

Linear classification with 0/1 loss

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

learning problem: $\underset{\theta, b}{\operatorname{argmin}} \sum_{i=1}^N \underbrace{\ell^{0,1}(f(x_i), y_i)}_{\mathbb{1}(f(x_i) \neq y_i)}$

error class 0 + err class 1

easy if classes are well separated

hard otherwise (NP-hard)

- logistic regression (replace the binary classification with sigmoid $\sigma(t) = \frac{1}{1+e^{-t}}$ a continuous function whose params we can opt on a convex opt pb)

$y_i = \{0, 1\}$
 $\hat{y}_i = \theta^T x_i + b \leftarrow \text{score for posit } x_i$

notation $\hat{y}_i = \sigma(\hat{y}_i) = \frac{1}{1+e^{-\hat{y}_i}} \in \text{the class probab entries (probits)}$

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

\hat{y}_i is interpreted as $P[y=1 | X=x_i, \theta, b]$

$(1-\hat{y}_i) \quad P[y=0 | \dots]$

Pb: learn θ, b

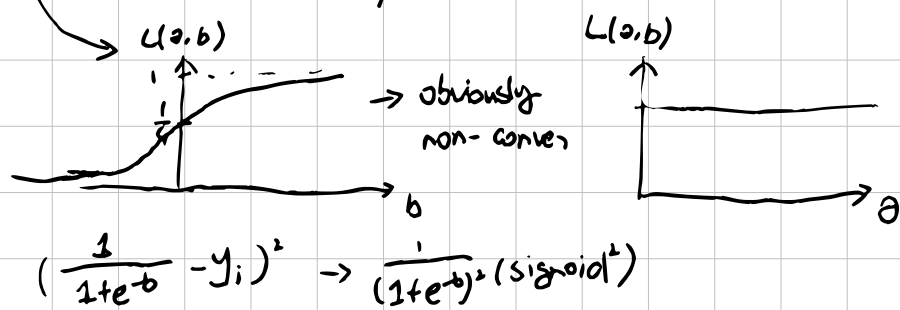
very naive idea: learn $\hat{\theta}, \hat{b} = \underset{\theta, b}{\operatorname{argmin}} \mathbb{1}[\hat{c}_i \neq y_i]$

identical to previous one (hard)

naive: $\underset{\theta, b}{\operatorname{argmin}} \sum_{i=1}^N (\hat{y}_i - y_i)^2$

the objective is continuous but not convex
 \Rightarrow hard to solve

Ex. show with a dataset with one example (0,0) that the objective is non-convex.



$(\frac{1}{1+e^b} - y_i)^2 \rightarrow \frac{1}{(1+e^b)^2} (\text{sigmoid}^4)$

The probabilistic interpretation:

The likelihood of model θ, b :

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

Neg log likelihood:

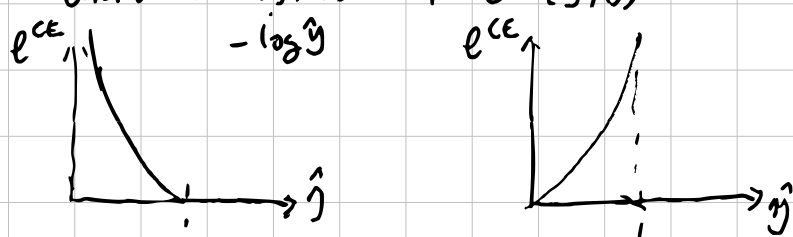
$\text{NLL} = -\sum_{i: y_i=1} \log \hat{y}_i - \sum_{i: y_i=0} \log(1-\hat{y}_i)$

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

logistic regression pb. cross entropy ℓ^{CE}

$\hat{\theta}, \hat{b} = \underset{\theta, b}{\operatorname{argmin}} \sum_{i=1}^N \ell^{CE}(\hat{y}_i, y_i)$

Ex. - draw $\ell^{CE}(\hat{y}, 1)$ and $\ell^{CE}(\hat{y}, 0)$



Multiclass Settings:

$\mathcal{Y} = \{1, \dots, k\}$

for each $k \in \mathcal{Y}$, $\theta_k \in \mathbb{R}^d$, $b_k \in \mathbb{R}$

score:

$(y_{i1}, \dots, y_{ik}) = (\theta_1^T x_i + b_1, \dots, \theta_k^T x_i + b_k)$

- predicted class: $\hat{c}_i = \underset{k \in \{1, \dots, k\}}{\operatorname{argmax}} \hat{y}_{ik}$

- softmax $\begin{pmatrix} e^{t_1} \\ \vdots \\ e^{t_k} \end{pmatrix} = \frac{1}{\sum_{i=1}^k e^{t_i}} \begin{pmatrix} e^{t_1} \\ \vdots \\ e^{t_k} \end{pmatrix}$

Estimated class pro: $\hat{y}_i = (\hat{y}_{i,1}, \dots, \hat{y}_{i,k}) = \text{softmax}(\dots)$

- multiclass CE loss:

$\ell^{CE}(\hat{y}_i, y_i) = \sum_{k=1}^k -\mathbb{1}[y_i=k] \log \hat{y}_{i,k}$

The logistic regression: the 2 views:

view 1: class probabilistic estimation loss

$\underset{\theta, b}{\operatorname{argmin}} \sum \ell^{CE}(\hat{y}_i, y_i)$

view 2: score loss

$\underset{\theta, b}{\operatorname{argmin}} \sum \underbrace{\ell^{\text{score}}(\hat{y}_i, y_i)}_{= \log(1+e^{-y_i \hat{y}_i})}$

ex. develop view 1 $\xrightarrow{\text{identical}}$ view 2

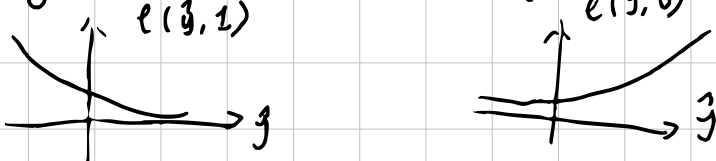
$\underset{\theta, b}{\operatorname{argmin}} \sum \ell^{CE}(\hat{y}_i, y_i)$

$= \dots = \ell^{CE}(\frac{1}{1+e^{-y_i \hat{y}_i}}, y_i)$

$= \dots = \sum -y_i \log(\frac{1}{1+e^{-y_i \hat{y}_i}}) - (1-y_i) \log(1 - \frac{1}{1+e^{-y_i \hat{y}_i}})$

$= \dots = \sum \log(1+e^{-y_i \hat{y}_i})$

$\log(1+e^{-y \hat{y}})$ is convex in \hat{y}



the overall objective

$\sum_{i=1}^N \ell^{\text{log}}(\hat{y}_i, y_i) = \sum_{i=1}^N \ell(\theta^T x_i + b | y_i)$ is convex in θ and b

Notation on this page: \hat{y} is score $\sigma(\hat{y})$

Neural nets : 2 views

1. outputs CPE

(cross prob estimation, CPE loss $\ell(\hat{y}, y)$)

2. score \hat{y} , scoring loss (multiclass logistic loss)
 $\ell(\hat{y}, y)$

Proper CPE loss $\hat{y} \in [0, 1]$ the risk of \hat{y}
 $R^\ell(\hat{y}(\cdot)) = \mathbb{E}_{x, y \sim p} [\ell(\hat{y}(x), y)]$ with CPE loss ℓ

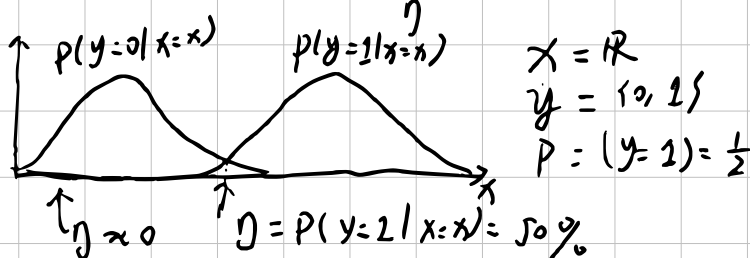
conditional risk: for $\eta \in [0, 1]$ and $\hat{y} \in [0, 1]$
 $C(\hat{y}, \eta) = \mathbb{E}_{y \sim \text{Ber}(\eta)} [\ell(\hat{y}, y)] = \eta \ell(\hat{y}, 1) + (1-\eta) \ell(\hat{y}, 0)$

note: $R^\ell(\hat{y}(\cdot)) = \mathbb{E}_x C(\hat{y}(x), \eta(x))$ where $\eta(x) = P(y=1|x)$

A good CPE loss will be a proper loss:

def: a CPE loss is proper iff

$\forall \eta \in [0, 1] \quad \eta \in \arg \min_{\hat{y}} C(\hat{y}, \eta)$



$$\begin{aligned} R^\ell(\hat{y}(\cdot)) &= \mathbb{E}_{x, y} [\ell(\hat{y}(x), y)] \\ &= \mathbb{E}_x [\mathbb{E}_y [\ell(\hat{y}(x), y) | x]] \\ &= \mathbb{E}_x [\eta(x) \ell(\hat{y}(x), 1) + (1-\eta(x)) \ell(\hat{y}(x), 0)] \\ &= \mathbb{E}_x [C(\hat{y}(x), \eta(x))] \end{aligned}$$

$\eta(x) = P(y=1|x)$

Ex. show that cross entropy is proper

$$\begin{aligned} C(\hat{y}, \eta) &= \eta \ell(\hat{y}, 1) + (1-\eta) \ell(\hat{y}, 0) \\ &= \eta (-\log \hat{y}) + (1-\eta) (-\log(1-\hat{y})) \end{aligned}$$

to show $\eta \in \arg \min_{\hat{y}} C(\hat{y}, \eta)$

$$\begin{aligned} \frac{\partial C}{\partial \hat{y}} &= -\frac{\eta}{\hat{y}} + \frac{1-\eta}{1-\hat{y}} \\ &= \frac{-\eta(1-\hat{y}) + \hat{y}(1-\eta)}{\hat{y}(1-\hat{y})} = 0 \end{aligned}$$

$$-\eta + \eta\hat{y} + \hat{y} - \hat{y}\eta = 0 \Rightarrow \hat{y} = \eta$$

$$\Rightarrow \eta \in \arg \min_{\hat{y}} C(\hat{y}, \eta)$$

ℓ^{CE} is (strictly) proper

$$\inf_{\hat{y}} R^\ell(\hat{y}) = \mathbb{E}_x \inf_{\hat{y}} C(\hat{y}, \eta(x)) = \mathbb{E}_x C(\eta(x), \eta(x))$$

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

$-\ell(\hat{y}, y)$ in $(0, 1)$, could also be interpreted as a parameter of a Bernoulli distri

ex. $\ell(10\%, 0)$

\downarrow Bern(10%) \downarrow Bern(0%)

Kullback - Leibler divergence: compare discrete distribution



two discrete distributions p and q on y

KL - divergence

$$KL(p||q) = \sum_{y \in Y} p(y) \log \frac{p(y)}{q(y)}$$

property:

$$KL(p||q) = 0 \quad \text{if } p=q$$

$$> 0 \quad \text{if } p \neq q$$

$$KL(p||q) \neq KL(q||p)$$

Compare $\text{Ber}(\hat{y})$ with $\text{Ber}(y)$

$$\begin{aligned} KL(\text{Ber}(y), \text{Ber}(\hat{y})) &= y \log \frac{y}{\hat{y}} + (1-y) \log \frac{1-y}{1-\hat{y}} \\ &= -y \log \hat{y} - (1-y) \log(1-\hat{y}) - \ell^{\text{CE}}(\hat{y}, y) \\ &\quad + y \log y + (1-y) \log(1-y) \end{aligned}$$

-Entropy(Ber(y)) = 0

Assume $0 \log 0 = 0$

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

Scoring loss

formulated: $y \in \{-1, 1\}$ $\ell(\hat{y}, y)$ logistic loss

$$\ell(\hat{y}, y) = \phi(\hat{y}y) \quad \phi - \text{loss or margin loss}$$

$$\ell^{0,1}(\hat{y}, y) = \mathbb{1}[\hat{y}y \leq 0] \quad \phi(\cdot) = \mathbb{1}[\cdot \leq 0]$$

$$\ell^{\text{logistic}}(\hat{y}, y) = \log(1 + e^{-\hat{y}y}) \quad \phi(z) = \log(1 + e^{-z})$$

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

The risk of scoring loss:

$$\eta(x) = P(y = 1 | x)$$

$\hat{y} \in \mathbb{R} \cup \{-\infty, \infty\}$ is a score and $\hat{y}(\cdot) : \mathcal{X} \rightarrow \mathbb{R} \cup \{-\infty, \infty\}$ is a scoring function.

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

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

$$= \eta \ell(\hat{y}, 1) + (1-\eta) \ell(\hat{y}, -1)$$

def: the risk for scoring func $\hat{y}(\cdot) : \mathbb{R}^d \rightarrow \mathbb{R}$ for loss ℓ is

$$R^\ell(\hat{y}(\cdot)) = \mathbb{E}_{x \sim p} [\ell(\hat{y}(x), y)] = \mathbb{E}_x [C(\hat{y}(x), \eta(x))]$$

Calibration of scoring losses:

what's "good" scoring loss?

def: A loss ℓ is *calibrational* if

$$\text{if } \eta \in [0, \frac{1}{2}] \text{ then } \int_{\hat{y} < 0} C(\hat{y}, \eta) < \int_{\hat{y} > 0} C(\hat{y}, \eta)$$

$$\text{if } \eta \in]\frac{1}{2}, 1] \text{ then } \int_{\hat{y} > 0} C(\hat{y}, \eta) < \int_{\hat{y} < 0} C(\hat{y}, \eta)$$

Thm: if ℓ is calibrated then

if $\hat{y}(\cdot)$ is the measurable function

minimizing R^ℓ ,

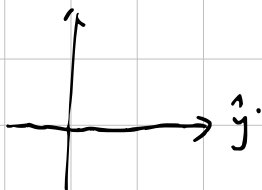
then it also minimizes $R^{0,1}$

Intuition: a measurable function learnt to minimize calibrational loss will also minimize 0/1 loss

$$C(\hat{y}, \eta) = \eta \ell(\hat{y}, 1) + (1-\eta) \ell(\hat{y}, -1)$$

Ex. $\ell^{\text{hinge}}(\hat{y}, y) = \max(0, 1 - \hat{y}y)$

draw $C(\hat{y}, \frac{1}{2}) = \eta(1 - \hat{y}) + (1-\eta)(1 + \hat{y}) =$



Thm: Any convex ϕ -loss with $\phi'(0) < 0$ is well-calibrated.
 \Rightarrow all scoring functions we saw are ...