

# Stefan Bielmeier Homework 3 - Nov 2nd 2021

## Problem 1)

(Goal: write  $P(Y=1|X) = \frac{1}{1 + \exp(-b + \omega^T X)}$  in terms of

- $\pi = P(Y=1)$  (Bernoulli variable outcome - binary)
- $\mu_{j,k}$  = mean of Gaussian for  $j \in \{0, 1\}$
- $X_j \sim \mathcal{N}(\mu_{j,k}, \sigma_j)$  for  $j = 1, \dots, d$  features
- Cond. independence of  $X_1, X_2, X_3$  given  $Y=k$

$$\Rightarrow \text{Bayes: } P(Y=1|X) = \frac{P(X|Y=1) \cdot P(Y=1)}{P(X)} \quad \text{where } P(X) = P(X|Y=1) \cdot P(Y=1) + P(X|Y=0) \cdot P(Y=0)$$

$$= \frac{P(X|Y=1) \cdot P(Y=1)}{P(X|Y=1) \cdot P(Y=1) + P(X|Y=0) \cdot P(Y=0)}$$

$$= \frac{1}{1 + \frac{P(X|Y=0) \cdot P(Y=0)}{P(X|Y=1) \cdot P(Y=1)}}$$

with conditional independence of features:

$$P(X|Y=0) \cdot P(Y=0) = P(Y=0) \cdot \prod_{j=1}^d P(X_j|Y=0)$$

$$\text{and } P(X|Y=1) \cdot P(Y=1) = P(Y=1) \cdot \prod_{j=1}^d P(X_j|Y=1) \quad (\text{chain rule})$$

$$= \frac{1}{1 + \frac{P(Y=0)}{P(Y=1)} \cdot \frac{\prod_{j=1}^d P(X_j|Y=0)}{\prod_{j=1}^d P(X_j|Y=1)}}$$

with  $P(X_j|Y=k) =$

$$= \frac{1}{1 + \left(\frac{1-\pi}{\pi}\right) \cdot \frac{\prod_{j=1}^d P(X_j|Y=0)}{\prod_{j=1}^d P(X_j|Y=1)}}$$

$$= \frac{1}{\sqrt{2\pi\sigma_j^2}} \cdot \exp\left(\frac{-(X_j - \mu_{j,k})^2}{2\sigma_j^2}\right)$$

Substituting, we get:

$$\begin{aligned}
 \Rightarrow \frac{\prod_{j=1}^d P(X_j | Y=0)}{\prod_{j=1}^d P(X_j | Y=1)} &= \frac{\prod_{j=1}^d \frac{1}{\sqrt{2\pi\sigma_j^2}} \cdot \exp\left(-\frac{(X_j - \mu_{j0})^2}{2\sigma_j^2}\right)}{\prod_{j=1}^d \frac{1}{\sqrt{2\pi\sigma_j^2}} \cdot \exp\left(-\frac{(X_j - \mu_{j1})^2}{2\sigma_j^2}\right)} \\
 &= \prod_{j=1}^d \exp\left(\frac{-(X_j - \mu_{j0})^2}{2\sigma_j^2} - \frac{(X_j - \mu_{j1})^2}{2\sigma_j^2}\right) \\
 &= \prod_{j=1}^d \exp\left(\frac{-X_j^2 + 2\mu_{j0}X_j - \mu_{j0}^2 + X_j^2 - 2\mu_{j1}X_j + \mu_{j1}^2}{2\sigma_j^2}\right) \\
 &= \prod_{j=1}^d \exp\left(\frac{\mu_{j0} - \mu_{j1}}{\sigma_j^2} \cdot X_j + \frac{\mu_{j1}^2 - \mu_{j0}^2}{2\sigma_j^2}\right)
 \end{aligned}$$

$\Rightarrow$  back into other equation:

$$\begin{aligned}
 P(Y=1|X) &= \frac{1}{1 + \frac{1-\pi}{\pi} \cdot \prod_{j=1}^d \exp\left(\frac{\mu_{j0} - \mu_{j1}}{\sigma_j^2} \cdot X_j + \frac{\mu_{j1}^2 - \mu_{j0}^2}{2\sigma_j^2}\right)} \quad \left| \begin{array}{l} \text{we want} \\ \frac{1}{1 + \exp(\dots)} \end{array} \right. \\
 &= \frac{1}{1 + \exp\left(\ln\left(\frac{1-\pi}{\pi}\right) + \prod_{j=1}^d \exp(\dots)\right)} = \frac{1}{1 + \exp\left(\ln\left(\frac{1-\pi}{\pi}\right) + \ln\prod_{j=1}^d \exp(\dots)\right)} \\
 &= \frac{1}{1 + \exp\left(\ln\left(\frac{1-\pi}{\pi}\right) + \sum_{j=1}^d \frac{\mu_{j0} - \mu_{j1}}{\sigma_j^2} \cdot X_j + \frac{\mu_{j1}^2 - \mu_{j0}^2}{2\sigma_j^2}\right)}
 \end{aligned}$$

$$= \frac{1}{1 + \exp\left(-\left(-C_0\left(\frac{1-\eta}{\eta}\right) - \sum_{j=1}^d \frac{\mu_{j1} - \mu_{j0}}{2\sigma_j^2}\right) + \sum_{j=1}^d \frac{\mu_{j0} - \mu_{j1}}{\sigma_j^2} \cdot x_j\right)}$$

w w j x j !  
"Vector"

$$= \frac{1}{1 + \exp(-b + w^\top \cdot x)} \quad \checkmark$$

Problem 2)a) i) Simple case of  $X_1$  and  $X_2$ 

$$\Rightarrow \text{Var}(X_1 + X_2) = \text{Var}(X_1) + 2\text{COV}(X_1, X_2) + \text{Var}(X_2)$$

knowing that  $\text{COV}(X_1, X_2) = 0$  because they're independent variables,  $\text{Var}(X_1 + X_2) = \text{Var}(X_1) + \text{Var}(X_2) = \sum_{i=1}^2 \text{Var}(X_i)$

$$\Rightarrow \text{Var}\left(\frac{\sum_{i=1}^n X_i}{n}\right) = \frac{1}{n^2} \text{Var}\left(\sum_{i=1}^n X_i\right) = \frac{1}{n^2} \sum_{i=1}^n \text{Var}(X_i) \quad \begin{matrix} \text{(all } X_i \text{ s)} \\ \cancel{\text{indep.}} \end{matrix}$$

with Average

$$\text{Var}(cX) = c^2 \text{Var}(X) = \frac{1}{n^2} \cdot n \cdot \sigma_i^2, \quad \text{with all variances the same: } \sigma_i^2 = \sigma^2$$

b) Simple:  $P(\text{"data point not chosen in 1 draw"}) = 1 - \frac{1}{M}$  (single random draw from M)

$$\text{Drawing } M \text{ times: } P(\text{OOB}) = \left(1 - \frac{1}{M}\right)^M$$

$$\text{for large } M: \lim_{M \rightarrow \infty} \left(1 - \frac{1}{M}\right)^M = \frac{1}{e} \approx 0.37 \quad \checkmark$$

$$= \frac{1}{1 + \exp\left(-\left(-\text{Cov}\left(\frac{1-r}{n}\right) - \sum_{j=1}^d \frac{\mu_{j1} - \mu_{j0}}{2\sigma_j^2}\right)\right)} + \sum_{j=1}^d \frac{\mu_{j0} - \mu_{j1}}{\sigma_j^2} \cdot x_j$$

"Vector"  $w_j$  in  $x_j$ !

$$\hat{=} \frac{1}{1 + \exp(-b + w^\top \cdot x)}$$

Problem 2)a) i) Simple case of  $X_1$  and  $X_2$ 

$$\Rightarrow \text{Var}(X_1 + X_2) = \text{Var}(X_1) + 2\text{COV}(X_1, X_2) + \text{Var}(X_2)$$

knowing that  $\text{COV}(X_1, X_2) = 0$  because they're independent variables,  $\text{Var}(X_1 + X_2) = \text{Var}(X_1) + \text{Var}(X_2) = \sum_{i=1}^2 \text{Var}(X_i)$

$$\Rightarrow \text{Var}\left(\frac{\sum_{i=1}^N X_i}{n}\right) = \frac{1}{n^2} \text{Var}\left(\sum_{i=1}^N X_i\right) = \frac{1}{n^2} \cdot \sum_{i=1}^n \text{Var}(X_i) \quad (\text{all } X_i \text{ s. indep.})$$

$\frac{1}{n^2} \cdot \sum_{i=1}^n \sigma_i^2$ , with all variances the same:  $\sigma_i^2 = \sigma^2$

with Average

$$\text{Var}(cX) = c^2 \text{Var}(X)$$

$$= \frac{1}{n^2} \cdot n \cdot \sigma^2 = \frac{\sigma^2}{n}$$

ii) Simple:  $P(\text{"data point not chosen in 1 draw"}) = 1 - \frac{1}{M}$  (single random draw from M)

$$\text{Drawing } M \text{ times: } P(\text{OOB}) = \left(1 - \frac{1}{M}\right)^M$$

$$\text{for large } M: \lim_{M \rightarrow \infty} \left(1 - \frac{1}{M}\right)^M = \frac{1}{e} \approx 0.37 \quad \checkmark$$

- a iii) The OOB error works the following way:
- 1) generate  $B$  samples of length  $n$  of dataset with length  $n$ .
  - 2) Train a model ~~itself~~ on each sample  $B_1, B_2$ , etc.
  - 3) Deliberately make models do a prediction on an observation that they didn't include in training set. sample.
  - 4) Average that prediction, and compare to real value  $\triangleq$  OOB error for ~~predict~~ single observation
  - 5) Do for all ~~all~~ observations  
 $\hookrightarrow \text{OOB error} = \frac{\text{Falsely pred. observations (OOB)}}{\text{all observations (OOB)}}$
- The OOB error decreases for a given model until with the more trees it has. It's better able to classify observations it "didn't know".
  - At some point ~~stabilizing~~, increasing the # of trees will become computationally more expensive, ~~fast~~ while the OOB error only decreases slightly / marginally less
  - Choose the # of trees when OOB isn't decreasing significantly anymore.

2a iv)

SB

Now for correlated random vars  $X_1, X_2, \dots, X_n$  Nor 3rd

2021

$$\Rightarrow \text{Var}\left(\frac{\sum_{i=1}^n X_i}{n}\right) = \frac{1}{n^2} \cdot \text{Var}\left(\sum_{i=1}^n X_i\right)$$

$$= \frac{1}{n^2} \cdot \left( \sum_{i=1}^n \text{Var}(X_i) + \sum_{i=1}^n \sum_{i < j} \text{Cov}(X_i, X_j) \right)$$

$$\text{upper bound: } \leq \frac{1}{n^2} \cdot (n \cdot \sigma_i^2 + n \cdot (n-1) \cdot \delta^2) \quad (\text{with } \text{Cov}(X_i, X_j) \leq \delta^2)$$

$$\leq \frac{1}{n^2} \cdot (n \cdot \sigma^2 + (n^2 - n) \delta^2) \text{ with } \sigma_i^2 = \sigma^2$$

$$\leq \frac{1}{n^2} \cdot (n \cdot \sigma^2 + (n^2 - n) \delta^2) \quad \text{OK}$$

$$\leq \cancel{\frac{\sigma^2}{n}} + \frac{n^2 \delta^2 - n \delta^2}{n^2} = \frac{\sigma^2}{n} + \delta^2 - \frac{\delta^2}{n}$$

$$\leq \frac{1}{n} \cdot \sigma^2 + \left(1 - \frac{1}{n}\right) \cdot \delta^2$$

v) with  $N \rightarrow \infty$ :

$$\lim_{n \rightarrow \infty} \frac{1}{n} \sigma^2 + \left(1 - \frac{1}{n}\right) \cdot \delta^2 = 0 + (1-0) \delta^2 = \delta^2$$

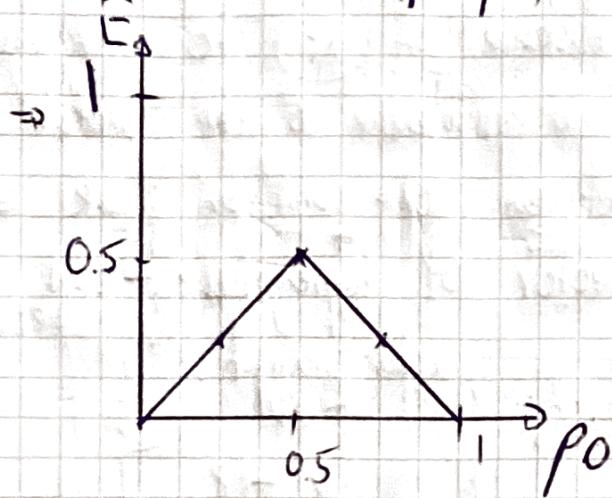
$$\text{from part i: } \lim_{n \rightarrow \infty} \frac{\sigma^2}{n} = 0$$

$\Rightarrow$  Bagging is limited by the covariance of the dataset's features or random variables (which we see in the trees' covariance).

$\Rightarrow$  Bagging is limited by the degree of correlation between the trees.

25)  $E = 1 - \max_{k \in \{0,1\}} p_k$  with  $k \in \{0,1\}$   
 $= 1 - \max(p_1, p_0)$  with  $p_0 = 1 - p_1$

SB  
Nr 3<sup>rd</sup>  
2021



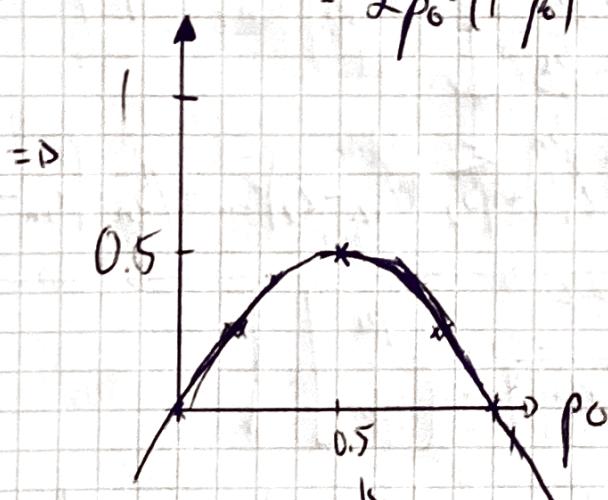
(classification error rate)

! we use  $p_0$  as first class (doesn't matter for plots)

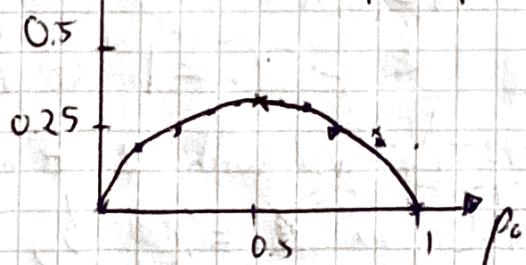
Gini-Index:

$$\sum_{k=1}^2 p_k(1-p_k) , \text{ starting at } k=0 \text{ with } k \in \{0,1\}$$

$$\Rightarrow \sum_{k=0}^1 p_k(1-p_k) = p_0(1-p_0) + p_1(1-p_1) = p_0 \cdot p_1 + p_1 \cdot p_0 = 2p_1 p_0 \\ = 2p_0 \cdot (1-p_0) = -2p^2 + 2p_0$$



Entropy:  $- \sum_{k=0}^1 p_k \log p_k = -(p_0 \log p_0 + p_1 \log p_1)$   
 $= -p_0 \log p_0 - (1-p_0) \log (1-p_0)$



### Problem 3a)

Mr 3rd  
SB

we want to calculate the distance of  $x_i$  to hyperplane (perpendicular distance). That is equivalent to the scalar projection of  $b$  (vector between some point on plane  $z$  and  $x_i$ ) on  $n$  (normal vector of hyperplane - unit normal vector)

$$d = |\text{comp}_n b| = \frac{|b \cdot n|}{\|n\|} \quad \text{with: scalar product.}$$

$$x_i = \begin{bmatrix} x_{i1} \\ x_{i2} \\ \dots \\ x_{in} \end{bmatrix}$$

$$z = \begin{bmatrix} z_1 \\ z_2 \\ z_3 \\ z_4 \\ \vdots \\ z_n \end{bmatrix}$$

$$\Rightarrow b = \begin{bmatrix} x_{i1} - z_1 \\ x_{i2} - z_2 \\ \dots \\ x_{in} - z_n \end{bmatrix}$$

and  $n = \frac{w}{\|w\|}$  with  $\langle w \rangle$  is  $\beta_1, \beta_2, \dots, \beta_n$  (weights of hyperplane)

$$\Rightarrow d = \frac{1}{\sqrt{\beta_1^2 + \beta_2^2 + \dots + \beta_n^2}} \cdot \left( \beta_1 \cdot (x_{i1} - z_1) + \beta_2 \cdot (x_{i2} - z_2) + \dots + \beta_n \cdot (x_{in} - z_n) \right)$$

$$\text{with } \left| \begin{bmatrix} \beta_1 \\ \beta_2 \\ \vdots \\ \beta_n \end{bmatrix} \right| := \left| \beta_1 \cdot x_{i1} + \beta_2 \cdot x_{i2} + \dots + \beta_n \cdot x_{in} - \beta_1 \cdot z_1 - \beta_2 \cdot z_2 - \dots - \beta_n \cdot z_n \right|$$

$$= \sqrt{\sum_{j=1}^n \beta_j^2}$$

$$\# \sqrt{\sum_{j=1}^n \beta_j^2} = 1 \quad (\text{constraint})$$

and with  $(-\beta_1 \cdot z_1 - \beta_2 \cdot z_2 - \dots - \beta_n \cdot z_n)$  with  $z \in \text{hyperplane}$

$$\Rightarrow \beta_1 \cdot z_1 + \beta_2 \cdot z_2 + \dots + \beta_n \cdot z_n + \beta_0 = 0 \Rightarrow \beta_0 = \beta_0$$

$$\Rightarrow d = |\beta_0 \cdot x_{i1} + \beta_1 \cdot x_{i2} + \dots + \beta_n \cdot x_{in} + \beta_c|$$

$$= |\beta_0 + \sum_{j=1}^n \beta_j \cdot x_{ij}|$$

now: max margin classifier : Classify based on sign of  $f(x_i)$  with  $f(x_i) = \beta_0 + \sum_{j=1}^n \beta_j \cdot x_{ij}$

$$\Leftrightarrow \text{sign}(y_i) = \text{sign}\left(\beta_0 + \sum_{j=1}^n \beta_j \cdot x_{ij}\right)$$

$$\Rightarrow d = y_i \cdot \left(\beta_0 + \sum_{j=1}^n \beta_j \cdot x_{ij}\right) = \left|\beta_0 + \sum_{j=1}^n \beta_j \cdot x_{ij}\right|$$

$$(-1 \cdot -1 = 1, 1 \cdot 1 = 1) \quad \checkmark$$

### Problem 35)

(i) given point  $x_i$ , with  $\varepsilon_i = 0$

$d = M \cdot (1 - \varepsilon_i) \Rightarrow d \geq M$   $x_i$  is on the correct side of margin (on either side) or, on the margin support vector ! at least  $\geq M$  away from hyperplane

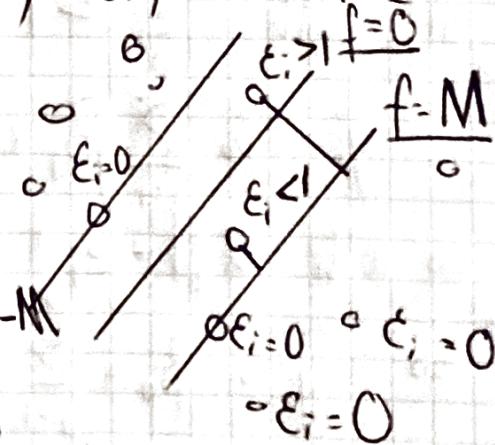
as  $x_i$  has a distance  $d$  to

Hyperplane, that is  $\geq M$ , we can

safely assume that it's classified

correctly  $\Rightarrow$  On correct side of hyperplane.

We know that the optimization problem will fulfill this constraint  $\text{sign}(y_i) \cdot \text{sign}(d_{\text{raw}}) \geq M \Rightarrow$  needs to be equal, sign  $\Rightarrow$  correct class!



3b ii)

$0 < \varepsilon_i < 1$  for point  $x_i$

$\Rightarrow$  is it on the correct side of the hyperplane?

$$\hookrightarrow y_i \cdot \left( \beta_0 + \sum_{j=1}^p \beta_j \cdot x_{ij} \right) \geq M \cdot (1 - \varepsilon_i) \quad \text{with } 0 < \varepsilon_i < 1$$

$$\Rightarrow 0 < y_i \cdot \left( \beta_0 + \sum_{j=1}^p \beta_j \cdot x_{ij} \right) < M$$

$\Rightarrow$  point  $x_i$  is classified correctly, (signs of  $y_i$  and  $\beta_0 + \sum_{j=1}^p \beta_j \cdot x_{ij}$  must match)  
 however, a loss incurred because  $\varepsilon_i > 0$ .

The point  $x_i$  is not on the correct side of the margin, but between the hyperplane and the margin support vector, with  $0 < d(x_i) < M$  to satisfy the inequality above.

iii) with  $\varepsilon_i > 1$  for  $x_i$  (one point.)

$$y_i \cdot \left( \beta_0 + \sum_{j=1}^p \beta_j \cdot x_{ij} \right) \geq M \cdot (1 - \varepsilon_i)$$

$$\Rightarrow \text{with } M \cdot (1 - \varepsilon_i) < 0 \text{ for } \varepsilon_i > 1$$

$$\Rightarrow y_i \cdot \left( \beta_0 + \sum_{j=1}^p \beta_j \cdot x_{ij} \right) < 0 \quad (\text{has a negative sign})$$

$\Rightarrow$  Misclassified, as  $\text{Sign}(y_i) \neq \text{Sign}(\beta_0 + \sum_{j=1}^p \beta_j \cdot x_{ij})$

$\Rightarrow$  Not on correct side of hyperplane,

$\Rightarrow$  Not on correct side of the margin!

SB  
Mar  
2021

Sc1

$$\text{Show that } f(x^*) = \beta_0 + \sum_{i=1}^n \alpha_i \cdot \langle x^*, x_i \rangle$$

$$\text{is equal to } f(x^*) = \beta_0 + \beta_1 x_1^* + \beta_2 x_2^* + \dots + \beta_p x_p^*,$$

with:  $\mathbb{R}^p$ , and  $x_i$  <sup>th</sup> notes the  $i^{\text{th}}$  training observation.

$\oplus \Rightarrow$  Find values for  $\beta_j$ ,  $j \in \{1, \dots, p\}$  in terms of  $\alpha_i$  &  $x_i$  for  $i \in \{1, \dots, n\}$

$$\begin{aligned} \Rightarrow f(x^*) &= \beta_0 + \sum_{i=1}^n \alpha_i \cdot \langle x^*, x_i \rangle \quad \text{with } \langle x^*, x_i \rangle \\ &= \beta_0 + \sum_{i=1}^n \alpha_i \left( x_1^* \cdot x_{i1} + x_2^* \cdot x_{i2} + \dots + \underbrace{x_p^* \cdot x_{ip}}_{(\text{Scalar product})} \right) \\ &= \beta_0 + \sum_{i=1}^n \alpha_i \cdot \sum_{j=1}^p x_j^* \cdot x_{ij} \quad (\text{we want } \sum_{j=1}^p \text{ outside}) \\ &= \beta_0 + \sum_{i=1}^n \sum_{j=1}^p \alpha_i \cdot x_j^* \cdot x_{ij} = \beta_0 + \sum_{j=1}^p x_j^* \sum_{i=1}^n \alpha_i \cdot x_{ij} \\ &\quad - \beta_0 + x_1^* \cdot \sum_{i=1}^n \alpha_i \cdot x_{i1} + x_2^* \cdot \sum_{i=1}^n \alpha_i \cdot x_{i2} + \dots + x_p^* \cdot \sum_{i=1}^n \alpha_i \cdot x_{ip} \\ \Rightarrow \underline{\beta_j} &= \sum_{i=1}^n \alpha_i \cdot x_{ij} \quad \text{for all } j \in \{1, \dots, p\} \\ &\stackrel{\cong}{=} \beta_0 + x_1^* \cdot \beta_1 + \dots + x_p^* \cdot \beta_p \quad \checkmark \checkmark \end{aligned}$$