

Ex 1: $x_t \in \mathbb{R}^d$

Each x_t represents a 128×128 pixel image ($d = 128^2$).

want to use $\{x_t: t=1, 2, \dots, n\}$ (training images) to learn a classifier

classif. $\xrightarrow{\text{def}}$ $f: \mathbb{R}^d \rightarrow \{+1, -1\}$
 \searrow \swarrow
 d.o.g. cat

if x contains a cat, we hope $f(x) = -1$

x contains a dog, hope $f(x) = +1$.

use $\{x_t: t=1, 2, \dots, n\}$ to learn f .

→ Actually not because we still need label $\{(x_t, y_t) : t=1 \dots n\}$

To each training image, we have a label $y_t \in \{+1, -1\}$.

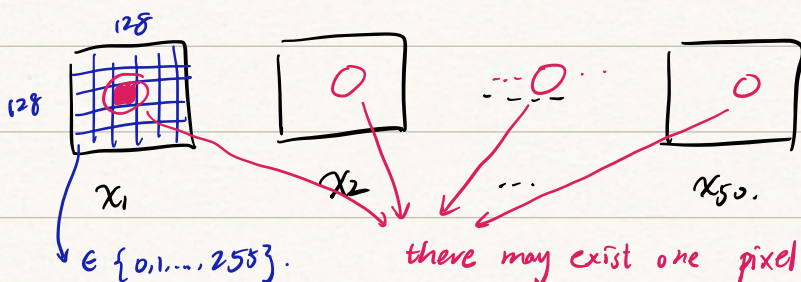
$$\hookrightarrow \mathcal{Q} = \{(x_t, y_t) : t=1, 2, \dots, n\}.$$

↳ use \mathbb{P} to learn a classifier $f \Rightarrow f: \mathbb{R}^d \rightarrow \{+1, -1\}$.

Ex of a rule for classifier f . $d=128^2=2^{14}$ pixels (dim of data)

$n = 50$ images

Each pixel value is an integer in $\{0, 1, \dots, 255\}$.



$x_{ti} \rightarrow i^{\text{th}}$ pixel in t^{th} image $x_t \rightarrow \begin{pmatrix} x_{t1} \\ \vdots \\ x_{td} \end{pmatrix}$

One classification Rule: (may not reasonable)

for x' , $f(x') = \{ y_t, \text{ if } x_{ti} = x'_i \text{ for some } t=1, \dots, n. \}$

1, else.

→ memorized the training sets.

Rmk: f classifies the training set perfectly!

⇒ Training loss = 0.

Q1: Is this classifier "GOOD"?

→ what is good?

main idea in machine learning is to generalize well.

→ Answer is not good.

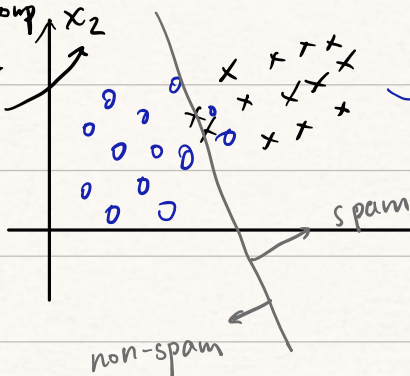
But this classifier f does not generalize well to new data.

Ex. (x_t, y_t) $x_t \in \mathbb{R}^2$ $y_t \in \{\text{spam}, \text{not spam}\}$

→ email example

$\begin{cases} x_{t1}: \# \text{ of words in Email } t \text{ are "loan"} \\ x_{t2}: \# \text{ of words in Email } t \text{ are "prize"} \end{cases}$

second comp
of feature
vector x_2



→ construct a good classifier

x_1 first comp of feature vector

Supervised Learning ⇒ x_t has label.

Notation

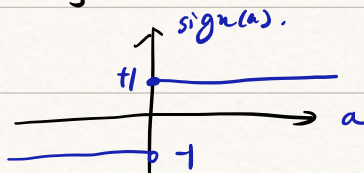
$\begin{cases} x_t: \text{training samples / feature vectors} \\ d: \# \text{ of features / dimension in } x_t \\ n: \# \text{ of training samples} \\ y_t: \text{labels} \end{cases} \begin{cases} \text{multi classification} \\ \text{binary classification} \end{cases} \quad y_t \in \{1, \dots, K\}$

Linear Classifiers & training error.

Def: A linear classifier $f: \mathbb{R}^d \rightarrow \{-1, +1\}$.
(linear classifier)

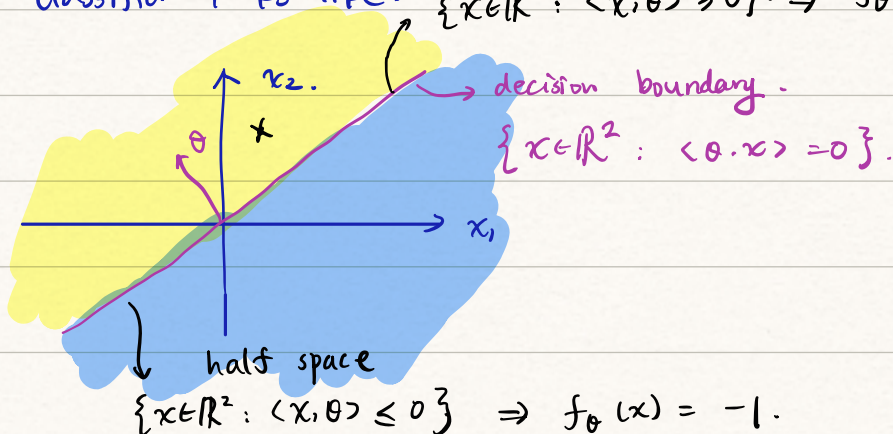
$$f_{\theta}(x) = \text{sign}\{\langle x, \theta \rangle\}.$$

$$\textcircled{1} \text{sign}(a) = \begin{cases} +1 & a \geq 0 \\ -1 & a < 0. \end{cases}$$



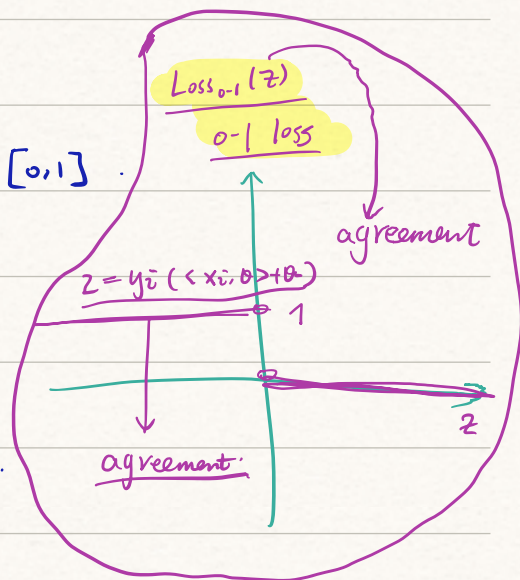
$$\textcircled{2} \langle x, \theta \rangle = \theta^T x = \sum_{i=1}^d \theta_i x_i$$

linear classifier looks like: half space.
 $\{x \in \mathbb{R}^2: \langle x, \theta \rangle \geq 0\} \Rightarrow f_{\theta}(x) = +1$



Def (training Error)

$$\hat{E}(\theta) = \frac{1}{n} \sum_{i=1}^n \text{Loss}(y_i, \underbrace{f_{\theta}(x_i)}_{\substack{\text{label.} \\ \hookrightarrow f(x_i; \theta)}}) \in [0, 1].$$



$$\text{Loss}(y, \hat{y}) = \mathbb{1}\{y \neq \hat{y}\} = \begin{cases} 1, & y \neq \hat{y} \\ 0, & y = \hat{y} \end{cases}$$

NOT Expectation
(Actually just a function)

\hookrightarrow 0-1 loss function

Indicative function

Many loss functions \Rightarrow $\begin{cases} \text{Entropy loss func} \\ \text{Hinge loss} \\ \text{Least square loss} \end{cases}$

★

Def: linearly separable - \mathcal{D}

A dataset $\mathcal{D} = \{(x_t, y_t) \mid t=1, 2, \dots, n\}$ is linearly separable

$$\Leftrightarrow \exists \theta \in \mathbb{R}^d \text{ st } \underline{y_t \langle \theta, x_t \rangle \geq 0} \text{ for } \forall t=1, 2, \dots, n$$

\hookrightarrow labeled correctly. (exists a hyperplane).

\hookrightarrow perceptron

Non-linearly separable \Rightarrow SVM.

Rmk: \mathcal{D} is linearly separable $\Leftrightarrow \exists \theta \in \mathbb{R}^d$ s.t. $\hat{E}(\theta) = 0$. (training error)
if $f_\theta(\cdot)$ is a linear classifier.

Rmk: for linearly classifier, i.e., $f_\theta(x) = \text{sign}(\langle \theta, x \rangle)$.

the decision boundary passes through origin

(we don't put the OFFSET!)

$\rightarrow \{x \in \mathbb{R}^d : \langle \theta, x \rangle = 0\} \rightarrow$ pass through the origin.

Perceptron Algorithm (Rosenblatt, 1960s).

linearly separable?

1). Arbitrarily initialize θ to be $\theta^{(0)}$

2). cycle through all n training samples

If $y_t \langle x_t, \theta^{(k)} \rangle \leq 0$ (i.e. there is a disagreement

between label y_t & prediction $\text{sgn}\{\langle x_t, \theta^{(k)} \rangle\}$.

Update $\theta^{(k+1)} = \theta^{(k)} + y_t x_t$

Some proof of Perceptron Algorithm

1. Considering \mathcal{D} contains only 1 sample, i.e., $\mathcal{D} = \{(x_1, y_1)\}$.

initialize $\theta^{(0)}$. ① if $y_1 \langle x_1, \theta^{(0)} \rangle > 0 \rightarrow$ over!

② if $y_1 \langle x_1, \theta^{(0)} \rangle \leq 0$

the simplest case

$|\mathcal{D}|=1$

then we have $\underline{\theta}^{(1)} = \underline{\theta}^{(0)} + y_1 \underline{x}_1$ (Update)

Check: $y_1 \langle \underline{x}_1, \underline{\theta}^{(1)} \rangle$

$$= y_1 \langle \underline{x}_1, \underline{\theta}^{(0)} + y_1 \underline{x}_1 \rangle$$

$$= y_1 \langle \underline{x}_1, \underline{\theta}^{(0)} \rangle + \boxed{y_1^2 \langle \underline{x}_1, \underline{x}_1 \rangle}$$

$$> y_1 \langle \underline{x}_1, \underline{\theta}^{(0)} \rangle + \|\underline{x}_1\|^2$$

To conclude, after updating - we have.

$$y_1 \langle \underline{x}_1, \underline{\theta}^{(w)} \rangle \geq y_1 \langle \underline{x}_1, \underline{\theta}^{(k)} \rangle$$

(increased by $\|\underline{x}_1\|^2$).

More generally, $y_1 \langle \underline{\theta}^{(k+1)}, \underline{x}_1 \rangle = y_1 \langle \underline{\theta}^{(k)}, \underline{x}_1 \rangle + \|\underline{x}_1\|^2$

→ After updating, $y_1 \langle \underline{\theta}^{(k+1)}, \underline{x}_1 \rangle$ increases by $\|\underline{x}_1\|^2 > 0$.

→ After $\left\lceil \frac{| \langle \underline{\theta}^{(w)}, \underline{x}_1 \rangle |}{\|\underline{x}_1\|^2} \right\rceil$ iterations, it will hold $y_1 \langle \underline{\theta}^{(k)}, \underline{x}_1 \rangle > 0$.

* In summary, if $|\mathcal{L}| = 1$, then perceptron always converges



after $\left\lceil \frac{| \langle \underline{\theta}^{(w)}, \underline{x}_1 \rangle |}{\|\underline{x}_1\|^2} \right\rceil$ # of steps.

this cannot be simply generalized to $|\mathcal{L}| = k$

because we cannot guarantee other training samples (\underline{x}_k, y_k) 's $y_k \langle \underline{\theta}^{(k)}, \underline{x}_k \rangle$ value. → it will be affected!

Perceptron Convergence Theorem (need this to prove)

Notation: (γ -linearly separable).

Def 1: Dataset $\mathcal{L} = \{(\underline{x}_i, y_i)\}_{i=1}^n$ is γ -linearly separable

if there exists $\underline{\theta}^* \in \mathbb{R}^d$, s.t. $y_t \langle \underline{\theta}^*, \underline{x}_t \rangle \geq \gamma > 0 \quad \forall t=1, 2, \dots, n$

$$\rightarrow \gamma = \min_{1 \leq t \leq n} y_t \langle \underline{\theta}^*, \underline{x}_t \rangle$$

Def 2: \mathcal{L} is R-bounded if $\|\underline{x}_t\| \leq R$ for every $t=1, 2, \dots, n$.

Rmk: Any dataset \mathcal{L} that is finite is R-bounded, (for some R).

→ take $R = \max_{1 \leq t \leq n} \|\underline{x}_t\|$ → while countable is not



the concept of limit \leftrightarrow sup.

↓
not bounded.

Thm: [Perceptron Convergence Thm]. (Novikoff 1962).

Initialize $\underline{\theta}^{(0)} = \underline{0} \in \mathbb{R}^d$. Assuming that \mathcal{D} is γ -linearly separable

& R -bounded. Then, the total # of updates of the Perceptron alg K satisfies

easy to satisfy

$$K \leq R^2 \frac{\|\underline{\theta}^*\|^2}{\gamma^2} \Rightarrow \text{give the upper bound of the \# of convergence.}$$

to get the upper bound of iteration

where $\underline{\theta}^*$ is the vector in the def. of γ -linearly separable.

Rmk: if γ is big \Leftrightarrow good separability

$\Rightarrow K$ tends to be small

\Rightarrow a few iterations.

★

Outline:

Classifier

$$f_{\theta}(\underline{x}) = \langle \underline{x}, \underline{\theta} \rangle$$

Linear classifier
(with no offset)

Update Parameter.



Algorithm



Training Error

$$\widehat{E}(\theta) = \frac{1}{n} \sum_{t=1}^n \text{Loss}(\bar{y}_t, f_{\theta}(\underline{x}_t))$$



0-1 Loss
hinge Loss
entropy Loss
(log loss).

$$\underline{L_{0-1}}(z) = (1-z)^+$$

$$L_{0-1}(z) = \begin{cases} z > 0, & 0 \\ z < 0, & 1 \end{cases}$$

$$z = \bar{y}_t \cdot \langle \theta, \underline{x}_t \rangle$$



agreement.