

Intro to ML (Lecture-6)

- ⊕ k-nearest neighbor classifier
- ⊕ Perceptron

Structure | Problem setting

① Training samples $\mathcal{X} = \{(x_1, y_1), (x_2, y_2), \dots\}$
(Binary classification) $x_i \in \mathbb{R}^d, y_i \in \mathbb{R} \quad \{0, 1\} \quad \{-1, 1\}$

② Test point $\hat{x} \in \mathbb{R}^d$

Aim:

Find \hat{y}
find function $f: \mathbb{R}^d \rightarrow \mathbb{R}$

Any guess / any function is as good as any other
Underlying Assumptions about: Inductive Bias

Applications:

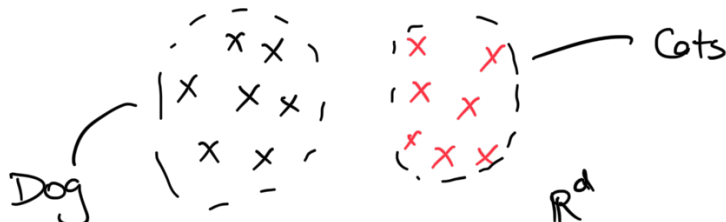
○ Image classification / Tagging

○ Spam classification

○ MP3 — Audio

① k-nearest neighbor

Bias / Assumption: Similar data points generally lie close to each other.



Algorithm / Pseudocode:

① For test point \hat{x} compute $\frac{d(\hat{x}, x_i)}{\text{distance b/w } \hat{x}, x_i}$ $\forall i \in \{1 \dots n\}$

$\| \quad \|_2$

$\| \quad \|_1$

② $i^* = \underset{i}{\operatorname{argmin}} d(\hat{x}, x_i)$

③ $\hat{y} = y_{i^*}$ ✓

1-nearest neighbor classifier

↓ extend

k-nearest neighbor

Running time each test point

n: training points

$O(nd)$

Problems?

① outliers

② label noise



③ $O(nd)$ large

$224 \times 224 \times 3$

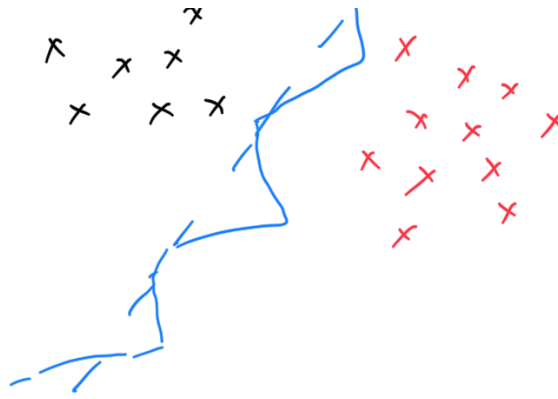
④ It ends up generating very "complex" decision boundaries

are not the best ones $\left\{ \begin{array}{l} \text{do not generalize well} \\ \text{- overfit} \end{array} \right\}$

"Occam's Razor"

x x ..





for i in range(N):
 $d(\tilde{x}, x_i) = \|x - x_i\|_2$
 if $d \leq \text{current-best}$
 $i\text{-star} = i$

$R=1$

$k = \text{anything}$

$$\hat{y} = y_i$$

⊗ Perceptron:

① Early attempts to solve classification

② 1950's Frank Rosenblatt

③

Aim/Goal:

Learn a linear separator b/w two classes

Setting:

$$\mathcal{X} = \{(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)\}$$

$$x_i \in \mathbb{R}^d \quad y_i \in \mathbb{R} \quad \&$$

Binary $y_i = \{-1, 1\}$

Aim

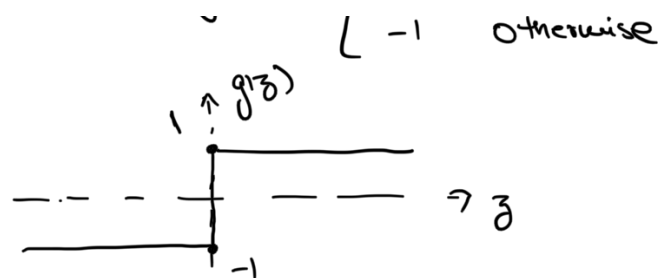
$$y = w^T x + b$$

$$= w^T x$$

Problem:

y is not really labels

$$g(z) = \begin{cases} 1 & \text{if } z \geq 0 \end{cases}$$



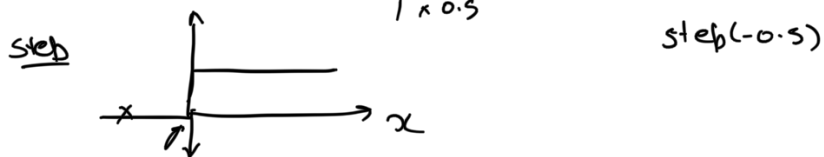
Loss: (Simplest possible)

$$L = 1 \quad \text{if misclassify}$$

$$L = 0 \quad \text{if classified correctly}$$

$$L = \text{step}(-y_i f(x_i)) \quad f(x_i) = w^T x$$

$$\frac{1 \times 0.5}{1 \times 0.5}$$



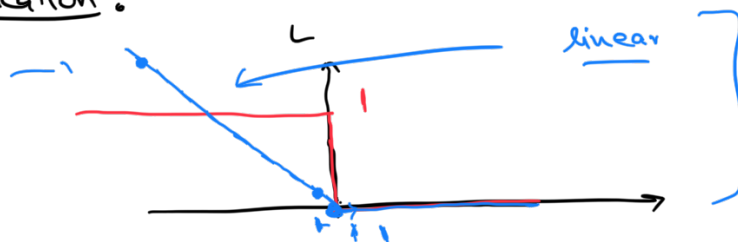
$y_i f(x_i)$

y_i	$f(x_i) = w^T x$	Loss
1	0.5	0
-1	-0.5	0
1	-0.5	1

Update: (gradient descent)

$$\frac{\partial L}{\partial w} = \begin{cases} 0 & \text{everywhere} \\ \bullet & \text{label flip} \end{cases}$$

Modification:



$$L(y, f(x, w)) = \begin{cases} -y_i f(x_i) & \text{correct} \\ 0 & \text{incorrectly classified} \end{cases}$$

$$\partial L \propto -y_i x_i \quad \text{incorrectly classified}$$

$$\overline{\Delta w} = \begin{cases} 0 & \text{otherwise} \end{cases}$$

$$w^{t+1} = w^t + \eta (y_i x_i)$$

Algorithm:

① $w_0 = 0$

② for epoch in range (10000):

for (x_i, y_i) in X :

if $f(x_i) \neq y_i$:

update

③ all data are correctly classified: terminate

Problems:

① previously correctly classified example can be misclassified.

② No. of iteration

Theorem: assumption:

① $\|x\| \leq 1$

② $b=0$

③ Data is perfectly separable

If there exists a solution then perceptron algorithm is guaranteed to find it.

[converge to 0 error in finite time]

Pf: $\rightarrow w^0 = 0$ $w^* = \text{optimum}$
 $\|x_i\| \leq 1$

margin $\rightarrow \gamma = \arg \min_i |w^T x_i|$ $\gamma > 0$

y_i	$w^T x$
1	0.1
1	0.5
1	10



$$\eta = 1$$

$$\begin{aligned} \textcircled{1} \quad \underline{(w^*)^T w^t} &= (w^*)^T (w^{t-1} + y_i x_i) \\ &= (w^*)^T w^{t-1} + y_i (w^*)^T x_i \\ &\geq \underline{(w^*)^T w^{t-1}} + \gamma \end{aligned}$$

↑
recursively

$$\boxed{(w^*)^T w^t \geq t\gamma} \leftarrow$$

$$\begin{aligned} \textcircled{2} \quad \underline{\|w^t\|^2} &= \|w^{t-1} + y_i x_i\|^2 \\ &= \|w^{t-1}\|^2 + \|y_i x_i\|^2 + \underbrace{2y_i w^{t-1} x_i}_{\substack{2y_i w^{t-1} x_i < 0 \\ \uparrow \quad \uparrow \\ \|w^{t-1}\| \|x_i\| = 1}} \\ &\leq \|w^{t-1}\|^2 + \underbrace{\|x_i\|^2}_{=1} \\ &\leq t \end{aligned}$$

$$\begin{aligned} \textcircled{3} \quad \text{Angle b/w optimal } w^* \text{ \& current solution} \\ 1 \geq \cos(w^*, w^t) &= \frac{(w^*)^T w^t}{\|w^*\| \|w^t\|} \end{aligned}$$

from ①

from ②

$$\begin{aligned} t\gamma &\leq \|w^*\| \|w^t\| \\ &\leq \|w^*\| \sqrt{t} \end{aligned}$$

$$\sqrt{t} \leq \frac{\|w^*\|}{\gamma}$$

$$t \leq \frac{1}{\gamma^2} \|w^*\|^2$$

$$\text{if } \|w^*\| = 1 \quad \boxed{t \leq 1/\gamma^2}$$

no. of update

Perceptron vs logistic:

① Differences:

→ output - probabilities

① if solution exist both of them will find it.

② logistic regression → motivated satisfies

