

Classification

data: (x, y) : y : discrete set
 $y = \{c_1, c_2, \dots, c_n\}$
 (Multi-class: $\{c_1, c_2, c_3, \dots\}$)
 $V = \{c_1, c_2\}$: Binary classification
 (no ordering on class)

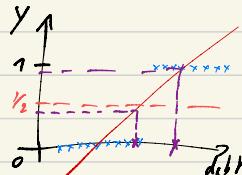
First example: use L-2 regression

Binary Classification Problem (0,1)

Given some data $S \rightarrow$ fit a regression fct \rightarrow if $f(x) > \frac{1}{2}$, then predict 1
 if $f(x) < \frac{1}{2}$, then predict 0

i.e.: S = credit card data (x, y)

x : current outstanding debt, $y = \begin{cases} 1 & \text{if the person default at the next payment} \\ 0 & \text{o.w.} \end{cases}$

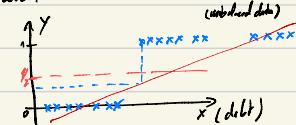


• However, classification is not only a special case of Lin.Reg.:

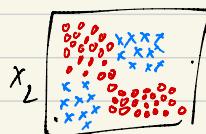
- The prediction is not continuous

- Sensitivity to unbalanced data:

(classification doesn't depend on L-2 loss, neither on accuracy)



Nearest Neighbor classification: (KNN...)



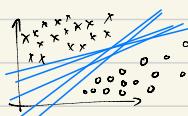
decision

$x: \text{dim } 2$

• The fundamental task of classification is to separate the space into decision regions
 \Rightarrow We can get some very complex boundaries

\Rightarrow This will only work on low dimensions. (Neighbors in high dimension?)

Linear decision boundaries:



• We assume that the data is linearly separable

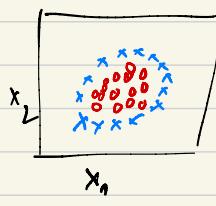


• minimal distance of the points to the hyperplane
 \rightarrow margin
 Idea: Maximize the margin

(SVM...)

• What if the data is not linearly separable?

Non linear classifier:



- We can do feature augmentation
 $(x, x^1, x^2, x^3, \dots)$
- Then do linear classification
→ Can be computationally exp
- We can use the Kernel method too

Bayes' prediction. (MAP rule)

$$\hat{y}(x) = \arg \max_{y \in Y} P(y|x)$$

$$(\hat{y}(x) = \arg \min_{g: Y \rightarrow Y} L_D(g))$$

$$P(\hat{y}|x) = \int p(x) P(\hat{y}(x)|x) dx$$

↳ prob of correct guess is the average form all x)

Thm: (Reg & classif)

let $y: X \rightarrow \mathbb{R}$ a reg. fct

$$g_y: X \rightarrow \{0, 1\}$$

$$x \mapsto \frac{1}{1 + e^{-(y(x))}}$$

$$L^c = \mathbb{E}_D [1_{g(x) \neq y}], \quad L^2 = \mathbb{E}_D [(y - y(x))^2]$$

$$\left\{ L^c(g_y) - L^c(g_y^*) \leq \sqrt{2(L^2(g_y) - L^2(g_y^*))} \right\}$$

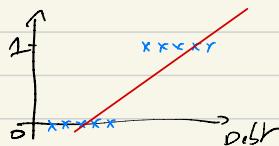
Logistic Regression

Binary classification.

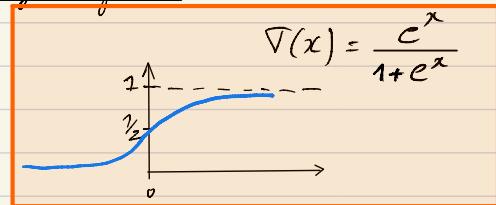
$(x, y) \sim D$
 x features $\in \mathcal{X}$
 y labels $\in \mathcal{Y} \subseteq \{0, 1\}$

Data $S = \{(x_n, y_n)\}_{n=1}^N$ iid
Given a new x we want to predict

Motivation:



Maps $(-\infty, +\infty)$ to $[0, 1]$
logistic function:



$$1 - \pi(x) = (1 + e^x)^{-1}$$

$$\pi'(x) = \frac{e^x}{(1 + e^x)^2} = \pi(x) \cdot (1 - \pi(x))$$

Logistic regression: $P(y=1|x) = \pi(x^T w)$
 $P(y=0|x) = 1 - \pi(x^T w)$ } Regression step

. logistic because of π , Regression because

- Quantize output: if $P(y=1|x) \geq 1/2$ predict 1
if $P(y=1|x) < 1/2$ predict 0

We assume that $x = \begin{pmatrix} 1 \\ x \end{pmatrix}$

(Scale w : $\|w\|$: factor in slope; w_0 : shift the transition (decision region) along w)
 $\|w\| \rightarrow \infty$, $\pi \rightarrow$ step function

Likelihood criterion:

$$S_{\text{Train}} = \{(x_i, y_i)\}_{i=1}^N \text{ iid}, \quad w^* = \arg \max_w P(y|x, w)$$

$$= \arg \max_w P(x|w) P(y|x, w)$$

Fact 1. $w \perp\!\!\!\perp x \Rightarrow P(x|w) = P(x)$

$$\begin{aligned} w^* &= \arg \max_w P(y|x, w) \\ &= \arg \max_w \prod_{i=1}^N P(y_i|x_i, w) \quad ((x_i, y_i)_{i=1}^N \text{ iid}) \\ &= \arg \max_w \prod_{i=1}^N \sigma(x_i^T w)^{y_i} (1 - \sigma(x_i^T w))^{1-y_i} \end{aligned}$$

- We prefer minimizing and maximization \Rightarrow minimize the negative log likelihood

$$\begin{aligned} w^* &= \arg \min_w - \sum_{i=1}^N y_i \log(\sigma(x_i^T w)) + (1-y_i) \log(1 - \sigma(x_i^T w)) \\ &= \arg \min_w \sum_{i=1}^N y_i \log \left(\frac{\sigma(x_i^T w)}{1 - \sigma(x_i^T w)} \right) - \log(1 - \sigma(x_i^T w)) \end{aligned}$$

$$\begin{aligned} (1 - \sigma(y)) &= \frac{1}{1 + e^y} \\ \left(\frac{\sigma(y)}{1 - \sigma(y)} \right) &= e^y \end{aligned}$$

$$w^* = \arg \min_w \sum_{i=1}^N -y_i \cdot x_i^T w + \log(1 + e^{x_i^T w})$$

$$\boxed{L(w) = \sum_{i=1}^N \log(1 + e^{x_i^T w}) - y_i \cdot x_i^T w}$$

($y \in \{0, 1\}$)

Gradient of $L(w)$:

$$\nabla L(w) = \sum_{i=1}^N \frac{e^{x_i^T w}}{1 + e^{x_i^T w}} \cdot x_i - y_i \cdot x_i = \sum_{i=1}^N x_i (\sigma(x_i^T w) - y_i)$$

$$X = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_N \end{bmatrix}_{N \times D}, \quad Y = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_N \end{bmatrix}_{N \times 1}$$

$$\boxed{\nabla L(w) = X^T [\sigma(Xw) - Y]}$$

- Some Gradient and LS back with $\nabla(\cdot)$ ($\nabla L(w) = X^T(Xw - Y)$)
- No closed form sol: $f_w \nabla L(w) = 0 \quad / \quad \underline{\text{L.R}}$
- Good news: L is convex

How to prove it?
 • Compute the Hessian
 • Look at open which preserves convexity (implies)

Convexity:
 - Product operator: positive sub: $(\cdot f + g)$ convex
 - composition:
 f convex and increasing $\Rightarrow f \circ g$ convex
 - $\lim_{t \rightarrow 0} \frac{f(tx) + f((1-t)x)}{t}$ is both convex and concave

$\Rightarrow y \mapsto \log(1 + e^y)$ convex?

$$\Rightarrow (\log(1 + e^y))^2 \cdot \frac{e^y}{1 + e^y} = \sigma(y)$$

$$\frac{\partial^2 \sigma(y)}{\partial y^2} = \sigma(y)(1 - \sigma(y))^{-2} = \sigma(y)(1 - \sigma(y))^{-1} > 0$$

$\Rightarrow \log(1 + e^y)$ is convex

$$\begin{aligned} \nabla^2 L(w) &= \sum_{i=1}^N x_i (\sigma(x_i^T w))' \\ &= \sum_{i=1}^N x_i \cdot x_i \cdot \sigma(x_i^T w) (1 - \sigma(x_i^T w))^{-1} \\ &= X^T S X \quad \text{where } S = \sum_{i=1}^N \sigma(x_i^T w) (1 - \sigma(x_i^T w))^{-1} \\ &\quad \text{since } S \geq 0, \quad \nabla^2 L(w) \gg 0 \quad (\text{L.R}) \end{aligned}$$

$\Rightarrow L$ is convex

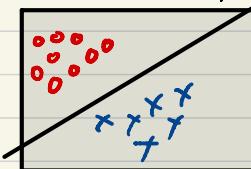
• How to minimize \mathcal{L} ? \Rightarrow Gradient descent:

$$w_k \in \mathbb{R}^D$$

$$w_{k+1} = w_k - \gamma_{k+1} \nabla \mathcal{L}(w_k)$$

$$\text{SGD: } w_{k+1} = w_k - \gamma_{k+1} \nabla \mathcal{L}_n(w_k)$$

\rightarrow What will happen if the data are linearly separable?



$$\mathcal{L}(w) = \sum_{i=1}^N \log(1 + e^{x_i^T w}) - y_i x_i^T w$$

$$\min \mathcal{L}(w)?$$

$$\inf \mathcal{L}(w) = 0 = \lim_{\|w\| \rightarrow \infty} \mathcal{L}(w)$$

The limit is not obtained by a finite w

\Rightarrow To avoid this we add L-2 regularization

$$\mathcal{L}(w) = \sum_{n=1}^N \log(1 + e^{x_n^T w}) - y_n x_n^T w + \lambda \|w\|_2^2$$

Newton alg. (2^{nd} order Alg.) $\left(\begin{array}{l} \text{for small } N, D \\ \hat{w}: \text{current estimate} \end{array}\right)$

$$\mathcal{L}(w) \sim \mathcal{L}(\hat{w}) + \nabla \mathcal{L}(\hat{w})(w - \hat{w}) + (w - \hat{w})^T \nabla^2 \mathcal{L}(\hat{w})(w - \hat{w}) + O(\|w - \hat{w}\|^3)$$

next pt: $w^* = \underset{w}{\operatorname{arg \min}} \text{quad. approx. } \mathcal{L}(w)$

$$\nabla \mathcal{L}(w) = \nabla \mathcal{L}(\hat{w}) + \nabla^2 \mathcal{L}(\hat{w})(w - \hat{w}) = 0$$

$$\Rightarrow w^* = \hat{w} - \nabla^2 \mathcal{L}(\hat{w})^{-1} \nabla \mathcal{L}(\hat{w})$$

$$\mathcal{L}(w) = \sum_{n=1}^N \log(1 + e^{x_n^T w}) - y_n x_n^T w \quad \bullet \text{ If our fat is quad. we will minimize it with one step of Newton Method}$$

$$= \sum_{n=1}^N \log \frac{(1 + e^{x_n^T w})}{e^{y_n x_n^T w}} = \sum_{n=1}^N \log \left(1 + e^{-y_n x_n^T w} \right) \quad \begin{aligned} y_n^1 &= 2y_n - 1 \\ y_n &= \frac{1}{2}(y_n + 1) \end{aligned}$$

if $y_n \in \{-1\} \hookrightarrow y_n \in \{0, 1\}$