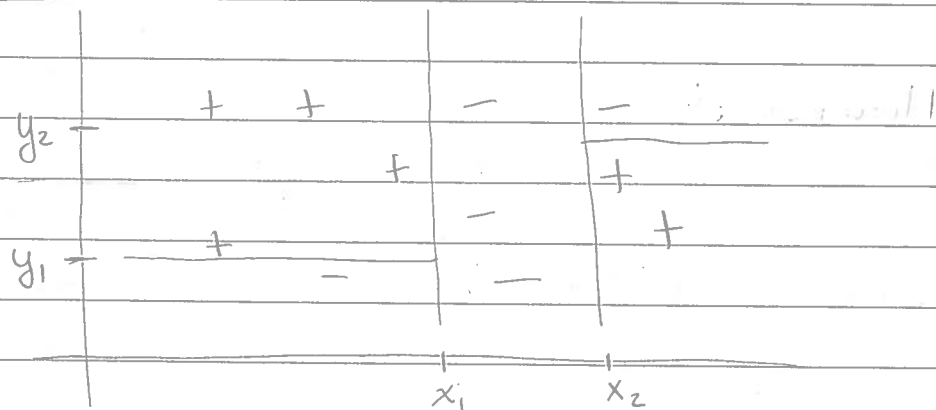
Idea: recursively partition the input space using axis-aligned splits, and fit a simple (usually constant) model in each partition.

Pros: easy to train, understandable by humans

Cons: lower accuracy, not good for very large input spaces (e.g. images)

CART/ID3: invented in parallel in AI & stats communities in the 80's.

Classification



(12-2)

We'd like to find a small tree (our regularizer) that predicts training data well. Can always be perfect! but may not want to be.

Tree can be described by:

- set of regions R_1, \dots, R_M
- output values O_1, \dots, O_M

Hypothesis: $h(x) = O_m$ if $x \in R_m$

Error in region: $E_m = |\{i \mid x^{(i)} \in R_m \text{ and } y^{(i)} \neq O_m\}|$

JP hard!
We'll do
greedy
approx.

Objective: find R, O to minimize
 $cM + \sum_{m=1}^M E_m$

Define empirical probability of an item from region R_m being in class k as

Multiclass
is easy

$$\hat{p}_{mk} = \frac{|\{i \mid x^{(i)} \in R_m \text{ and } y^{(i)} = k\}|}{N_m}$$

$$N_m = |\{i \mid x^{(i)} \in R_m\}|$$

BuildTree(D, K): D = data, K = min leaf size

if $|D| < K$:

return Leaf($\arg \max_k |\{i \mid x^{(i)} \in D, y^{(i)} = k\}|$)

else:

find dimension v and split-point s that

yields D_L and D_R :

$$D_L = \{(x^{(i)}, y^{(i)}) \in D \mid x_v^{(i)} < s\}$$

$$D_R = \{(x^{(i)}, y^{(i)}) \in D \mid x_v^{(i)} \geq s\}$$

that minimizes

$$|D_L| \cdot \text{impurity}(D_L) + |D_R| \cdot \text{impurity}(D_R)$$

return Node(v, s , BuildTree(D_L, K), BuildTree(D_R, K))

(12-3)

Measures of impurity

• Entropy $Q(D) = H(\{y^{(i)} \mid (x^{(i)}, y^{(i)}) \in D\})$

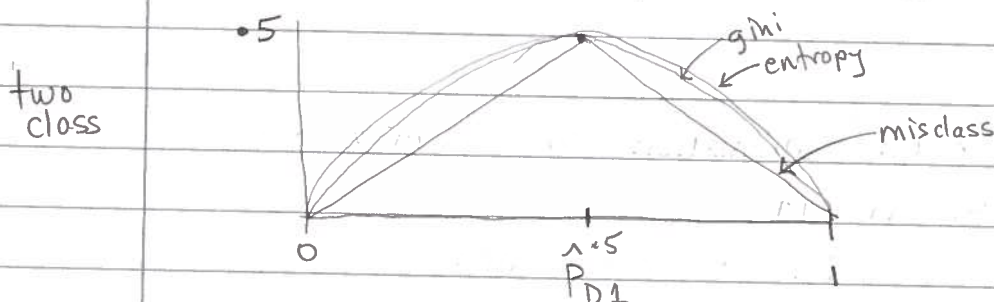
$$= - \sum_K \hat{P}_{DK} \log \hat{P}_{DK}$$

• Gini index

$$\sum_K \hat{P}_{DK} (1 - \hat{P}_{DK})$$

• Misclassification error $E_m/N_m = 1 - \hat{P}_{D \text{ maj } K}$

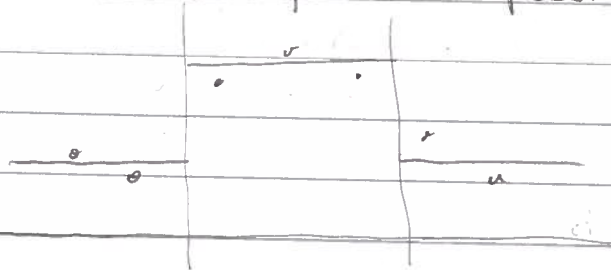
People usually use entropy. Not a big difference.



Generally

- For best results, grow a very big tree (small k) and then prune it.
- Easy to incorporate different loss functions (greedily minimize expected loss at each stage)
- Discrete features easy
- Each iteration takes time $O(d \cdot n \cdot n)$
 - consider split
 - compute $n = \text{number of data points} = \# \text{ possible splits}$

Regression



$h =$ Piecewise constant function.

(12-5)

Random forest last step

Combine votes of all trees to get resulting hypothesis.

Works surprisingly well. Always try as baseline!

Nearest neighbor

Lazy! Don't do any work at training time — just remember your data

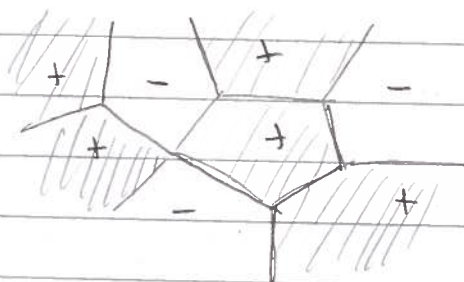
Need distance metric $d(x, x')$

- $d(x, x) = 0$
- $d(x, x') \geq 0$
- $d(x, x') = d(x', x)$
- $d(x, x'') \leq d(x, x') + d(x', x'')$

Hypothesis:

$$h(x) = y^{(i)} \text{ where } i = \arg \min_i d(x, x^{(i)})$$

that is, the output of the closest training point.



Voronoi partition

Variations

- Find K nearest neighbors; output average (regression) or majority (classification)
- Fit locally linear regression models

Use good data structures (e.g. ball trees) to find neighbors efficiently.

Works poorly in high dimensions, often.

(12-4)

Prediction at leaf

$$O_m = \text{average } y^{(i)}_{\{i | x^{(i)} \in R_m\}}$$

Error at leaf

$$E_m = \sum_{\{i | x^{(i)} \in R_m\}} (y^{(i)} - O_m)^2$$

Pick j, s split that minimizes

$$E_{D_{\text{left}}} + E_{D_{\text{right}}}$$

Random forests

Why do trees often have high prediction error? Because they have high estimation error (variance).

→ easy to overfit

→ highly sensitive to data points

Idea: reduce variance (of any learning algorithm) by bagging (bootstrap aggregation):

- randomly sample B subsets of size n , with replacement, from D
- train a new hypoth. on each one $\Rightarrow h_i$
- return

$$h(x) = \frac{1}{B} \sum_i h_i(x) \quad \text{for regression}$$

$$h(x) = \text{majority}_i (h_i(x)) \quad \text{for classification}$$

Random forest

for $b = 1 \dots B$:

D_b = bootstrap sample of size n from D

T_b = decision tree but at each iteration

- select m dimensions at random
- pick best split among them