

Statistics and Data Science 365 / 565

# Data Mining and Machine Learning

February 15

Yale

# Outline

- Recap: Empirical risk minimization  
*Different motivation for logistic loss*
- { Geometric interpretations: least squares and margin  
*↳ Probability*
- Maximum likelihood      *mle*
- Notebook

## Recap

Minimize population risk

$$R(f) = \mathbb{E} \ell(f(x), y)$$

Can't do that. Do empirical risk minimization

coz don't know  $\mathbb{E} \ell(\cdot, \cdot)$

$$\hat{R}(f) = \frac{1}{n} \sum_{i=1}^n \ell(f(x_i), y_i)$$

when data identically distributed

$$\mathbb{E}[\hat{R}(f)] = R(f)$$

**Binary classification surrogate losses:** Boosting, logistic, hinge

**Regression losses:** squared, absolute, (there are more)

# Geometric interpretation

نحوی بزه  
برای

$\theta \in \mathbb{R}^2$

$X \in \mathbb{R}^{n \times 2}$

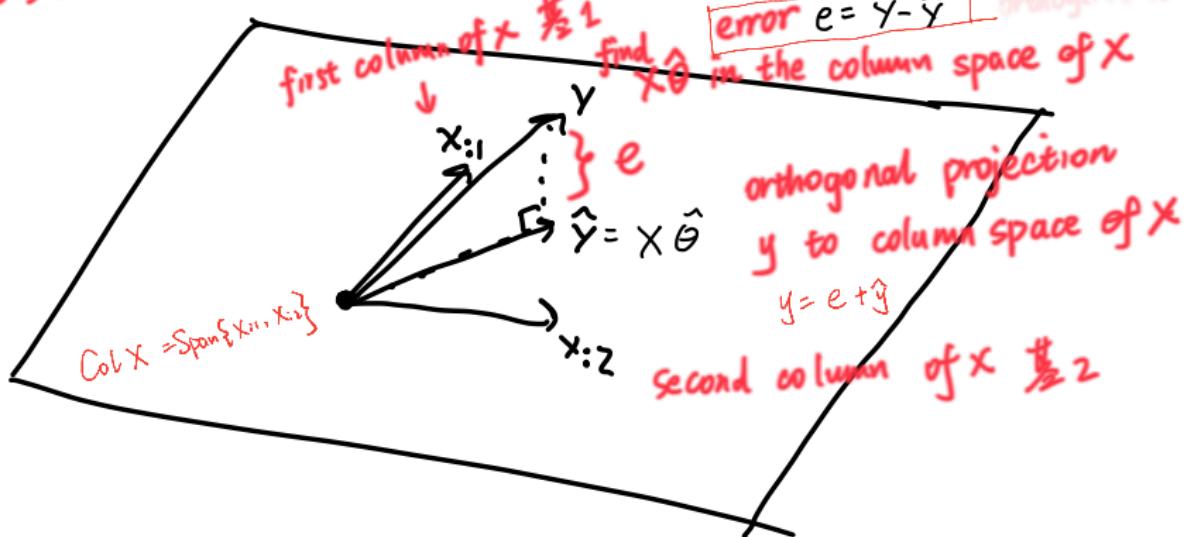
Linear Regression:  $y = X\theta + w$

OLS solution  $\hat{\theta}$  is to make sure  $e \perp \text{col } X$   
so  $e$  is shortest distance between  $y$  and  $\hat{y}$

proof in the 2

OLS:  $\min_{\theta} \|X\theta - y\|_2^2$

col  $X$



# Geometric interpretation: margin

$$= y f(x)$$

focus on classification

{ blue: +1  
red: -1

label

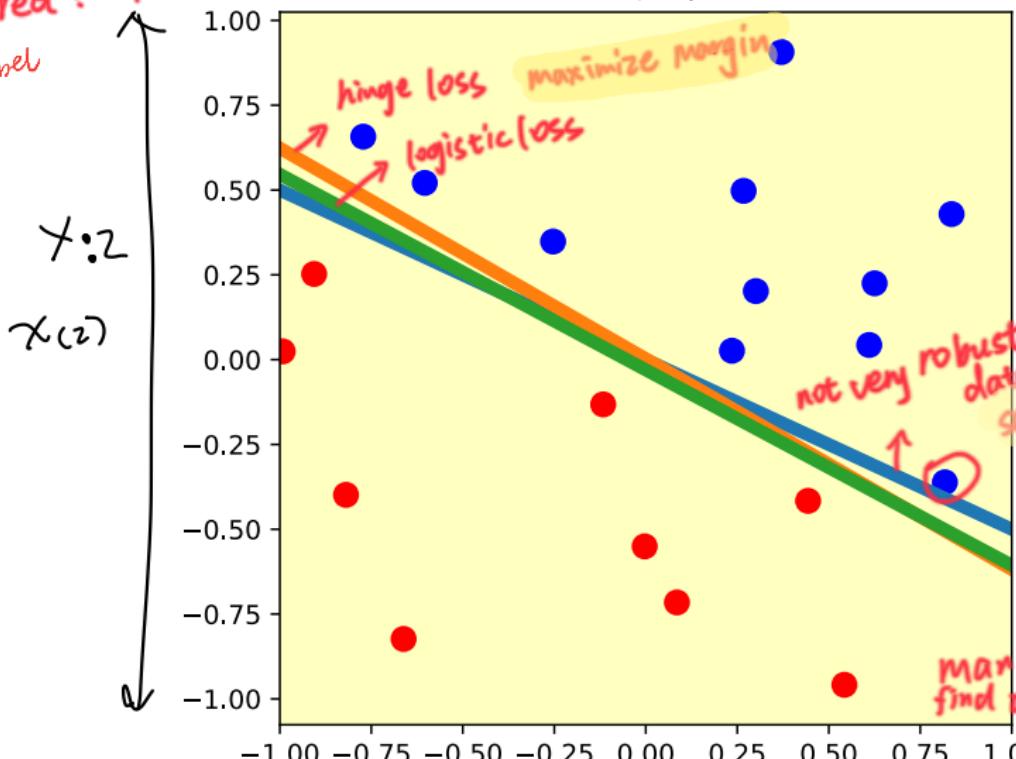
$x_{(1)}$

$x_{(2)}$

label  $y$

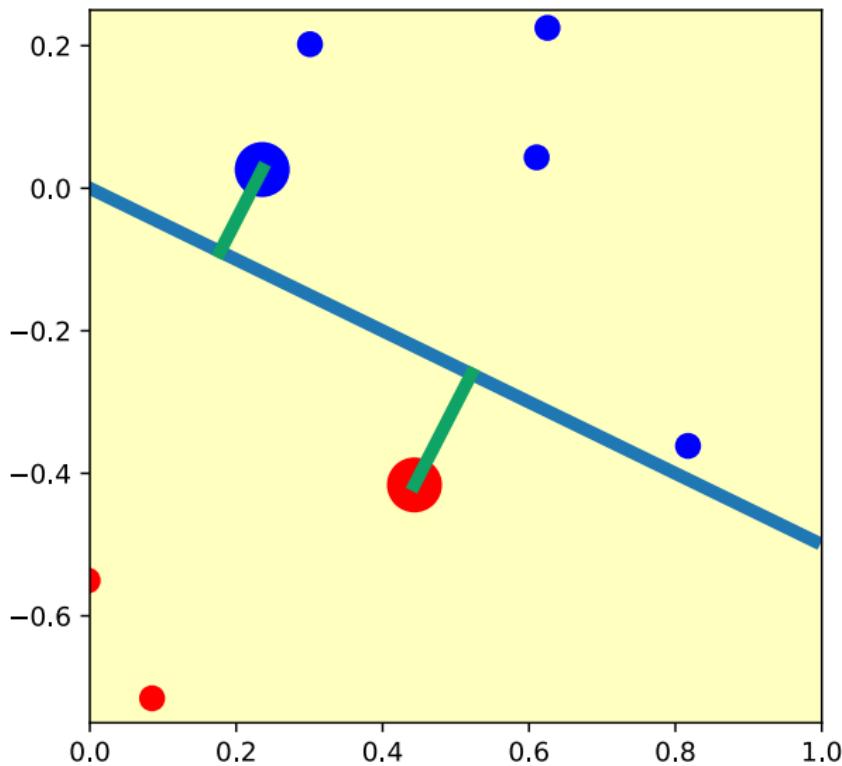
$\downarrow$  how "certain"  
 $x \in \mathbb{R}^{n \times 2}$

my function  
it should be  
think for  
one value or  
another value  
for binary  
classification



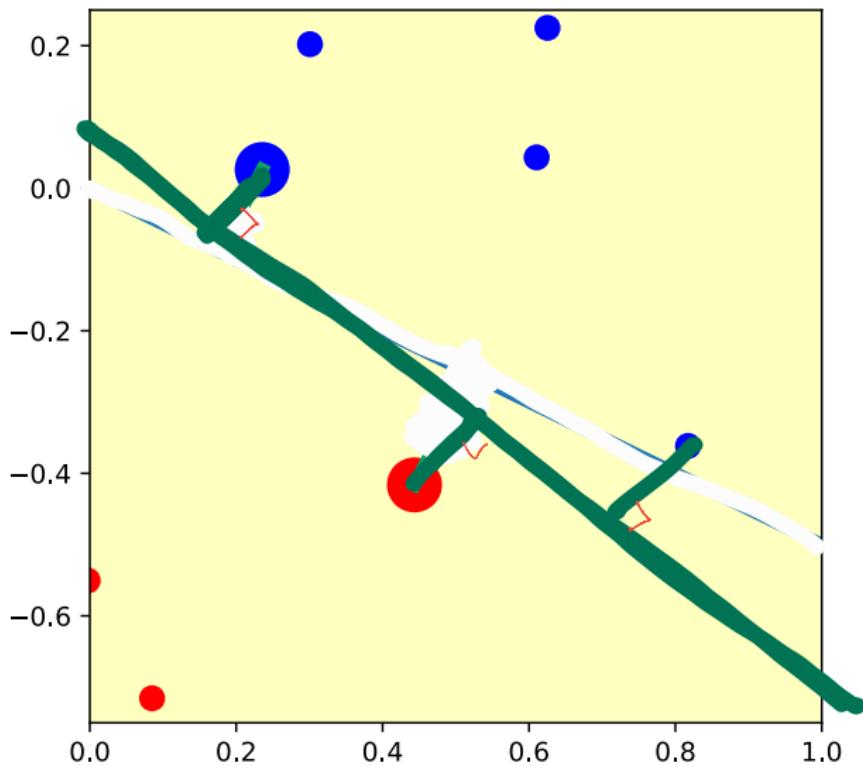
# Geometric interpretation: margin

*zoom in  
green line: margin distance of point between decision boundary*



# Geometric interpretation: margin

a more preferable decision boundary



# Geometric interpretation: margin

if we don't have  $\| \theta \|_2$ , take  $\theta$  and  $\theta_0$  very large, margin  $\rightarrow +\infty$

make margin formula scale free

decision boundary is dictated by  $\theta$ , but don't change with size of  $\theta$   
does change is "certainty", as  $\theta, \theta_0$  gets larger means we are more certain

Take  $\mathcal{Y} = \{-1, +1\}$ ,  $x_i \in \mathbb{R}^d$ ,  $\theta \in \mathbb{R}^d$ ,  $\theta_0 \in \mathbb{R}$  talk calibration later

$\theta_{\text{null}}(\theta_0)$ : intercept of  $\theta$  is real number

↓ inner product of  $x_i, \theta$

$$\text{margin}(\theta) \equiv \min_i \frac{y_i(\langle x_i, \theta \rangle + \theta_0)}{\|\theta\|_2}$$

Hw 3 proof for perfect separable data  
 $y_i(\langle x_i, \theta \rangle + \theta_0)$  is the smallest distance between any point and decision line  $\nparallel \|\theta\|_2$ ; renormalization



# Probability interpretation: classification

to be easier

$$\mathcal{Y} = \{0, 1\}$$

$$R(f) = \mathbb{E} \mathbb{1}(f(x) \neq y)$$

$$f^* = \arg \min_f R(f).$$

$$f^*(x) = \mathbb{1}\left(\mathbb{P}(y=1|x) > 0.5\right)$$

Rule:  $\mathbb{E}[\mathbb{1}(\dots)] = \mathbb{P}(\dots)$  expected value of indicator = Probability of ...

$$\mathbb{E} \mathbb{1}(f(x) \neq y) = \mathbb{P}(f(x) \neq y)$$

$$\begin{aligned} & \text{law of iterated expectation} \\ & \mathbb{E}[Y] = \mathbb{E}[\mathbb{E}[Y|X]] \\ & \text{given } + \quad \downarrow \\ & = \mathbb{E}\left[\mathbb{E}[\mathbb{1}(f(x) \neq y)|x]\right] \quad \text{2 } y \in \{0, 1\} \\ & = \mathbb{E}\left[\left[\mathbb{1}(f(x) \neq 1)\mathbb{P}(y=1|x) + [\mathbb{1}(f(x) \neq 0)\mathbb{P}(y=0|x)\right]\right] \quad \downarrow \quad \text{penalty ①} \quad \text{①} + \text{②} = 1 \end{aligned}$$

Focus on the inside. of  $E$  coz when inside is smallest  $E[\cdot]$  is also smallest

无论  $f(x)$  是 1 或 0，总要 pay 2 个 penalty 中的一个，因为有 noise，可能分类错

try to pick the smallest penalty of 2 penalties e.g. 0.4, 0.6) pick 0.4

$R(f) = \begin{cases} \mathbb{E} \mathbb{P}(y=1|x) & f(x)=0 \\ \mathbb{E} \mathbb{P}(y=0|x) & f(x)=1 \end{cases}$

# Probability interpretation: classification

$$f^* = \arg \min_f R(F).$$

$$f^*(x) = \mathbb{1}\left(\mathbb{P}(y = 1|x) > 0.5\right)$$

*Hamming loss*

$$\mathbb{E} \mathbb{1}(f(x) \neq y) = \mathbb{P}(f(x) \neq y)$$

$$= \mathbb{E}\left[\mathbb{E}[\mathbb{1}(f(x) \neq y)|x]\right]$$

$$= \mathbb{E}\left[[\mathbb{1}(f(x) \neq 1)]\mathbb{P}(y = 1|x) + [\mathbb{1}(f(x) \neq 0)]\mathbb{P}(y = 0|x)\right]$$

Focus on the inside.

In general,

But

Define  $\eta(x) = \mathbb{P}(y = 1|x)$ . Don't have in general.

So

Create a model:

assume true model is this  $\theta^*$ ; true  $\theta$



$$\eta(x) = \frac{\exp(\theta^{*T}x)}{1 + \exp(\theta^{*T}x)}$$

$$\theta^{*T}x = \langle \theta^*, x \rangle$$

inner product

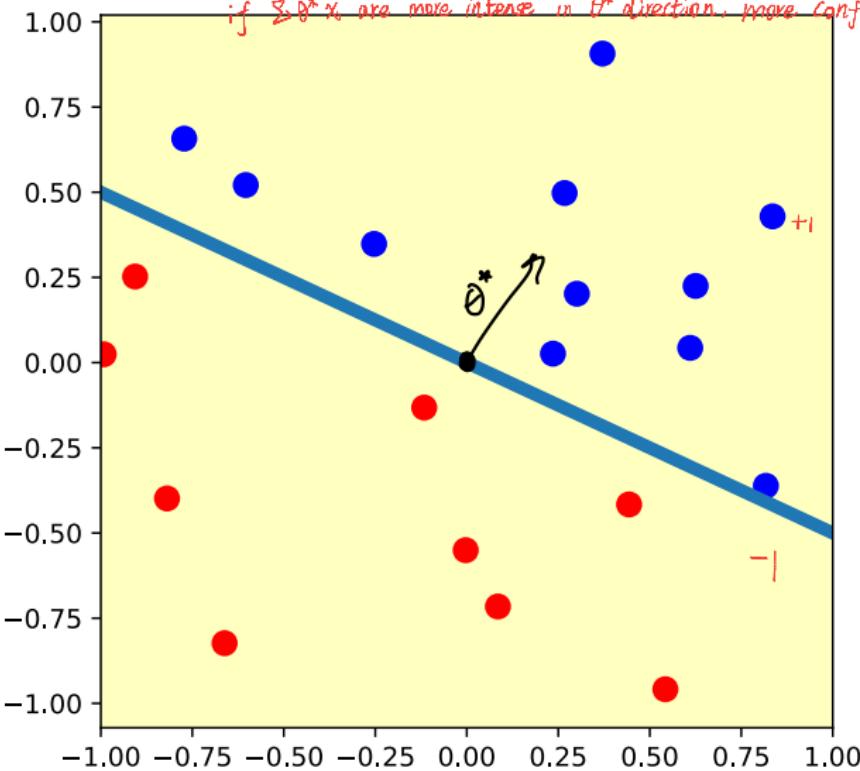
Assuming ignore intercept  $\theta_0$  or  $\theta_0$  automatically built in

# Logistic

$$\text{true } \theta = \theta^* = (1, 2)$$

If data points are more aligned with direction of  $\theta^*$ , more confident label is positive

if  $\sum \theta^* x_i$  are more intense in  $\theta^*$  direction, more confident label is positive

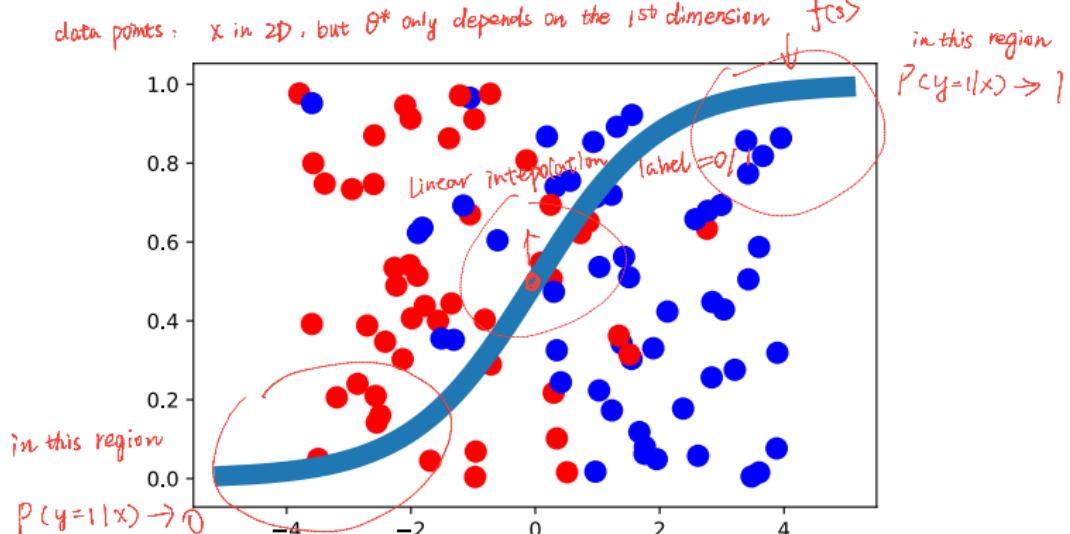


# Logistic

Sigmoid function  $f(s) = \frac{1}{1 + \exp(-s)} = \frac{\exp(s)}{1 + \exp(s)}$   $= P(y=1|x)$   
probability that label = +1

meaning of noise

data points :  $x$  in 2D, but  $\theta^*$  only depends on the 1st dimension  $f(x)$

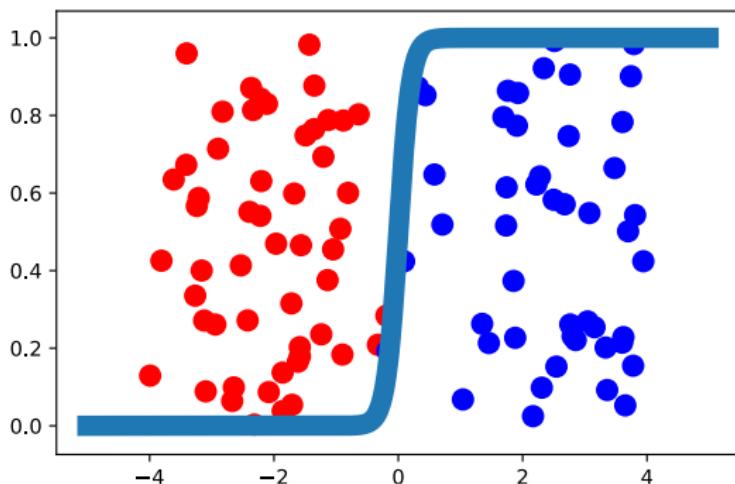


# Logistic

f(s)

Scaling will decrease noise

more certain where label actually are



## NLLLOSS

**CLASS** `torch.nn.NLLLoss(weight: Optional[torch.Tensor] = None, size_average=None, ignore_index: int = -100, reduce=None, reduction: str = 'mean')`

[\[SOURCE\]](#)

The negative log likelihood loss. It is useful to train a classification problem with C classes.

If provided, the optional argument `weight` should be a 1D Tensor assigning weight to each of the classes. This is particularly useful when you have an unbalanced training set.

The `input` given through a forward call is expected to contain log-probabilities of each class. `input` has to be a Tensor of size either  $(minibatch, C)$  or  $(minibatch, C, d_1, d_2, \dots, d_K)$  with  $K \geq 1$  for the K-dimensional case (described later).

Obtaining log-probabilities in a neural network is easily achieved by adding a `LogSoftmax` layer in the last layer of your network. You may use `CrossEntropyLoss` instead, if you prefer not to add an extra layer.

The target that this loss expects should be a class index in the range  $[0, C - 1]$  where  $C = \text{number of classes}$ ; if `ignore_index` is specified, this loss also accepts this class index (this index may not necessarily be in the class range).

The unreduced (i.e. with `reduction` set to `'none'`) loss can be described as:



$$\ell(x, y) = L = \{l_1, \dots, l_N\}^\top, \quad l_n = -w_{y_n} x_{n, y_n}, \quad w_c = \text{weight}[c] \cdot 1\{c \neq \text{ignore\_index}\}$$

where  $x$  is the input,  $y$  is the target,  $w$  is the weight, and  $N$  is the batch size. If `reduction` is not `'none'` (default `'mean'`), then

# Maximum likelihood/minimum negative log-likelihood

Assumption of Logistic Regression

①  $X$  are i.i.d

Need to estimate  $\theta^*$  from (assume iid) data  $\{(x_i, y_i)\}$

Candidate parameter  $\hat{\beta}$       ②  $P(Y=1|X)$  for logistic  $f = f(s)$

Example  $i$  we have guess  $p_i = f(x_i^T \hat{\beta})$ . Likelihood of guess

probability Mass Function of Bernoulli Distribution  
i.i.d product

{pdf continuous data  
pmf discrete data}

assume  $\beta$  is independent of  $x_i$

$$L(\beta) = \prod_{i=1}^n p_i^{y_i} (1 - p_i)^{1-y_i} = P_{\beta}(\{x_i, y_i\} | \hat{\beta})$$

$$= \prod_{i=1}^n P_{\beta}(y_i | x_i)$$

$$= P(\{x_i, y_i\} | \{x_i\}, \beta)$$

$$= P(\{x_i, y_i\} | \{x_i, \beta\})$$

$$= \left\{ \prod_{i=1}^n p_i^{y_i} \right\}_{\substack{p_i = 1 \\ p_i = 0}}^{(y_i=1) \\ (y_i=0)}$$

Negative log-likelihood

i.e. minimize NLL

$$\ell(\beta) = \sum_{i=1}^n \underbrace{-y_i \log(p_i) - (1 - y_i) \log(1 - p_i)}_{\text{This is the loss for example } i}$$

Let

$$\hat{\theta} = \arg \min_{\beta} \ell(\beta) \quad (\text{MLE})$$

$\hat{\theta}$  is the maximum likelihood estimate for  $\theta^*$  (this <sup>true</sup> model)

Max likelihood/min negative log-likelihood fits into loss/risk framework

$$\frac{\exp(\hat{\theta}^T x)}{1 + \exp(\hat{\theta}^T x)}$$

Not necessarily a good estimate for  $\mathbb{P}(y = 1|x)$  i.e. not good

Calibration  $\rightarrow$  how good is my estimate *ansatzing the true probability*

Good classification performance

# Logistic Regression

i.e. Sigmoid function

$$\text{In logistic model } P_{\beta}(y=1|x) = \frac{\exp(x^T \beta)}{1 + \exp(x^T \beta)} \Rightarrow P_{\beta}(y=0) = \frac{1}{1 + \exp(x^T \beta)}$$

$L(\beta)$  可似簡成

$$P_{\beta}(y|x) = \frac{\exp(yx^T \beta)}{1 + \exp(x^T \beta)} =$$

(Next time)  $P(y|x) = \frac{\exp(yg(x))}{1 + \exp(g(x))}$

Likelihood

$$\prod_{i=1}^n \frac{\exp(y_i x_i^T \beta)}{1 + \exp(x_i^T \beta)} = \prod_{i=1}^n \left( \frac{\exp(x_i^T \beta)}{1 + \exp(x_i^T \beta)} \right)^{y_i} \times \left( \frac{1}{1 + \exp(x_i^T \beta)} \right)^{1-y_i}$$

Log-likelihood

loss for label  $y \in \{0, 1\}$

$$\sum_{i=1}^n y_i x_i^T \beta - \log(1 + \exp(x_i^T \beta))$$

Last lecture for  $\mathcal{Y} = \{-1, +1\}$  logistic loss

$$\ell(\beta^T x, y) = \log(1 + \exp(-2y\beta^T x)) / \log(2)$$

same as log likelihood

See supplementary notes Logistic Regression.ipynb to verify they are (up to scaling) the same

# Conclusion

Logistic regression as maximum likelihood estimation

Given features  $x_i$ , the larger probability of get  $y_i$  ( $P(y_i|x_i)$ ), the more accurate our model

Geometric interpretation of least squares and margin

Next time:

generative modeling

non-linear decision boundaries

exponential families *Lec 7*

optimization