

Learning in the Classification Setting

From the 0/1 loss to Convex Surrogate Losses

1) Linear Classification

- Learning with the 0/1 loss
- A surrogate learning problem: Logistic Regression
- The two views of logistic regression

2) CPE-Losses (Class Probability Estimate Losses)

- Definition of CPE-losses
- Proper losses
- the Cross-Entropy Loss

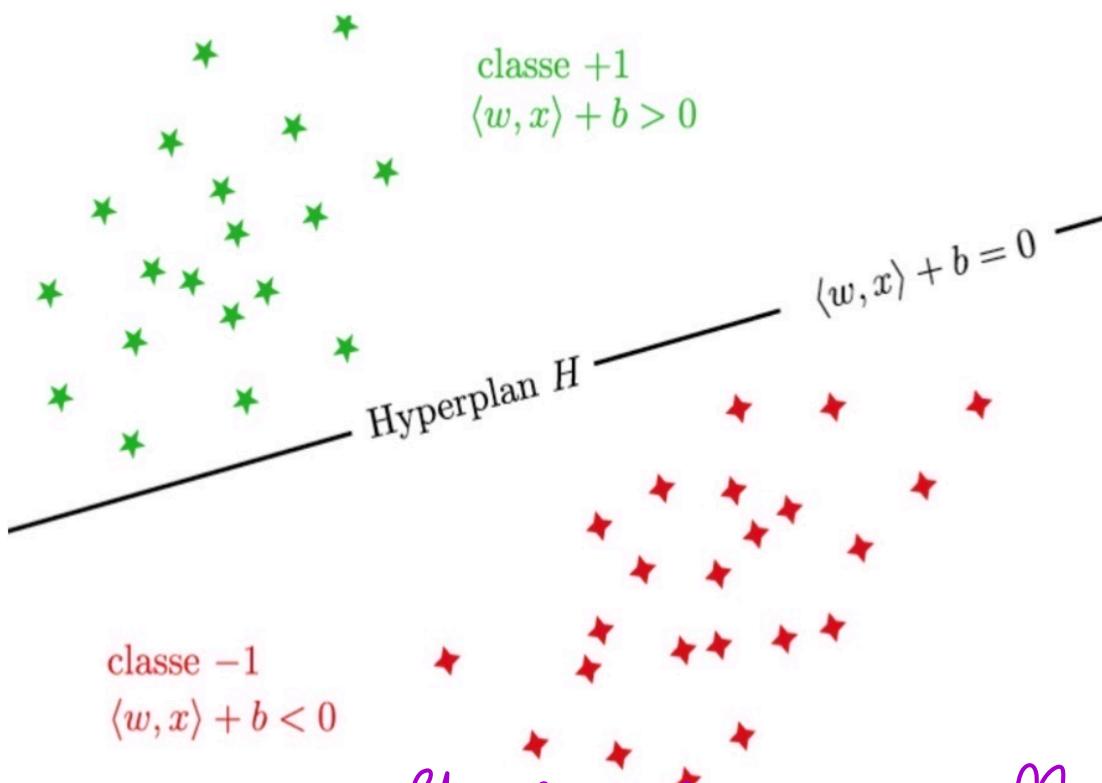
3) Scoring Losses

- Definition of Scoring Losses
- Well Calibrated Losses

1) Linear classification

Linear classification with the 0/1 loss

- $S = \{(x_1, y_1), \dots, (x_n, y_n)\}$, $\mathcal{X} = \mathbb{R}^d$, $y = \hat{y} = \{0, 1\}$
- $\mathcal{F} = \{x \mapsto \mathbb{1}_{[\theta^T x + b > 0]} : \theta \in \mathbb{R}^d, b \in \mathbb{R}\}$
- The 0/1 loss is $\ell^{0/1}(\hat{y}, y) = \mathbb{1}_{[\hat{y} \neq y]}$



⚠️ easy pb iff the classes are well separated
Hard (NP-hard) otherwise

notation:

$$\begin{aligned}\theta^T x &= \langle \theta, x \rangle \\ &= \sum_{j=1}^d \theta^{(j)} x^{(j)}\end{aligned}$$

learning problem:

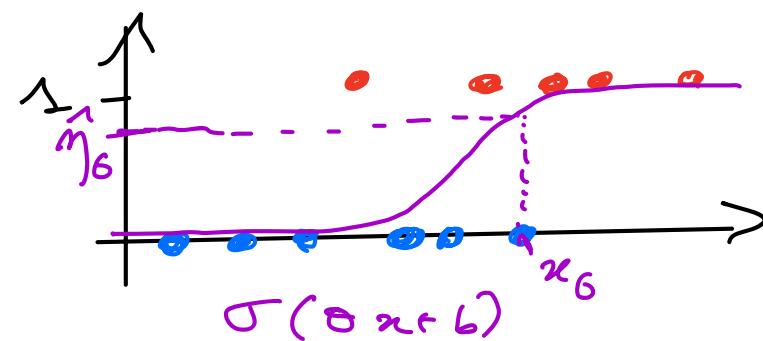
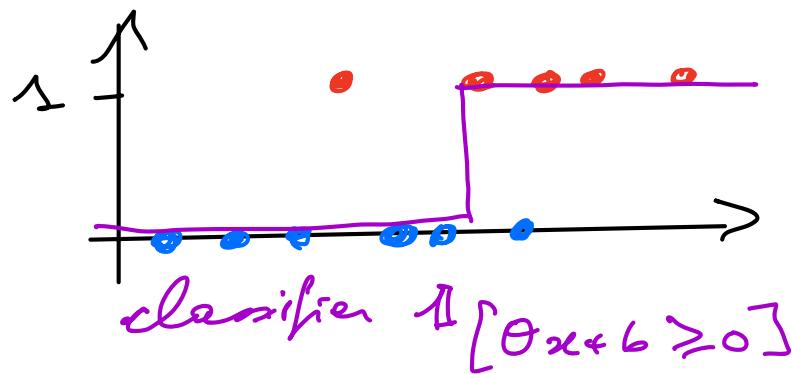
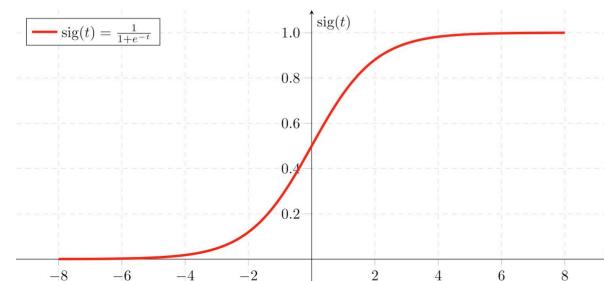
$$\underset{f \in \mathcal{F}}{\operatorname{argmin}} \sum_{i=1}^n \underbrace{\ell^{0/1}(f(x_i), y_i)}_{\mathbb{1}_{[f(x_i) \neq y_i]}}$$

$$\underset{\theta, b}{\operatorname{argmin}} \sum_{\substack{h_i : y_i = 1 \\ h_i : y_i = 0}} \mathbb{1}_{[\theta^T x_i + b < 0]} + \sum_{\substack{h_i : y_i = 1 \\ h_i : y_i = 0}} \mathbb{1}_{[\theta^T x_i + b \geq 0]}$$

Linear classification with the Logistic regression framework

- Idea: replace the binary classifier with a continuous function whose parameters we can optimize as a convex opt pb.
- The sigmoid function

$$\sigma(t) = \frac{1}{1+e^{-t}}$$



- why the sigmoid function?

Linear classification with logistic regression : which loss ?

$$y = \{0, 1\}$$

Notations: (dependency on θ, b , omitted for brevity)

$$\hat{y}_i = \theta^T x_i + b \leftarrow \text{Score for point } x_i$$

$$\hat{\eta}_i := \sigma(\hat{y}_i) = \frac{1}{1+e^{-\hat{y}_i}} \leftarrow \text{the class probability estimate (probits)}$$

$$\hat{c}_i := \mathbb{1}[\hat{y}_i \geq 0] = \mathbb{1}[\hat{\eta}_i \geq \frac{1}{2}] \leftarrow \text{predicted class.}$$

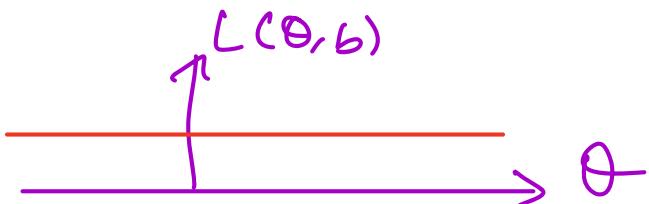
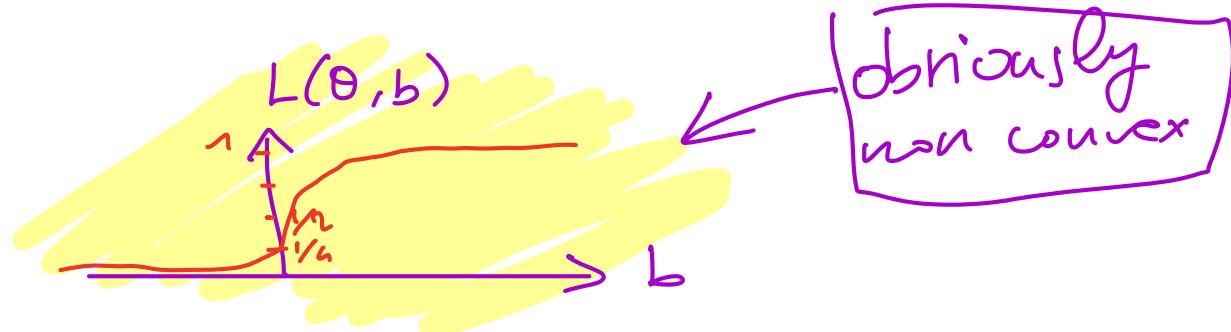
- $\hat{\eta}_i$ is interpreted as $P[y=1 | X=x_i; \theta, b]$
- $(1 - \hat{\eta}_i)$ " " " $P[y=0 | X=x_i; \theta, b]$

- Pb : learn θ, b

- Very naive Idea : learn $\hat{\theta}, \hat{b} = \operatorname{argmin} \sum_{i=1}^N \mathbb{1}[\hat{c}_i \neq y_i]$
Naive because identical to previous problem, so hard
- Naïve Idea : learn $\hat{\theta}, \hat{b} = \operatorname{argmin} \sum_{i=1}^N (\hat{\eta}_i - y_i)^2$
Naïve because the objective is continuous, but not convex \Rightarrow hard to solve.

Exercise

Show with a data set of one example $(0, 0)$ that the objective is not convex-

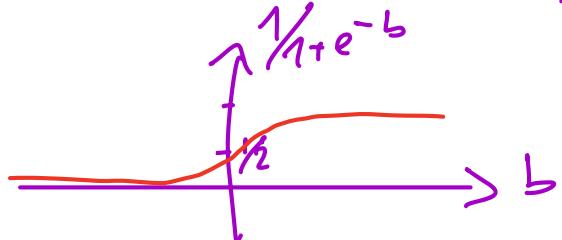


$$\text{argmin} \sum_{i=1}^N (\hat{y}_i - y_i)^2$$

recall: $\hat{y}_i = \theta^T x_i + b$

$$\hat{y}_i = \sigma(\hat{g}_i) = \frac{1}{1+e^{-\hat{g}_i}}$$

$$\text{Here, } L(\theta, b) = \left(\frac{1}{1+e^{-b}} \right)^2$$



Linear classification with logistic regression : which loss ?

- We will use the probabilistic interpretation

of logistic regression to get a good loss.

- The likelihood of the model θ, b is

$$L(\theta, b) = \prod_{i=1}^N P(Y=y_i | X=x_i; \theta, b) = \prod_{\{i: y_i=1\}} \hat{\pi}_i \times \prod_{\{i: y_i=0\}} (1-\hat{\pi}_i)$$

- Making the data the most probable w.r.t. θ, b amounts to maximize $L(\theta, b)$ or to minimize the negative log likelihood $NLL(\theta, b)$

$$NLL(\theta, b) = -\log L(\theta, b)$$

$$= - \sum_{\{i: y_i=1\}} \ln \hat{\pi}_i - \sum_{\{i: y_i=0\}} \ln (1-\hat{\pi}_i)$$

$$= \sum_{i=0}^N -y_i \ln \hat{\pi}_i - (1-y_i) \ln (1-\hat{\pi}_i)$$

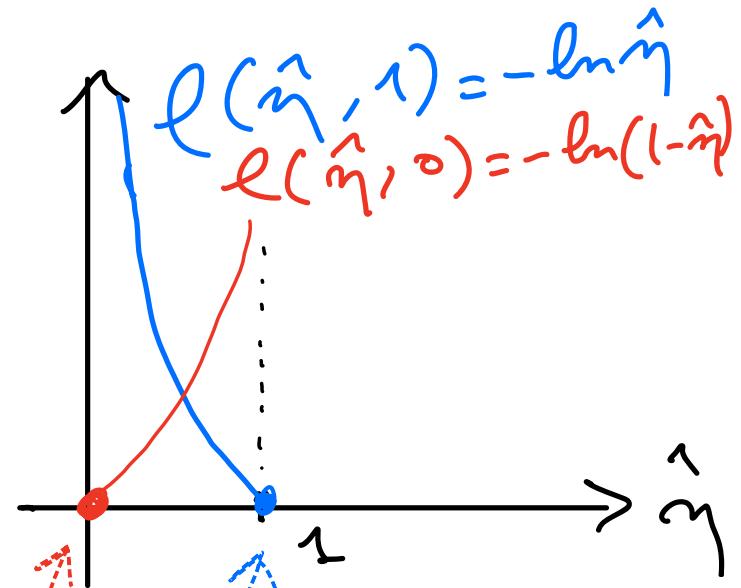
$\mathcal{L}_{CE}(\hat{\pi}_i, y_i)$ = cross-entropy loss

Logistic regression pb

$$\hat{\theta}, \hat{b} = \arg \min_{\theta, b} \sum_{i=1}^N \mathcal{L}_{CE}(\hat{\pi}_i, y_i)$$

Exercise

- draw $l^{CE}(\hat{\eta}, 1)$ and $l^{CE}(\hat{\eta}, 0)$
- Based on the drawing,
show $\underset{\hat{\eta}}{\operatorname{argmin}} l^{CE}(\hat{\eta}, 0) = 0$
and $\underset{\hat{\eta}}{\operatorname{argmin}} l^{CE}(\hat{\eta}, 1) = 1$

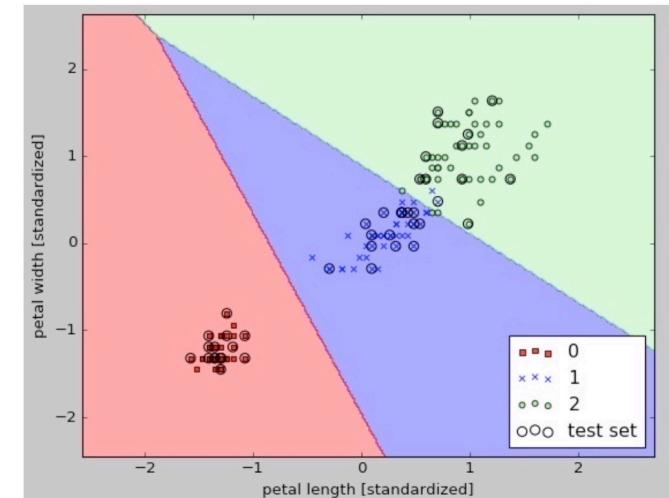


logistic regression : the multiclass setting

- $y = \{1..K\}$
- for each $b_k \in \mathcal{Y}$, we have $\theta_k \in \mathbb{R}^d$, $b_k \in \mathbb{R}^d$
- for each example x_i , the scores of x_i for each class is
 $(\hat{y}_{i,1}, \dots, \hat{y}_{i,K}) = (\theta_1^T x_i + b_1, \dots, \theta_K^T x_i + b_K)$
- predicted class: $\hat{c}_i = \arg \max_{k \in \{1..K\}} \hat{y}_{i,k}$
- replace sigmoid by softmax:

$$\text{softmax}\left(\begin{matrix} t_1 \\ \vdots \\ t_K \end{matrix}\right) = \frac{1}{\sum_{k=1}^K e^{t_k}} \left(\begin{matrix} e^{t_1} \\ \vdots \\ e^{t_K} \end{matrix} \right)$$
- Estimated class probability $\hat{\eta}_i = (\hat{\eta}_{i,1}, \dots, \hat{\eta}_{i,K}) = \text{softmax}(\hat{y}_{i,1}, \dots, \hat{y}_{i,K})$
- multiclass cross entropy loss

$$\ell_{\text{CE}}(\hat{\eta}_i, y_i) = - \sum_{k=1}^K \mathbb{1}_{[y_i=k]} \ln \hat{\eta}_{i,k}$$



logistic regression : the multiclass setting, example

three classes: $\mathcal{Y} = \{\text{dog, cat, mouse}\}$
 $x = (2, 3)$

parameters

$$\theta_{\text{dog}} = (1, 1)^T$$

$$b_{\text{dog}} = 0$$

$$\theta_{\text{cat}} = (-1, 0)^T$$

$$b_{\text{cat}} = 0$$

$$\theta_{\text{mouse}} = (1, -1)^T$$

$$b_{\text{mouse}} = 0$$

Scores

$$\hat{y}_{\text{dog}} = 5$$

$$\hat{y}_{\text{cat}} = -2$$

$$\hat{y}_{\text{mouse}} = -1$$

CPRs

$$\exp(\hat{y}_{\text{dog}}) = 148$$

$$\exp(\hat{y}_{\text{cat}}) = 0,1$$

$$\exp(\hat{y}_{\text{mouse}}) = 0,4$$

$$\sum_c \exp(\hat{y}_c) = 149$$

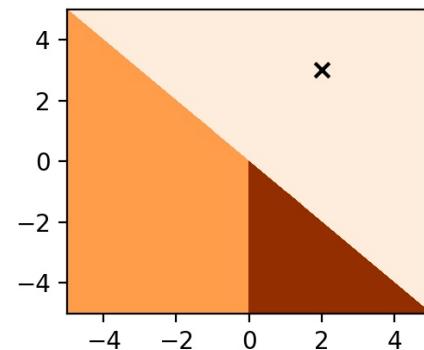
$$\hat{y}_{\text{dog}} = 99,6\%$$

$$\hat{y}_{\text{cat}} = 0,1\%$$

$$\hat{y}_{\text{mouse}} = 0,2\%$$

predicted class

$$\hat{c} = \text{dog}$$



The logistic regression: the two views (in the bi-class setting)

View 1: Working with cross-entropy loss $l^{CE}(\hat{y}_i, y_i)$ with $y \in \{0, 1\}$
which is a CPE-loss (class probability estimate)

$$\hat{\theta}, \hat{b} = \underset{\theta, b}{\operatorname{arg\,min}} \sum_i l^{CE}(\hat{y}_i, y_i)$$

View 2: work with the logistic loss on scores with $y \in \{-1, 1\}$
The logistic loss $l^{\text{logistic}}(\hat{y}_i, y_i)$ is
a scoring loss

$$\hat{\theta}, \hat{b} = \underset{\theta, b}{\operatorname{arg\,min}} \sum_i l^{\text{logistic}}(\hat{y}_i, y_i)$$

$$\text{with } l^{\text{logistic}}(\hat{y}_i, y_i) = \ln(1 + e^{-y_i \hat{y}_i})$$

Both views are identical, and are convex in θ, b

Exercise: develop view 1 to show it is identical to view 2.

The logistic Regression: the two views

Exercise: develop view 1 to show it is identical to view 2.

Solution:

The logistic Regression: the two views

Exercise: develop view 1 to show it is identical to view 2.

Solution:

$$\ell^{CE}(\hat{\eta}, 1) = -\ln \hat{\eta} = -\ln \frac{1}{1+e^{-\hat{\eta}}} = \ln(1+e^{-\hat{\eta}})$$

$$\ell^{CE}(\hat{\eta}, 0) = -\ln(1-\hat{\eta}) = -\ln\left(1 - \frac{1}{1+e^{-\hat{\eta}}}\right) = -\ln \frac{e^{-\hat{\eta}}}{1+e^{-\hat{\eta}}}$$

$$= -\ln\left(\frac{1}{1+e^{\hat{\eta}}}\right) = \ln(1+e^{\hat{\eta}})$$

$$\ell^{CE}(\hat{\eta}, y) = \ln(1+e^{-y\hat{\eta}})$$

Cross
entropy loss

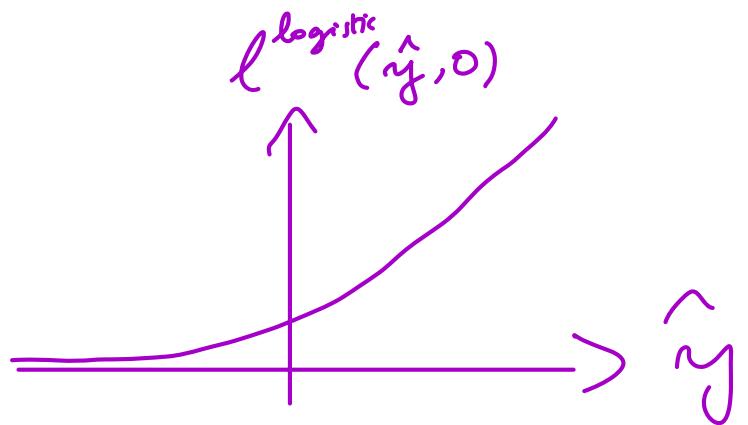
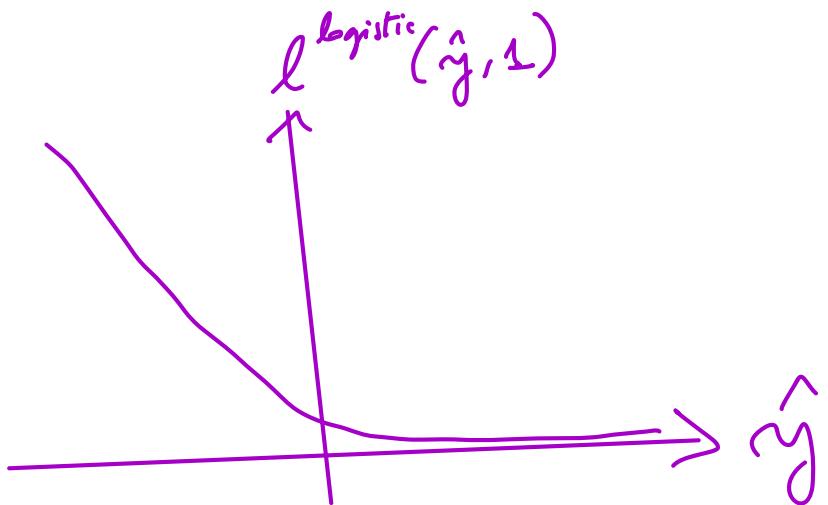
$E_{0,1}$

\uparrow
 $\in \{-1, 1\}$

Logistic loss

Convexity of the learning problems

$l^{\text{logistic}}(\hat{y}, y) = \ln(1 + e^{-y\hat{y}})$ is convex in \hat{y}
 $\in [-1, 1]$



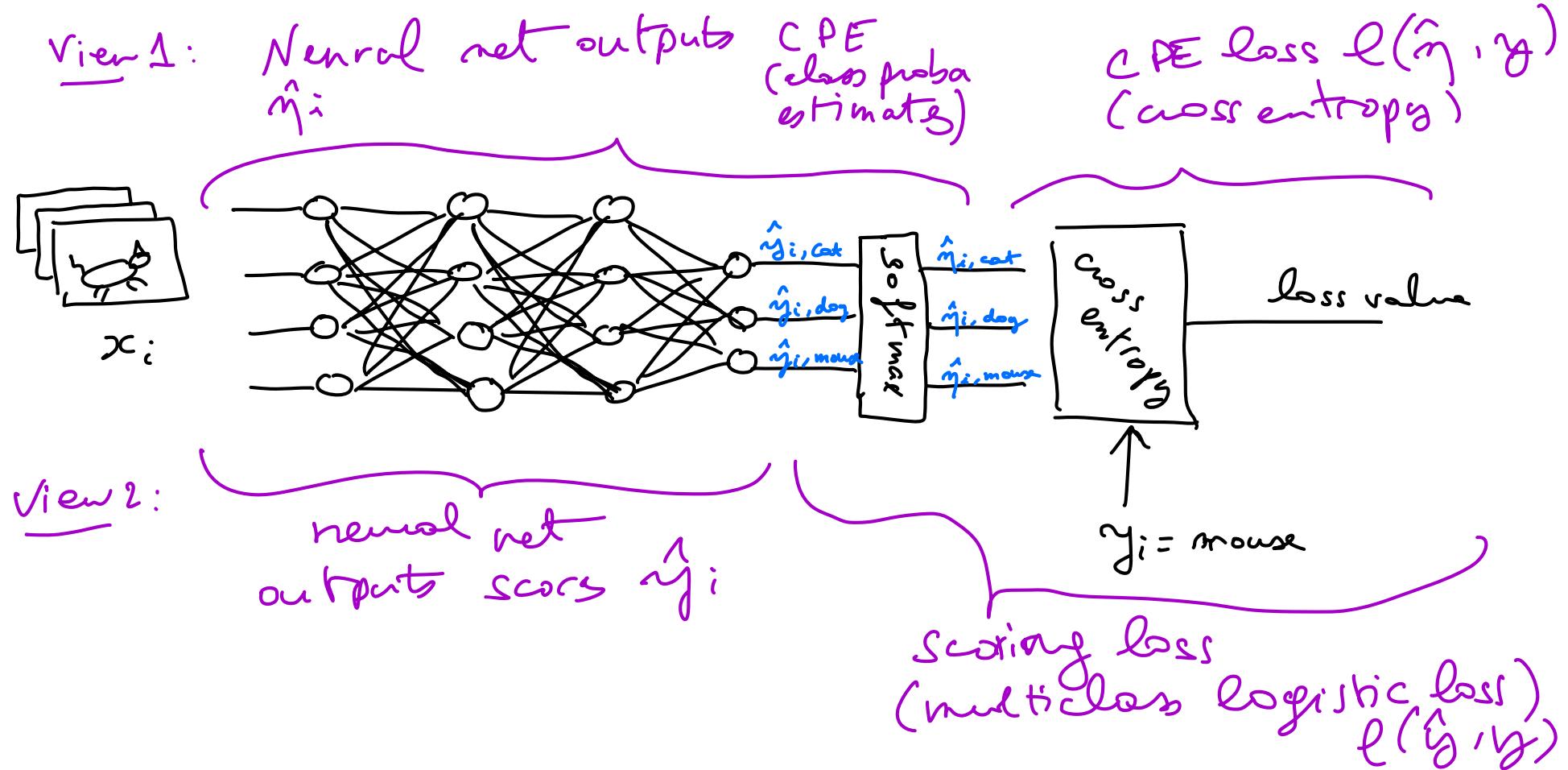
The overall objective

$$\sum_{i=1}^N l(\hat{y}_i, y_i) = \sum_{i=1}^N l^{\text{logis}}(\theta^T x_i + b, y_i)$$

is convex in θ and b

From logistic regression to neural nets: the two views

- We train a neural net on 64×64 images with $y = \{\text{cat, dog, mouse}\}$
 $X = \mathbb{R}^{64 \times 64}$



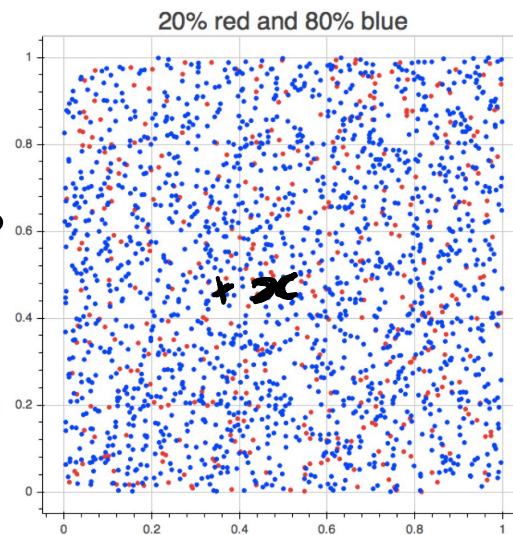
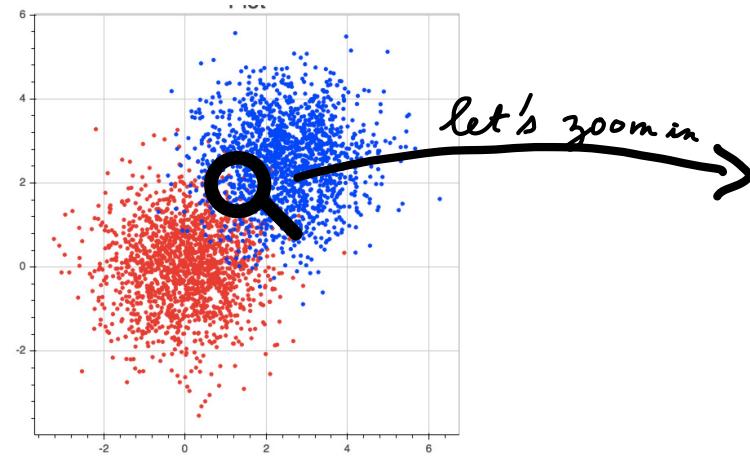
2) CPE losses

Class
Probability
Estimated

2) CPE-losses (for $\gamma = \{0, 1\}$)

- Recall that a CPE loss is of the form $\ell(\hat{\gamma}, \gamma)$ (e.g. cross entropy).
- What are good CPE-losses?

$$\gamma = \{1, 0\}$$



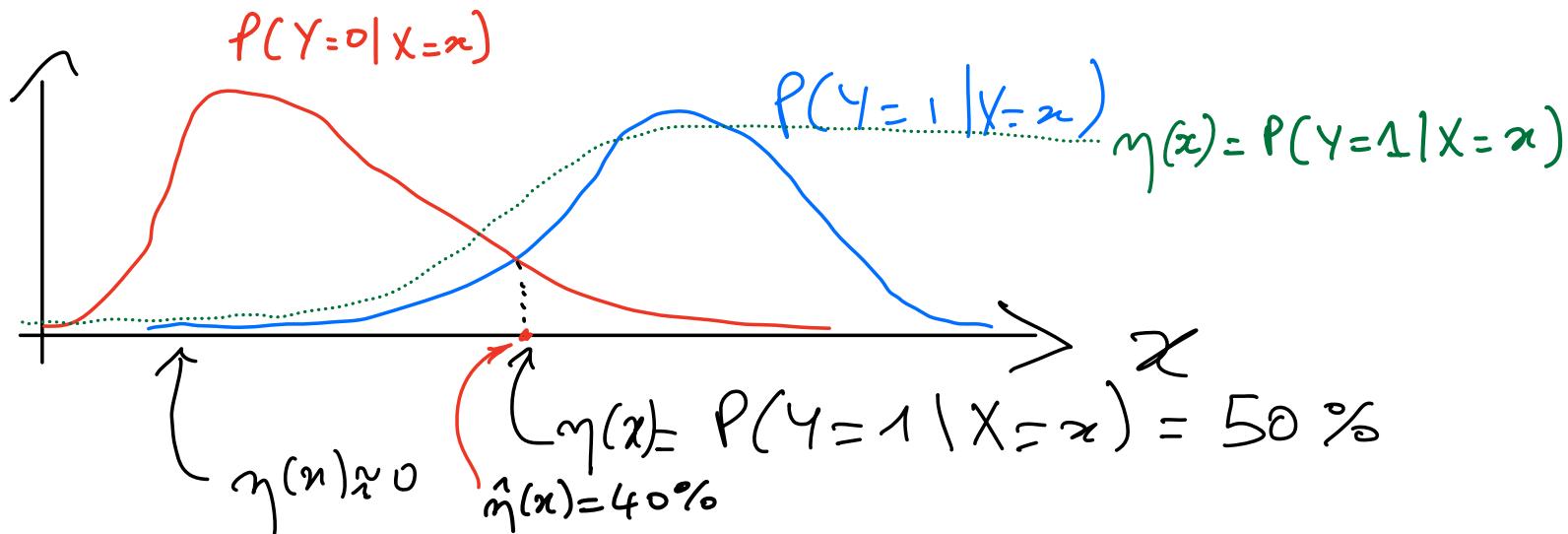
- The point x is in a region where $\gamma = p(\gamma=1|x) = 80\%$
- So a neural network trained with a "good" CPE loss should predict the correct probability $\hat{\gamma} = 80\%$
- The average loss incurred by points in this region will be $\gamma \ell(\hat{\gamma}, 1) + (1-\gamma) \ell(\hat{\gamma}, 0)$

3) Risk of CPE-losses

Notation:

- $\hat{\eta} \in [0, 1]$ is a probability
 - $\hat{\eta}(\cdot) \in \mathbb{R}^d \mapsto [0, 1]$ is a CPE prediction function (e.g. a neural net with softmax)
 - let us first adapt our definition of the risk to new losses
- def: the risk of a prediction function $\hat{\eta}(\cdot)$ with CPE-loss ℓ is
- $$R^\ell(\hat{\eta}(\cdot)) = \mathbb{E}_{x,y \sim P} [\ell(\hat{\eta}(x), y)]$$
- To capture the intuition of the previous slide,
we need the conditional risk

4) conditional risk of CPE losses



definition: for $\eta \in [0,1]$ and $\hat{\eta} \in [0,1]$, the conditional risk is:

$$C(\hat{\eta}, \eta) = \mathbb{E}_{Y \sim \text{Ber}(\eta)} [\ell(\hat{\eta}, Y)] = \eta(x) \times \ell(\hat{\eta}(x), 1) + (1-\eta(x)) \times \ell(\hat{\eta}(x), 0)$$

$$\begin{aligned} R^\ell(\hat{\eta}(\cdot)) &= \mathbb{E}_{X,Y \sim P} [\ell(\hat{\eta}(X), Y)] \\ &= \mathbb{E}_X \left[\mathbb{E}_Y [\ell(\hat{\eta}(X), Y) | X] \right] \\ &= \mathbb{E}_X \left[\eta(x) \times \ell(\hat{\eta}(x), 1) + (1-\eta(x)) \times \ell(\hat{\eta}(x), 0) \right] \\ &= \mathbb{E}_X [C(\hat{\eta}(x), \eta(x))] \end{aligned}$$

5) Proper CPE losses

- A "good" CPE loss will be a proper loss

def: a CPE loss is proper iff

$$\forall \gamma \in [0,1] \quad \gamma \in \arg\min_{\hat{\gamma}} L(\hat{\gamma}, \gamma)$$

Exercises

1) - show the cross entropy is proper

let us show that $\eta = \underset{\hat{\eta}}{\operatorname{arg\min}} C(\hat{\eta}, \eta)$

$$C(\hat{\eta}, \eta) = -\eta \ln \hat{\eta} - (1-\eta) \ln (1-\hat{\eta})$$

$$\frac{d}{d\hat{\eta}} C = -\frac{\eta}{\hat{\eta}} + \frac{(1-\eta)}{1-\hat{\eta}} \triangleq 0 \Rightarrow \frac{\eta}{\hat{\eta}} = \frac{1-\eta}{1-\hat{\eta}}$$

$$\Rightarrow \eta(1-\hat{\eta}) = \hat{\eta}(1-\eta) \Rightarrow \eta = \hat{\eta} \Rightarrow l^{CE} \text{ is (strictly) proper}$$

2) - show (latter) that the MSE mean square error is strictly proper.

3) - Write the Bayes risk ($= \min_f R^e(f)$)

for the case l is proper. Apply to cross entropy

$$\inf_f R^l(f) = \mathbb{E}_X \inf_{\hat{\eta}} C(\hat{\eta}, \eta(x)) = \mathbb{E}_X C(\eta(x), \eta(x))$$

cross entropy, $\mathbb{E}_X C(\eta(x), \eta(x)) = \mathbb{E}_X H(\eta(x))$ where $H(\cdot)$ is entropy

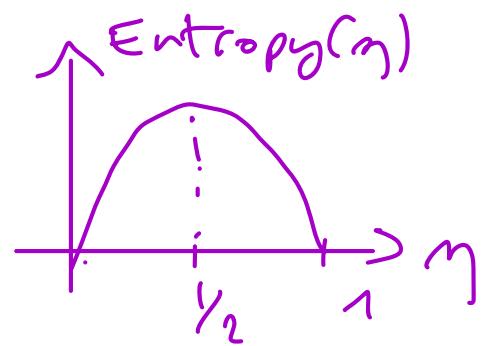
6) Note on CPE-losses and the cross-entropy

- $l(\hat{\eta}, y)$ in [40, 14], could also be interpreted as a parameter of Bernoulli:
parameter of a Bernoulli distribution

ex: $l(10\%, 0)$ can be interpreted as a dissimilarity between two distributions $\text{Ber}(10\%)$ and $\text{Ber}(0\%)$

- There is a standard way to compare discrete distributions: Kullback Leibler divergence for two discrete distributions p and q on \mathcal{Y} ,
$$KL(p \parallel q) = \sum_{y \in \mathcal{Y}} p(y) \ln \frac{p(y)}{q(y)}$$
- property: $KL(p \parallel p) = 0$, $KL(p \parallel q) > 0$ if $p \neq q$
 $KL(p \parallel q) \neq KL(q, p)$

$$KL(p \parallel q) = \sum_{y \in Y} p(y) \ln \frac{p(y)}{q(y)}$$



I compare $\text{Ber}(\hat{y})$ with $\text{Ber}(y)$

$$\begin{aligned} KL(\text{Ber}(y), \text{Ber}(\hat{y})) &= y \times \ln \frac{\hat{y}}{y} + (1-y) \ln \frac{1-\hat{y}}{1-y} \\ &= \underbrace{-y \ln \hat{y} - (1-y) \ln (1-\hat{y})}_{\ell^{CE}(\hat{y}, y)} + y \ln y + (1-y) \ln (1-y) \\ &\quad - \underbrace{\text{Entropy}[\text{Ber}(y)]}_{=0} = 0 \end{aligned}$$

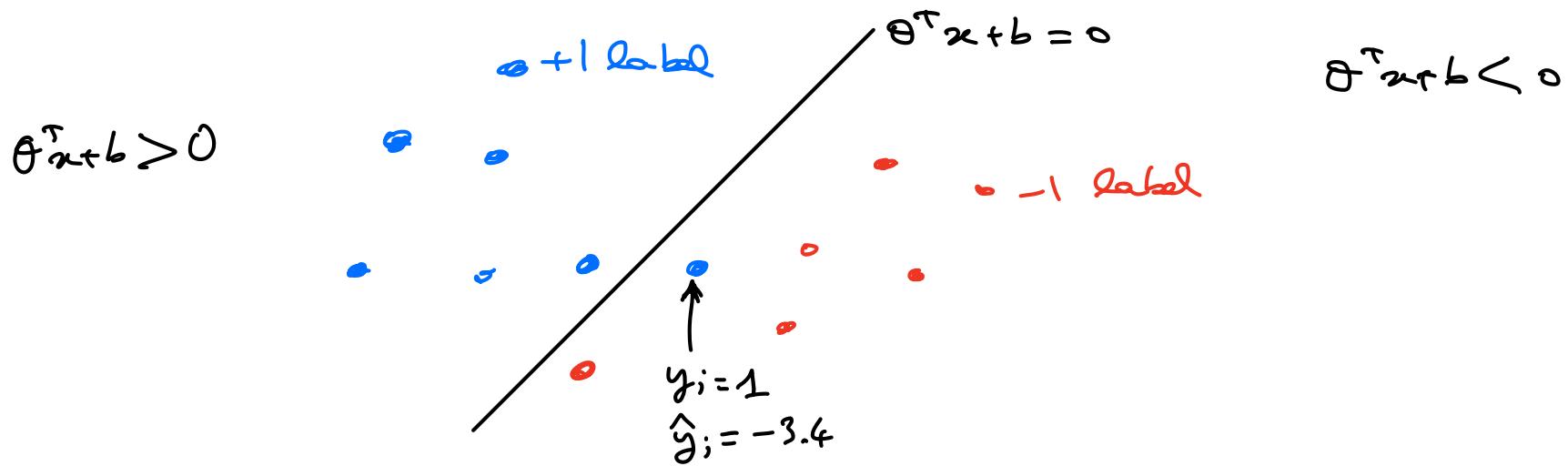
(Assumption $0 \ln 0 = 0$)

$$\ell^{CE}(\hat{y}, y) = KL(\text{Ber}(y), \text{Ber}(\hat{y}))$$

3. Scoring Losses

3) Scoring losses

- $y = \{-1, 1\}$
- $\ell(\hat{y}, y)$ ex: logistic loss
- With linear classifiers, $\hat{y}_i = \theta^T x_i + b$



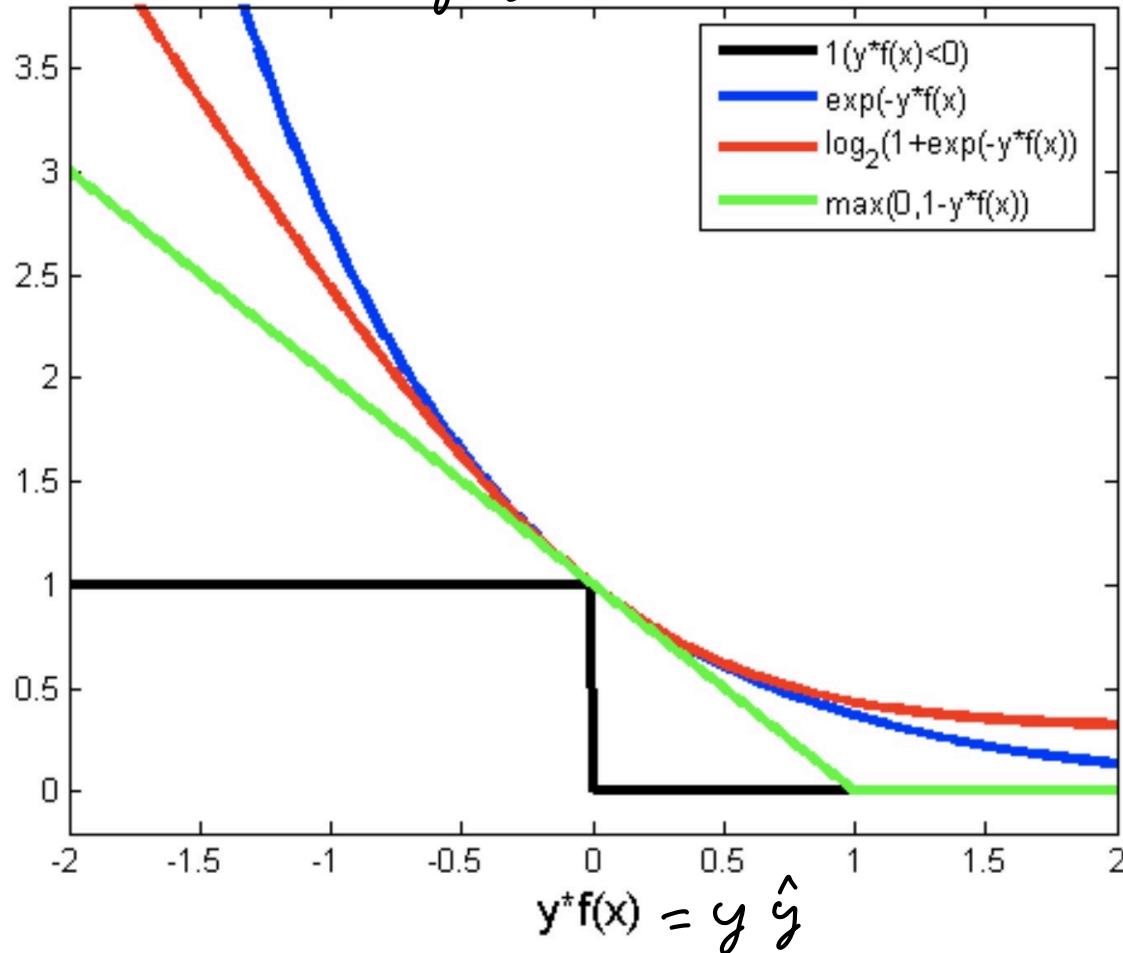
- Most scoring losses can be formulated this way: $\ell(\hat{y}, y) = \phi(\hat{y} - y)$ for some function ϕ . These are called ϕ -losses or margin losses
- Examples of ϕ -losses

Scoring losses

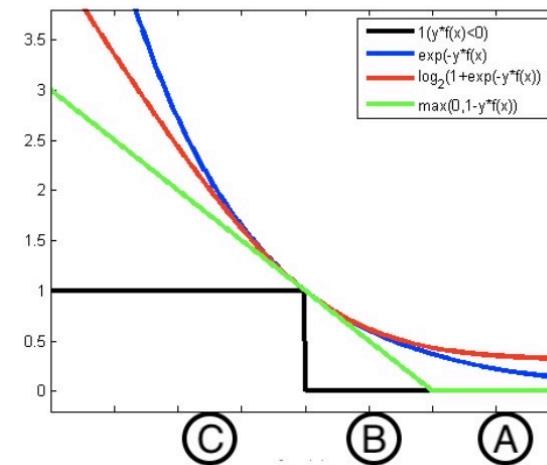
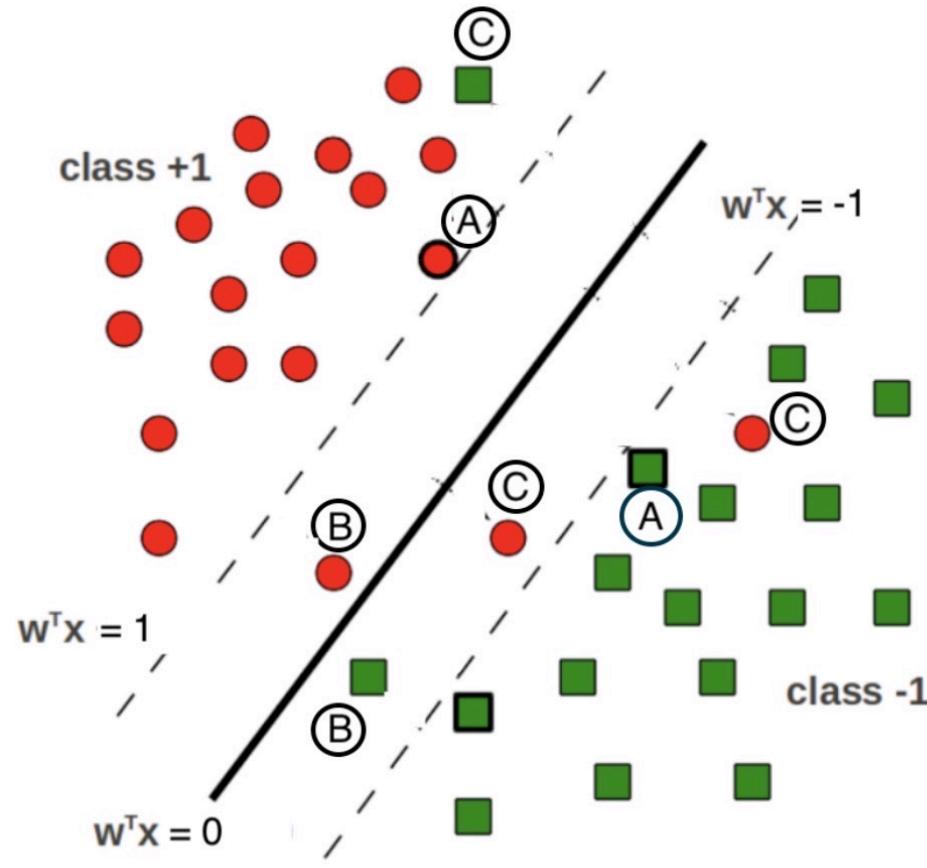
$\ell_{0/1}(\hat{y}, y) = \mathbb{1}_{[y\hat{y} < 0]}$ is a ϕ -loss with
 $\phi(\cdot) = \mathbb{1}_{[\cdot < 0]}$

$\ell_{\text{logistic}}(\hat{y}, y) = \ln(1 + e^{-y\hat{y}})$ is a ϕ -loss

$\ell_{\text{hinge}}(\hat{y}, y) = \max(1 - y\hat{y}, 0)$



Scoring losses



The risk of a scoring loss

We note $\eta(x) = P(Y=1 | X=x)$

$\hat{y} \in \mathbb{R} \cup \{-\infty, \infty\}$ is a score and $\hat{\eta}(\cdot) \in \mathcal{X} \mapsto \mathbb{R} \cup \{-\infty, \infty\}$ is a scoring function (e.g. a neural net)

def: For any $\hat{y} \in \mathbb{R}$, the conditional risk for scoring loss l is

$$C(\hat{y}, \eta) = \mathbb{E}_{Y \sim \begin{cases} 1 \text{ w.p. } \eta \\ -1 \text{ w.p. } 1-\eta \end{cases}} [l(\hat{y}, Y)] \\ = \eta l(\hat{y}, 1) + (1-\eta) l(\hat{y}, -1)$$

def: The risk of scoring function $\hat{\eta}(\cdot): \mathbb{R}^d \rightarrow \mathbb{R}$ for loss l is

$$R^l(\hat{\eta}(\cdot)) = \mathbb{E}_{(X, Y) \sim P} [l(\hat{\eta}(x), y)] = \mathbb{E}_X C(\hat{\eta}(x), \eta(x))$$

Calibration of scoring losses

Same as before: what are "good" scoring losses?

def: A loss l is calibrated if the following holds:

- if $\eta \in [0, \frac{1}{2}[$ then $\inf_{\hat{y} < 0} C(\hat{y}, \eta) < \inf_{\hat{y} \geq 0} C(\hat{y}, \eta)$
- if $\eta \in]\frac{1}{2}; 1]$ then $\inf_{\hat{y} > 0} C(\hat{y}, \eta) < \inf_{\hat{y} \leq 0} C(\hat{y}, \eta)$

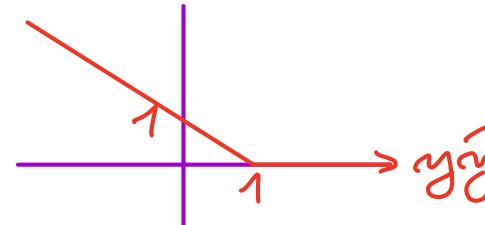
Thm: if l is calibrated then

if $\hat{y}(\cdot)$ is the measurable function
minimizing R^l ,
then it also minimizes R^{ors}

Intuitively: A measurable function learnt to minimize a
calibrated loss will also minimize the O/I loss

Calibration of scoring losses

Exercise:

- let $l^{\text{hinge}}(\hat{y}, y) = \max(0, 1 - \hat{y}y)$ = 

- we will consider 3 cases in this exercise :

case 1: $\eta < \frac{1}{2}$

case 2: $\eta = \frac{1}{2}$

case 3: $\eta > \frac{1}{2}$

- Draw $C(\hat{y}, \eta)$ as a function of \hat{y} in each of these three cases .
- In each case, show which value of \hat{y} minimizes $C(\hat{y}, \eta)$
- Also show which predicted class corresponds to these \hat{y}
- If, instead of the hinge loss, we used the 0/1 loss, which class would be the optimal one in these 3 cases ?
- Using the definition of calibration, show the hinge loss is calibrated

Calibration of scoring losses

Thm: Any convex ϕ -loss with
 $\phi'(0) < 0$ is well calibrated

\implies all scoring losses we saw
are well-calibrated