

# Machine Learning Summary

## Supervised

Types of models:

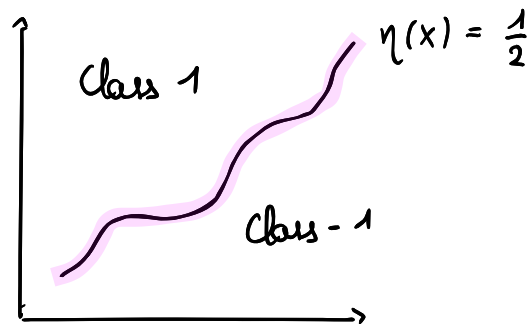
- **discriminant**: estimate decision frontier directly  
 $P(Y=1 | X=x)$
- **generative**: estimate  $P(X=x_c | Y=1)$  and  $P(X=x_c | Y=-1)$

If we know the decision frontier  $\eta(x) = P(Y=1 | X=x)$ , there is no need for simulation. This is the **Bayes classifier**:

$$g(x) = \begin{cases} +1 & \text{if } \eta(x) \geq \frac{1}{2} \\ -1 & \text{if } \eta(x) < \frac{1}{2} \end{cases}$$

$$\Leftrightarrow g(x) = 2 \mathbb{1} \left\{ \eta(x) \geq \frac{1}{2} \right\} - 1$$

However, this classifier can never be reached.



Usually, we try to approach the Bayes classifier using the so-called **plug-in estimator**:

↳ Estimate  $\eta(x)$

• Plug-it in  $\hat{g}(x) = 2 \mathbb{1} \left\{ \hat{\eta}(x) \geq \frac{1}{2} \right\} - 1$

# I. Generative models.

## 1. Linear discriminant analysis (LDA)

Let:  $\cdot G = \mathcal{L}(X | Y=1)$   
 $H = \mathcal{L}(X | Y=-1)$  } distribution of the 2 classes

Define a likelihood ratio  $\Phi(X) = \frac{dG}{dH}(x) = \frac{P(X=x | Y=+1)}{P(X=x | Y=-1)}$

$$= \frac{P(Y=1 | X=x) \cdot \cancel{P(X=x)} / P(Y=1)}{P(Y=-1 | X=x) \cdot \cancel{P(X=x)} / P(Y=-1)} = \frac{1-p}{p} \cdot \frac{\eta(x)}{1-\eta(x)}$$

$$\hookrightarrow \frac{p \Phi(x)}{(1-p) + p \Phi(x)} = \eta(x)$$

Gaussian underlying hypothesis:

$$\left. \begin{array}{l} G = N_+(\mu_+, \Gamma) \\ H = N_-(\mu_-, \Gamma) \end{array} \right\} N(\mu, \sigma^2) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{1}{2\sigma^2}(x-\mu)^2}$$

$$\text{If } \eta(x) \geq \frac{1}{2} \Rightarrow \Phi(x) \geq \frac{1-p}{p}$$

$$\hookrightarrow {}^t x \Gamma^{-1} (\mu_+ - \mu_-) + \frac{1}{2} ({}^t \mu_- \Gamma^{-1} \mu_- - {}^t \mu_+ \Gamma^{-1} \mu_+) \geq \log \frac{p}{1-p}$$

$$\Leftrightarrow \alpha + {}^t \beta x \geq 0$$

$$\text{where: } \cdot \beta = \Gamma^{-1} (\mu_+ - \mu_-)$$

$$\cdot \alpha = \frac{1}{2} ({}^t \mu_- \Gamma^{-1} \mu_- - {}^t \mu_+ \Gamma^{-1} \mu_+) - \log \left( \frac{p}{1-p} \right)$$

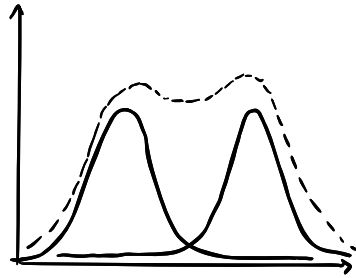
$$\hookrightarrow \hat{p} = \frac{n_+}{n}, \quad \hat{\mu}_+ = \frac{1}{n_+} \sum_{\{x_i=+1\}}^n x_i, \quad \hat{\mu}_- = \frac{1}{n_-} \sum_{\{x_i=-1\}}^n x_i$$

## 2. Quadratic discriminant analysis (QDA).

Relax the Gaussian hypothesis  $\Rightarrow$  Mixture of Gaussians.

$$\delta F_1 + (1-\delta) F_2$$

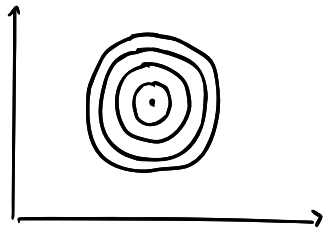
Good! It can be bi-modal.



## 3. Naive Bayes classifier

We now suppose independence among our variables:

$$\Gamma = \prod_{1 \leq i \leq d} \sigma_i^2$$



## II. Discriminant models

### 1. Perceptron

$$g(x) = \text{sgn}(d + {}^t \beta x)$$

$$\min_{\alpha, \beta} \hat{L}_n(g), \quad \hat{L}_n(g) = \frac{1}{n} \sum_{i=1}^n \mathbb{1}\{-(d + {}^t \beta x_i) y_i > 0\}$$

$$\begin{array}{c} d \\ x^0 \\ \vdots \\ x^d \end{array} \begin{array}{c} \diagdown \\ \beta_0 \\ \diagup \end{array} \begin{array}{c} \beta_1 \\ \vdots \\ \beta_d \end{array} \begin{array}{c} \diagdown \\ \diagup \end{array} L(x, \beta_1, \dots, \beta_d) \begin{array}{c} \diagup \\ \diagdown \end{array} \begin{array}{c} +1 \\ -1 \end{array}$$

Due to lack of derivability of "sign" function:

↳ Apply stochastic gradient descent → Steps evolving at some point.

Limitation: • Data need to be linearly separable

• Does not care about margin of the hyperplane.

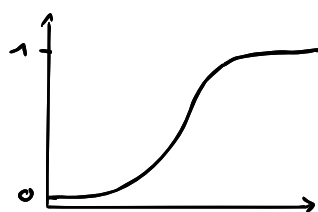
## 2. Logistic regression

• Regression:  $\eta(x) = \mathbb{E}(Y|x) = \alpha + \beta x$

• Classification:  $\eta(x) = \frac{e^{\alpha + \beta x}}{1 + e^{\alpha + \beta x}} = \frac{1}{1 + e^{-(\alpha + \beta x)}}$

Find  $\hat{\alpha}$  and  $\hat{\beta}$  by conditional maximum log-likelihood, and plug-in:

$$(\hat{\alpha}, \hat{\beta}) \Rightarrow \hat{\eta}_{\hat{\alpha}, \hat{\beta}}(x) \Rightarrow \hat{g}(x) = 2 \mathbb{1}\{\hat{\eta}_{\hat{\alpha}, \hat{\beta}}(x) \geq \frac{1}{2}\} - 1$$



$$= 2 \mathbb{1}\left\{\frac{e^{\hat{\alpha} + \hat{\beta}x}}{1 + e^{\hat{\alpha} + \hat{\beta}x}} \geq \frac{1}{2}\right\} - 1$$

$$= 2 \mathbb{1}\{\hat{\alpha} + \hat{\beta}x \geq 0\} - 1$$

## 3. K-Nearest Neighbors

Let  $D$  be a distance metric (e.g. Euclidean) and  $\sigma_x$  classifies points.

$$D(x, x_{\sigma_x(1)}) \leq D(x, x_{\sigma_x(2)}) \leq \dots$$

↳ Only keep  $K$  first neighbors.

$$g_{KNN}(x) = g_{KNN}(x, y_{\sigma_x(1)}, \dots, y_{\sigma_x(k)})$$

$$= \begin{cases} +1 & \text{if } \sum_{i=1}^k \mathbb{1}\{y_{\sigma_x(i)}\} > \frac{k}{2} \\ -1 & \text{if } \sum_{i=1}^k \mathbb{1}\{y_{\sigma_x(i)}\} < \frac{k}{2} \end{cases}$$

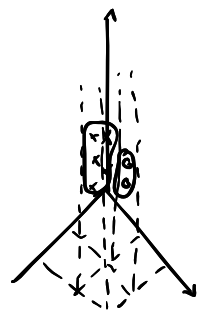
$$= 2 \cdot \mathbb{1}\left\{\sum_{i=1}^k y_{\sigma_x(i)} \geq 0\right\} - 1$$

- Limitations :
- Sensitive to metric  $D$
  - Sensitive to choice of  $K$
  - Computationally expensive
  - Does not scale

- Extreme cases :
- $k = 1 \rightarrow$  Training error is 0
  - $k = n \rightarrow$  Training error is  $\min \left\{ \frac{n_+}{n}, \frac{n_-}{n} \right\}$

#### 4. Local averaging

- Divide feature space in regions  $C_1 \cup \dots \cup C_k$
- $\hat{\eta}(x) = \sum_{k=1}^K \mathbb{1}\{x \in C_k\} \frac{\sum_{i=1}^n \mathbb{1}\{Y_i = 1, X_i \in C_k\}}{\sum_{i=1}^n \mathbb{1}\{X_i \in C_k\}}$  Majority rule.  
 $\hookrightarrow$  Nadaraya-Watson
- $\hat{g}(x) = 2 \mathbb{1}\{\hat{\eta}(x) \geq \frac{1}{2}\} - 1$



Limitations : Depends highly on the regions we consider

#### 5. Decision tree

Build  $g(x) = \sum_{l=1}^m \mathbb{1}\{x \in C_l\} (\arg\max_c [\sum_{i: x \in C_l} \mathbb{1}\{Y_i = c\}])$

where  $C_1 \dots C_m$  automatic partitions

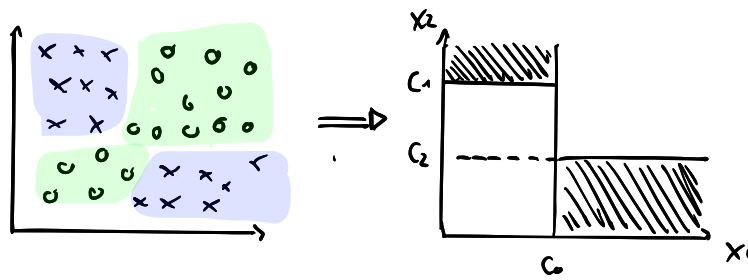
- Partition rule should minimize the local criteria

1) Crossed entropy :  $H(s) = - \sum_{l=1}^c p_l(s) \log p_l(s)$

2) Gini index :  $H(s) = \sum_{l=1}^c p_l(s) (1 - p_l(s))$

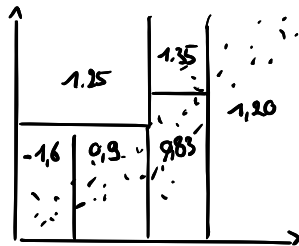
3) Classification error :  $H(s) = 1 - p_c(s)$

where  $c$  # classes and  $p$  proportion of well classified data for each point :  $p_c(s) = \frac{1}{n} \sum_{i=1}^n \mathbb{1}\{Y_i = c\}$



- Advantage:**
- Interpretable, consistent, multiclass
  - No pre-processing
- Limitations:**
- Large variance, unstable
  - No global optimization.

Can also be applied for regression trees:

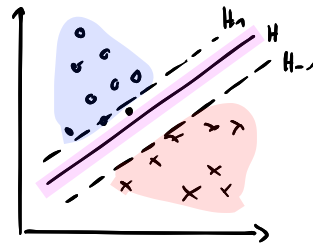


$$L(t_j, \mathcal{Z}, S) = \frac{n_d}{n} \text{Var}_{\text{emp}}(D(j, \mathcal{Z}, S)) + \frac{n_g}{n} \text{Var}_{\text{emp}}(G(j, \mathcal{Z}, S))$$

$$\text{where } \text{Var}_{\text{emp}}(S) = \frac{1}{|S|} \sum_{x_i \in S} (y_i - \bar{y})^2$$

## 6. Support Vector Machine

### a. Linear SVM



• **Separable case**

$$\min_{w, b} \frac{1}{2} \|w\|^2 \text{ s.t. } 1 - y_i (w^T x_i + b) \leq 0, \quad i = 1 \dots n$$

$$\min_{w, b, \alpha} \Rightarrow \mathcal{L}(w, b, \alpha) = \frac{1}{2} \|w\|^2 + \sum_i \alpha_i (1 - y_i (w^T x_i + b)), \quad \alpha_i \geq 0$$

Find  $w$ ,  $b$  and  $\alpha$  by quadratic solver

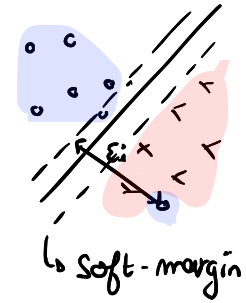
$$\hookrightarrow g(x) = \text{sign} \left( \sum_{i=1}^n \alpha_i y_i x_i^T x + b \right)$$

## • Non-separable case

Define slack variables as misclassified:  $\varepsilon_i$

$$\min_{w, b, \varepsilon} \frac{1}{2} \|w\|^2 + C \sum_{i=1}^n \varepsilon_i$$

$$\text{s.t. } y_i (w^T x_i + b) \geq 1 - \varepsilon_i, \quad i = 1 \dots n$$



$$\mathcal{L}(w, b, \varepsilon, d, \mu) = \frac{1}{2} w^T w + C \sum \varepsilon_i + \sum_{i=1}^n d_i (1 - \varepsilon_i - y_i (w^T x_i + b)) - \sum \mu_i \varepsilon_i$$

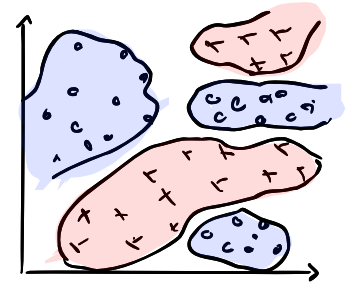
$C$  regulates the importance of misclassified observations.

- If:
- $d_i^* = 0 \rightarrow$  No error, does not influence hyperplane
  - $0 < d_i^* < C \rightarrow$  On the frontier
  - $d_i^* = C \rightarrow x_i^*$  is misclassified and is support

## b. Non-linear SVM

Using kernel trick:

$$g(x) = \text{sign} \left( \sum_i d_i y_i \underbrace{k(x_i, x)}_{\text{instead of } x_i^T x} + b \right)$$



Idea: Transform data through non-linear function  $\phi$

Thanks to Moore-Aronszajn, guarantee that:

- $\mathcal{F}$  Hilbert space: feature space
- $\phi: \mathcal{X} \rightarrow \mathcal{F}$ : feature map

$$\text{s.t. } \langle \phi(x), \phi(x') \rangle_{\mathcal{F}} = k(x, x')$$

$$\min_{w, b, \varepsilon} \frac{1}{2} \|w\|^2 + C \sum_{i=1}^n \varepsilon_i \quad \text{s.t. } y_i (w^T \phi(x_i) + b) \geq 1 - \varepsilon_i \quad i = 1 \dots n$$

$$\varepsilon_i \geq 0$$

$$\hookrightarrow g(x) = \text{sign} \left( \sum_i d_i y_i \phi(x_i)^T \phi(x) + b \right)$$

$$= \text{sign} \left( \sum_i d_i y_i k(x_i, x) + b \right)$$

What kernels can we use?

- Linear:  $k(x, x') = x^T x'$
- polynomial:  $k(x, x') = (x^T x' + c)^d$
- gaussian:  $k(x, x') = \exp(-\gamma \|x - x'\|^2)$  (RBF kernels)

Why do we use kernels?

- work in spaces of infinite dimensions
- gaussian kernels allow approximation of most problems
- There is an algebra for kernels

Limitations:

- Find parameter  $C$  and  $\gamma$  for kernel
- Multiclass does not really work  $\rightarrow$  several binary SVMs instead.

### C. Support Vector Regressor

Define an  $\epsilon$ -tube.

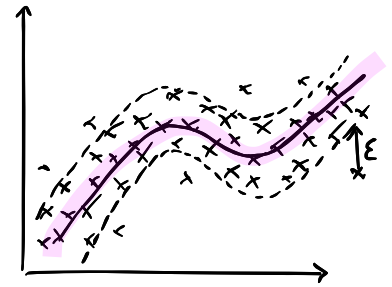
$$\min_{w, b, \epsilon} \frac{1}{2} \|w\|^2 + C \sum_i (\epsilon_i + \epsilon_i^*)$$

$$\text{s.t. } y_i - f(x_i) \leq \epsilon + \epsilon_i$$

$$f(x_i) - y_i \leq \epsilon + \epsilon_i^*$$

$$\epsilon_i^* \geq 0$$

where in the general case:  $f(x) = w^T \phi(x) + b$



Limitations:

- Quadratic program computationally expensive
- Does not really scale



## 7. Bagging

Comes from Bootstrap Aggregating.

- Draw  $T$  training independent samples  $\{S_1 \dots S_T\}$
- Learn a model  $f_t \in \mathcal{F}$  from each sample  $S_t$
- Compute the average model  $f_{\text{ens}}(x) = \frac{1}{T} \sum_{t=1}^T f_t(x)$

Why use bagging?

- Bias remains the same:  $\mathbb{E}_{S_1 \dots S_T} [f_{\text{ens}}(x)] = \mathbb{E}_S [f_S(x)]$
- Variance divided by  $T$ : 
$$\mathbb{E}_{S_1 \dots S_T} [(f_{\text{ens}}(x) - \mathbb{E}_{S_1 \dots S_T} [f_{\text{ens}}(x)])^2] = \frac{1}{T} \mathbb{E}_S [(f_S(x) - \mathbb{E}_S [f_S(x)])^2]$$
- Useful if unstable learning algorithm such as trees.

Limitations:

- Obtained model is complex
- Outperformed by random forest and gradient boosting.

## 8. Random Forest

We have a train set  $S_{\text{train}}$  and  $p$  features.

- Create  $T$  samples (bootstrap)  $\{S_1 \dots S_T\}$
- For  $t = 1 \dots T$ :
  - Select  $f$  features s.t.  $f \ll p$  randomly without replacement.
  - $h_{\text{tree}}^{(t)}$  decision tree learned from  $S_t$  with best split of features.
- $H^T = \frac{1}{T} \sum_t h_{\text{tree}}^{(t)}$

### 3. Extra Trees

We won't use bootstrap here, but  $\text{Strain}$  each time.

• for  $t = 1 \dots T$ :

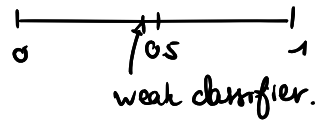
- Always use  $\text{Strain}$
- Select  $f$  features randomly without replacement  $\rightarrow \text{Split}(S, i)$ 
  - Let  $a'_{\max}$  and  $a'_{\min}$  denote maximal value of  $x_i$  in  $S$
  - Draw uniformly a cut-point  $a_c$  in  $[a'_{\max}, a'_{\min}]$
  - ↳ Choose the best split
- $h_{\text{tree}}^{(t)}$  randomized tree learned from  $\text{Strain}$ .
- $H^T = \frac{1}{T} h_{\text{tree}}^{(t)}$ .

Advantage: Fast, parallelizable, easy to tune

Limitations: • Overfitting if large tree  
• Lost interpretability

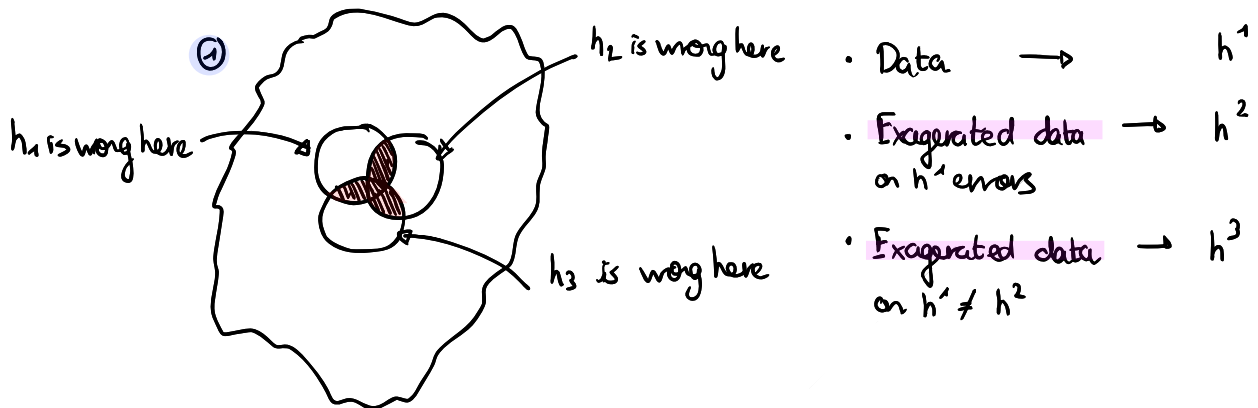
## 10. Boosting

Let  $h$  be a binary weak classifier  $\Rightarrow [-1, +1]$

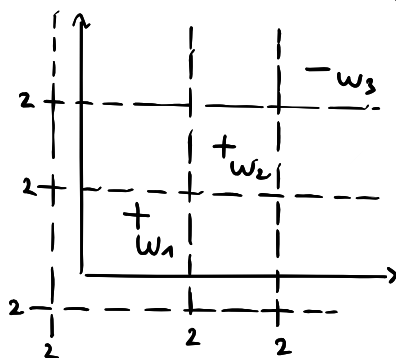


### • Adaboost.

Boosting is a way to recursively combine classifiers that are weak learners. It is based on several key concepts:



### ② Decision tree stumps



A stump defines what we can do with 1 test. Each possible test is a classifier

↳ 12 possible tests here

The error rate is:

$$\epsilon = \sum_{y_i \neq h(x_i)} w_i \quad \text{s.t.} \quad \sum w_i = 1$$

The weights account for the exaggeration on some data.

$$H(x) = \text{sign} ( \alpha^1 h^1(x) + \alpha^2 h^2(x) + \alpha^3 h^3(x) + \dots )$$

General idea of Adaboost:

1) Let  $w_i^1 = \frac{1}{N}$ ,  $N = \# \text{ samples}$

2) Pick  $h^t$  that minimizes  $\epsilon^t \rightarrow$  pick  $d^t$

3) Calculate  $w_k^{t+1}$  s.t:

$$w_i^{t+1} = \frac{w_i^t}{Z} e^{-\alpha^t h^t(x) y(x)}$$

where:

- $Z$  normalizing constant  $\Rightarrow \sum w_i = 1$
- $h^t(x)$  predicted classifier value
- $y(x)$  label of training observation.

↳ If  $y_i$  and  $h^t(x_i)$  have same sign, exponential is negative and weight on observation is small.

We want to choose  $d^k$  such that it minimizes the error:

$$\alpha_t = \frac{1}{2} \log \left( \frac{1 - \epsilon_t}{\epsilon_t} \right) \text{ is an error bound.}$$

$$\Rightarrow w_i^{t+1} = \frac{w_i^t}{2} \times \begin{cases} \sqrt{\frac{\epsilon_t}{1-\epsilon_t}} & \text{if well classified} \\ \sqrt{\frac{1-\epsilon_t}{\epsilon_t}} & \text{if misclassified} \end{cases}$$

What is the value of the constant  $Z$ ?

$$\sqrt{\frac{\epsilon^t}{1-\epsilon^t}} \cdot \sum_{y_i = \text{hrr}(x_i)} w_{ji}^t + \sqrt{\frac{1-\epsilon^t}{\epsilon^t}} \sum_{y_i \neq \text{hrr}(x_i)} w_{ji}^t = z = 2\sqrt{\epsilon^t(1-\epsilon^t)}$$

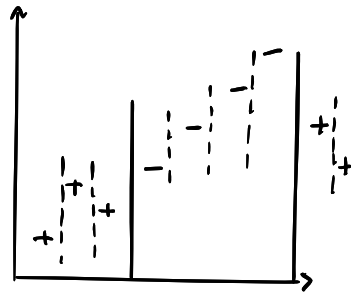
We can replace:  $w_i^{t+1} = \begin{cases} \frac{w^t}{2} \cdot \frac{1}{1-\epsilon} & \text{if correct} \\ \frac{w^t}{2} \cdot \frac{1}{\epsilon} & \text{if wrong} \end{cases}$

$\Rightarrow \sum_{\text{correct}} w^{t+1} = \frac{1}{2}$  and  $\sum_{\text{wrong}} w^{t+1} = \frac{1}{2}$  whereas number of misclassified observations decreases (still larger weights).

Used for computations:

- no  $Z$
- no  $d$
- no exponential.

To reduce computation time, also notice that a classifier between 2 well classified points will never be optimal  $\Rightarrow$  No need to compute it.



Limitations: • Adaboost does overfit at some point  
↳ early stop or L1-reg.

## • Gradient boosting

Idea: At each boosting step, we would like to find how to correct  $h^{t-1}$  by gradient descent.

$$-g_t(x_i) = \frac{\partial \ell(y_i, h_{t-1}(x_i))}{\partial h_{t-1}(x_i)}$$

But gradient only defined for  $x_i$ 's values, not all  $x$ .  $\Rightarrow$  Take best approximation.

$$(h_t, d_t) = \arg\min_{h, d} \sum_{i=1}^n (\ell(y_i, h_{t-1}(x_i) + d h)) = L(y_i, h_{t-1}(x_i) + d h(x_i))$$

$$\approx \ell(y_i, h_{t-1}(x_i)) + d \langle \nabla \ell(h_{t-1}), h \rangle$$

- Steps:**
- Replace minimization step by a gradient descent type step
  - Choose  $h_t$  as the best possible descent direction.
  - Choose  $\lambda_t$  that minimizes  $L(y, H + \lambda h_t)$

$\Rightarrow$  Gradient boosting and Adaboost can be shown to be equivalent.

### Anyboost / Forward Stagewise Additive model

1) Set  $t=0$ ,  $H_0 = 0$

2) For  $t=1, \dots, T$ :

$$\bullet h_t, \lambda_t = \underset{h, \lambda}{\operatorname{argmin}} \sum_{i=1}^n l(y_i, H_{t-1}(x_i) + \lambda h(x_i))$$

$$\bullet H_t = H_{t-1} + \lambda_t h_t$$

3) Output  $H_T = \sum_{t=1}^T \lambda_t h_t$ .

Losses might be:

- **Adaboost**:  $l(y, h) = e^{-yh}$
- **Logit boost**:  $l(y, h) = \log(1 + e^{-yh})$
- **L2-boost**:  $l(y, h) = (y-h)^2$
- **L1-boost**:  $l(y, h) = |y-h|$
- **Huberboost**:  $l(y, h) = |y-h|^2 \mathbb{1}_{|y-h| \leq \epsilon} + (2\epsilon|y-h| - \epsilon^2) \mathbb{1}_{|y-h| > \epsilon}$

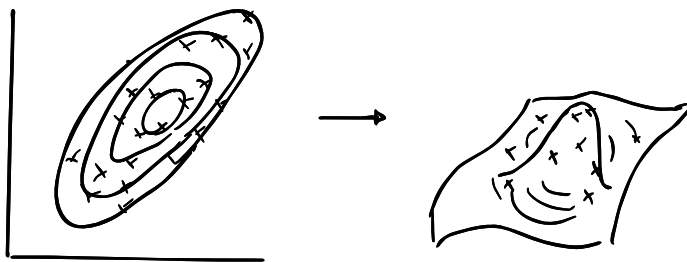
## VII. Depth and depth-based classification

### 1. Depth metrics.

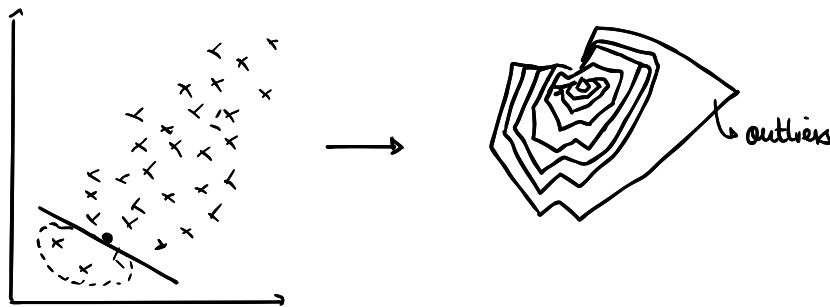
Pick a centroid  $X$  and compute depth of point  $x$  to  $X$ . Adds a dimension to plot the data in 3D.

• Mahalanobis depth  $D^{\text{Mah}(n)}(x|X) = \frac{1}{1 + (x - \mu_X)^T \Sigma_X^{-1} (x - \mu_X)}$

↳ Contours are however always ellipses, and not robust to outliers

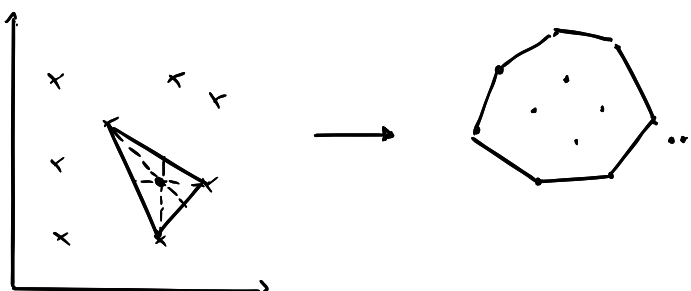


• Tukey depth:  $D^{\text{Tukey}}(x|X) = \frac{1}{n} \min_{u \in S^{d-1}} \# \{i: u^T X_i \geq u^T x\}$  number of closed halfspace linked to a point  $x$

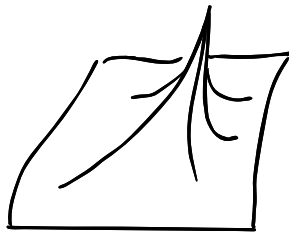


Robust to outliers, and non-parametric.

• Zonoid depth: Number of centers that can be built from e.g. 3 observations



• Projection depth:  $D^{PD}(x | X) = \frac{1}{1 + \sup_{u \in S^{d-1}} \frac{|x^T u - \text{med}(X^T u)|}{\text{MAD}(X^T u)}}$

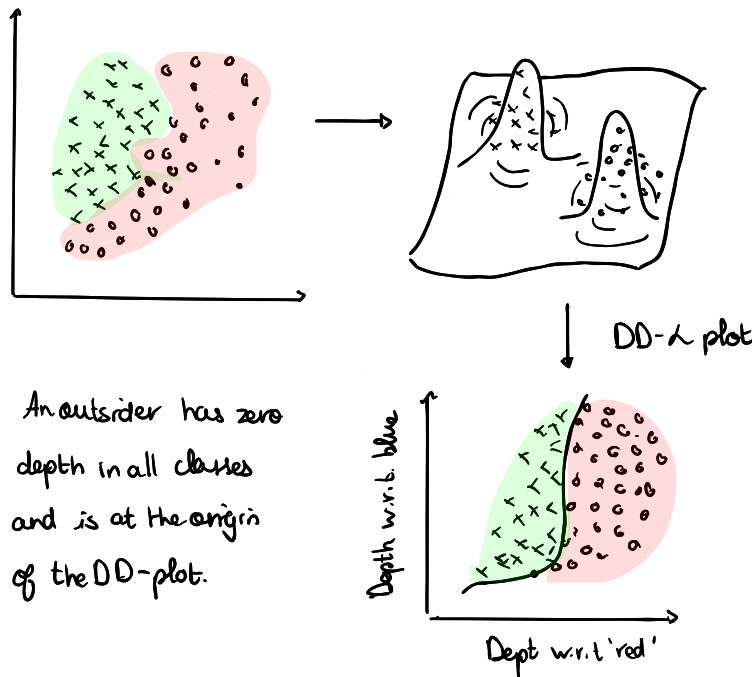


## 2. Depth-based classification

Suppose  $Q$  conditional r.v of  $Q$  classes,  $X_1 \dots X_Q$  corresponding training samples containing  $n_1 \dots n_Q$  observations. Maximal depth classifier.

$$cl_n^{(n)}(x) = \arg \max_{q \in \{1, \dots, Q\}} D^{(n)}(x | X_q)$$

↳ Assign to cluster (with centroid) giving largest depth.



An outsider has zero depth in all classes and is at the origin of the DD-plot.

## 3. Depth KNN

Define depth-based neighborhood  $R_x^{\beta(n)}$  as the smallest depth region  $D_\alpha^{(n)}(Y)$  containing  $k$  observations from  $X_0 \cup X_1$  where  $\beta(n) = \frac{k}{n_1 + n_2}$

$$cl_n^{\text{DKNN}}(x) = \begin{cases} 0 & \text{if } k_{x_0} \geq k_{x_1} \\ 1 & \text{if } k_{x_0} < k_{x_1} \end{cases} \text{ where } k_{x_i} = \#\{z: z \in X_i \cap R_x^{\beta(n)}\}$$