# SUPERVISED LEARNING

## CLASSIFICATION

$$\{ x_1, \ldots, x_n \} \qquad x_i \in \mathbb{R}^d$$

$$y_1, \ldots, y_n \qquad y_i \in \{0, 1\} \ / \ y_i \in \{+1, -1\}$$

Goal: $h : \mathbb{R}^d \rightarrow \{0, 1\}$

$$\mathbb{1}(z) := \begin{cases} 1 & \text{if } z \text{ is true} \\ 0 & \text{o/w} \end{cases}$$

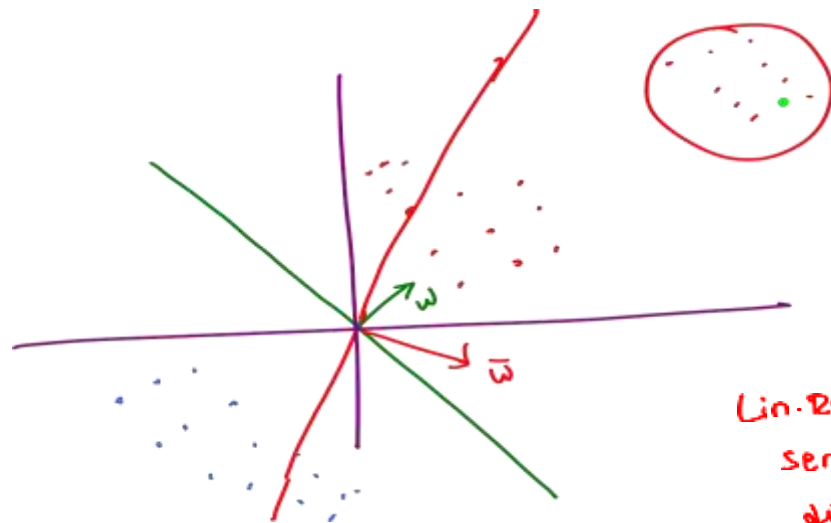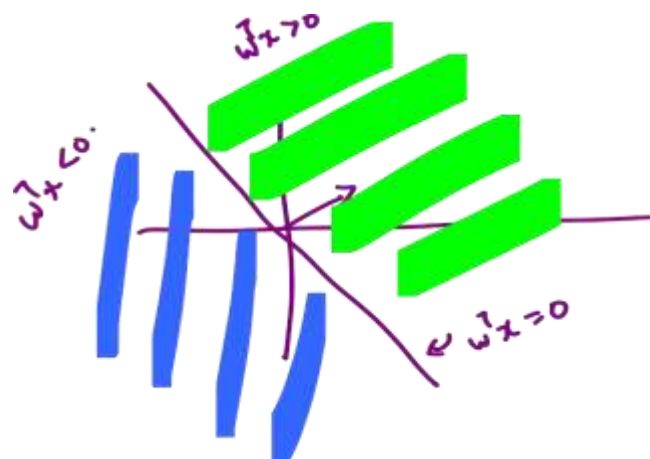ERROR: $\sum_{i=1}^{n} \mathbb{1}\left( h(x_i) \neq y_i \right)$

$$\mathcal{H}_{linear} = \left\{ h_w : \quad h_w(x) = \begin{cases} 1 & \text{if} \quad \vec{w}^T x \geq 0 \\ 0 & \text{o/w} \end{cases} \right.$$

$$\min_{h \in \mathcal{H}_{linear}} \sum_{i=1}^{n} \mathbb{1}(h(x_i) \neq y_i)$$

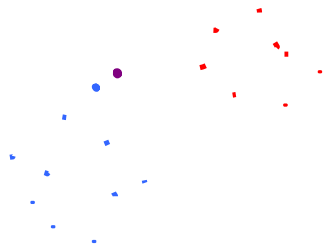$$sign(z) = \begin{cases} 1 & \text{if} \\ & z \geq 0 \\ 0 & \text{o/w.} \end{cases}$$

$$= \min_{w \in \mathbb{R}^d} \sum_{i=1}^{n} \mathbb{1}\left( sign(\vec{w}^T x_i) \neq y_i \right)$$

$$\longrightarrow \quad \underline{NP \ HARD \ Problem}$$

$\vec{w}x > 0$

$\vec{w}x < 0$.

$\leftarrow \vec{w}x = 0$
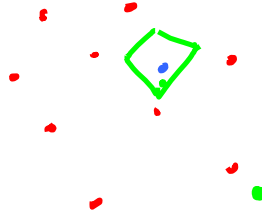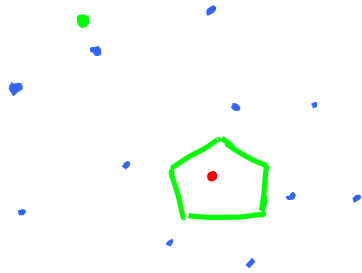
$\vec{w}$

$\vec{w}$

Lin·Reg is sensitive to distances.

# SIMPLEST POSSIBLE ALGORITHM

- Given $x_{test}$, find $x^*$ – the closest point to $x_{test}$ in the training set
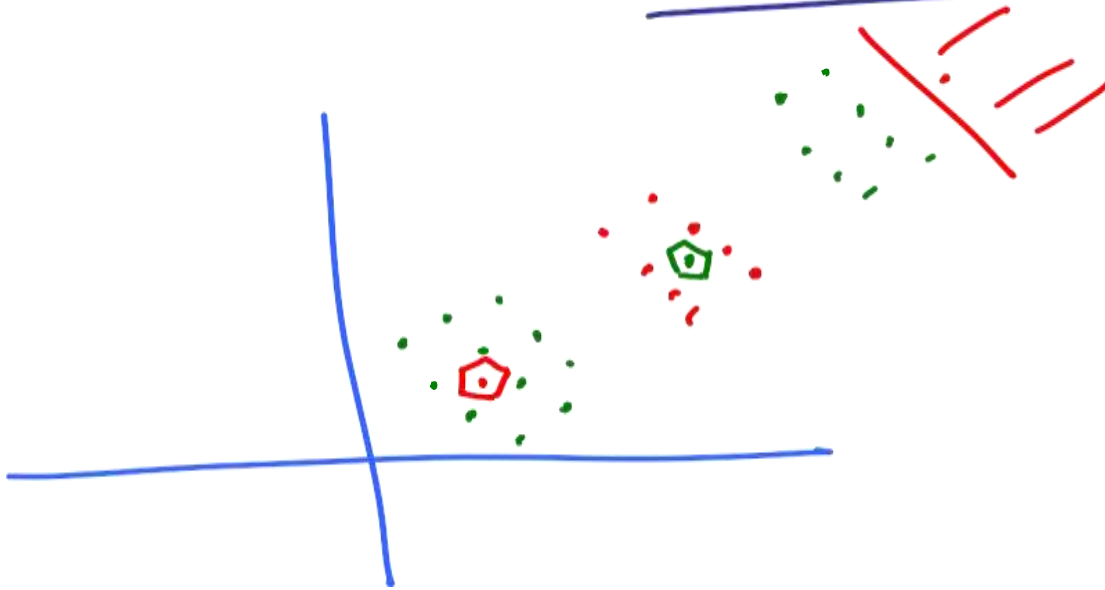
Predict $\hat{y}_{test} = y^*$.

Issue: outliers.
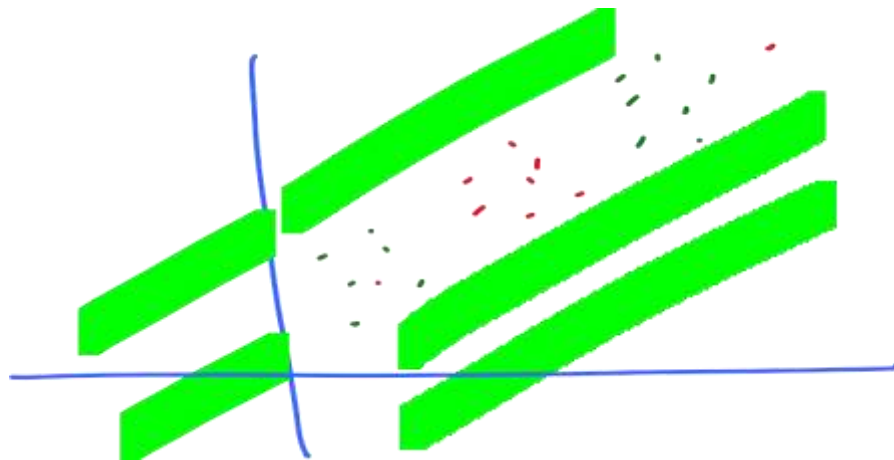
Fix: check more neighbours

K - Nearest Neighbour

→ Given $x_{test}$, look at $k$- nearest neighbours $x_1^*, x_2^*, \ldots, x_k^*$

→ $y_{test} = $ majority $(y_1^*, \ldots, y_k^*)$
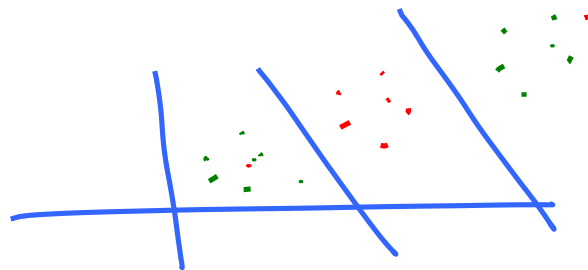
# DECISION BOUNDARY .

K=1

$k = n$

Choosing $k$

Cross validate.

$k^*$

# ISSUES with K-NN

- PREDICTION IS COMPUTATIONALLY EXPENSIVE.

- NO MODEL is learnt. Cannot throw away data after "learning"

# DECISION TREES

INPUT : Dataset $\{ (x_1, y_1) \cdots , (x_n, y_n) \}$

OUTPUT : DECISION TREE

## Decision tree

QUESTION

```
    Y        N
   1         0
```

PREDICTION: Given $X_{test}$, traverse through the tree to reach a leaf node. Predict $y_{test}$ = answer in leaf node

QUESTION:

A question is a (feature, value) pair.

Eg: height $\leq$ 180 cm ?

- How to measure "goodness" of a Question?

$$D \qquad \{ (x_1, y_1) \cdots \cdots (x_n, y_n) \}$$

Y

N

$$D_{yes} \quad \{ (x_5, y_5), (x_8, y_8) \cdots \}$$

$$\{ (x_1, y_1), (x_2, y_2), \cdots \} \quad D_{No}$$

$\{ y_1, y_2, \ldots, y_n \}$  $y_i \in \{0, 1\}$.

MEASURE of Impurity.

ENTROPY



entropy($p$)

$= -\left( p \log p + (1-p) \log 1-p \right)$

$[\log(0) = 0]$

0    0.25    0.5    0.75    1

$p \rightarrow$ Fraction of $1^s$.

$$\text{Information gain} \begin{pmatrix} \text{feature,} \\ \text{value} \end{pmatrix} = \text{ENTROPY} \left( D \right) - \left( \frac{\gamma \; \text{ENTROPY} \left( D_{yes} \right) + (1-\gamma) \; \text{ENTROPY} \left( D_{no} \right)}{2} \right)$$

$$\gamma = \frac{|D_{yes}|}{|D|} \qquad\qquad (1-\gamma) = \frac{|D_{no}|}{|D|}$$

## ALGORITHM

→ Discretize each feature in $[min, max]$ range.

→ Pick Question $(f_k \leq \theta)$ that has highest information gain

→ Repeat for $D_{yes}$ & $D_{no}$.

$F_t < \theta_t$ ?

$F_\lambda < \theta_\lambda$

$F_Y < \theta_Y$

Y    N

Y    N      Y    N

Points

- Depth is a hyper-parameter.

- Can also stop growing if node is "sufficiently" smooth.

# Decision Boundary

# Formal Treatment of the Learning Problem
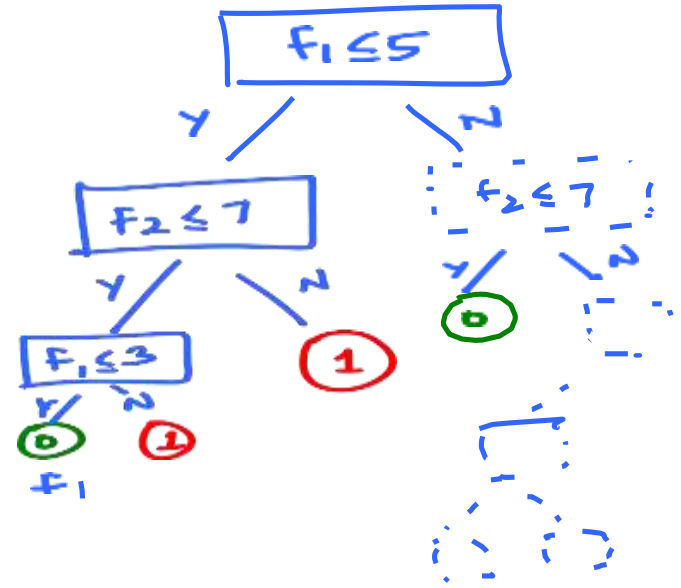
Instance space

Label space.

Assume $D \sim X \times Y$

$R^d$

$\{0,1\}$

Both training and test data are drawn from the same distribution!

## If one is given access to D, what is the best classifier?

$$h^* = \arg\min_{h \in \mathcal{H}} \mathbb{E}_{(x,y) \sim D}[\mathbb{I}(h(x) \neq y))]$$

In words, h* is the classifier that minimizes the average **test** error

## Formal Treatment of the Learning Problem

$$h^* = \arg\min_{h \in \mathcal{H}} \mathbb{E}_{(x,y) \sim D}[\mathbb{I}(h(x) \neq y))]$$

Expectation of a indicator function is probability. So,

$$h^* = \arg\min_{h \in \mathcal{H}} \mathbb{P}_{(x,y) \sim D}[h(x) \neq y)]$$

**P(x) is Uniform over {x1,x2,x3,x4,x5}**



Prob of y/x = -1

Prob of y/x = 1

Data points

**A typical dataset from this distribution**

| x5 | 1 |
|----|-----|
| x3 | -1 |
| x2 | 1 |
| x4 | 1 |
| x3 | -1 |
| x1 | 1 |

Prob of y/x = -1

Prob of y/x = 1

A sub-optimal classifier

Prob of y/x = -1

Prob of y/x = 1

Optimal classifier – Gives 0 test error!!

What if there is no classifier can make zero error?

**P(x)  is  Uniform over {x1,x2,x3,x4,x5}**



**A typical dataset
from this distribution**

| x5 | 1 |
|----|----|
| x3 | -1 |
| x2 | 1 |
| x4 | 1 |
| x2 | -1 |
| x4 | -1 |

Error = + + + +

Is this classifier optimal?

**How should the "best" classifier predict?**

Error =

# Formal Treatment of the Learning Problem

$$h^* = \arg\min_{h \in \mathcal{H}} \mathbb{E}_{(x,y) \sim D}[\mathbb{I}(h(x) \neq y))]$$

Expectation of a indicator function is probability. So,

$$h^* = \arg\min_{h \in \mathcal{H}} \mathbb{P}_{(x,y) \sim D}[h(x) \neq y)]$$

BAYES OPTIMAL CLASSIFIER

$$h^*(x) = \{1 \ \ if \ \ P(y/x) \geq 0.5 \ \ and \ \ -1 \ \ otherwise\}$$

## GOOD NEWS

We know the form of the best classifier

## BAD NEWS

We don't know the distribution D over X,Y

We will make assumptions about distribution generating data

TYPES OF MODELING

$\rightarrow$ GENERATIVE MODEL

$\rightarrow$ DISCRIMINATIVE MODEL

Gen. Model

$$P(x, y)$$

Discriminative model

$$P(y|x)$$

Eg: K-NN

: Decision-tree.

$$y|_x \sim N(\vec{\beta}x, \overset{2}{\sigma})$$

**Note that in both models, we only need P(y/x)**

# Generative Models

$$P(x,y) \;=\; \boxed{P(x)} \cdot \boxed{P(y/x)} \;=\; \boxed{P(y)} \; \boxed{P(x/y)}$$

Data: $\left\{ (x_1, y_1) \;\cdots\; (x_n, y_n) \right\}$

$x_i \in \{0, 1\}^d \qquad y_i \in \{0, 1\}$

Eg: Spam-Classification.

\# words in dictionary

"Hello, how are you?"

ARE   HELLO   HOW   YOU
$\begin{bmatrix} 0 & 1 & 0 & \cdots & 0 & 1 & 0 \cdots 0) & \cdots & 1 \end{bmatrix}$

# GENERATIVE STORY

- **STEP 1:** Decide the label by tossing a coin

$$P(y_i = 1) = p$$

- **STEP 2:** Decide features given label using $P(x_i / y_i)$

$\underline{d=3}$

Lottery　Money　Hello

1　　1　　0

$y=1$
Spam

$y=0$
non-spam



110

000　010

110

How many parameters?　　$1 + (2^d - 1) + (2^d - 1)$

$= 2^{d+1} - 1$

Issue:　• Too many parameters

• Need Alternate Story.

$P(y=1) = p$

$y=1$

$y=0$

$P(x|y=1)$

$P(x|y=0)$

$y=1$

$P_1'$  $P_2'$  $\cdots$  $P_d'$

$P_1$  $P_2$  $\cdots$  $P_d$

from $1^{st}$

$d$

$d$

# parameters ?

$1 + d + d$

$= \underline{2d+1}$

Step 1: $\qquad P(y=1) = p$

$f_i \in \{0,1\}$

Step 2: $\qquad P\left( x = [\overset{\overset{p_i^y}{\downarrow}}{f_1} \; f_2 \cdots \; f_d] \, / \, y \right)$

$$= \prod_{k=1}^{d} \left(p_k^y\right)^{f_k} \left(1 - p_k^y\right)^{(1-f_k)} \qquad \leftarrow$$

ASSUMPTION: Features are "CONDITIONALLY INDEPENDENT"
given label

Parameters to estimate :

$$p, \quad \{ p_1^1, \ldots, p_d^1 \}, \quad \{ p_1^0, \ldots, p_d^0 \}$$

## Maximum Likelihood estimators

$$\hat{p} = \frac{1}{n} \sum_{i=1}^{n} y_i \quad \leftarrow \left[ \text{Fraction of Spam emails in data} \right]$$

label

$j^{th}$ word

$$\hat{p}_j^y = \frac{\sum_{i=1}^{n} \mathbb{1}\left( f_j^i = 1, \, y_i = y \right)}{\sum_{i=1}^{n} \mathbb{1}\left( y_i = y \right)} \qquad \left[ \begin{array}{l} \text{Fraction of} \\ y\text{-labeled} \\ \text{emails that} \\ \text{contain } j^{th} \text{word} \end{array} \right]$$

Given $x_{test} \in \{0, 1\}^d$, [test email]

What is $y_{test}$?

Predict 1 if $P\left(y_{test} = 1 \,/\, x_{test}\right) > P\left(y_{test} = 0 \,/\, x_{test}\right)$

Predict 0 otherwise.

## BAYES RULE

$$P\left(y^{test} = 1 \,/\, x^{test}\right) = \frac{P\left(x^{test} \,/\, y^{test}_{=1}\right) \cdot P\left(y^{test} = 1\right)}{P\left(x^{test}\right)}$$

Say $\quad x^{test} = \begin{bmatrix} f_1 & f_2 & \cdots & f_d \end{bmatrix} \quad \in \{0, 1\}^d$

① $\quad P\left(y^{test} = 1 \mid x^{test}\right) \propto \left[ \prod_{k=1}^{d} \left(\hat{p}_k^1\right)^{f_k} \left(1 - \hat{p}_k^1\right)^{(1-f_k)} \right] \cdot \hat{p}$

② $\quad P\left(y^{test} = 0 \mid x^{test}\right) \propto \left[ \prod_{k=1}^{d} \left(\hat{p}_k^0\right)^{f_k} \left(1 - \hat{p}_k^0\right)^{(1-f_k)} \right] \cdot (1-\hat{p})$

If ① > ②, predict $y^{test} = 1$

predict $y^{test} = 0$ otherwise.

MODEL uses 2 key things

- CLASS CONDITIONAL INDEPENDENCE

- BAYES THEOREM

  - Naive assumption.
  - may not hold.
  - works well in practice.

NAIVE - BAYES algorithm

# Pitfall to watch out for

$\rightarrow$ If a word does not appear in the train set, but appears in the test set, both $\hat{p}_j^1 = 0$ and $\hat{p}_j^0 = 0$, this can't be predicted.

## Possible fix

$\rightarrow$ Add 2 pseudo emails to data.

$x = [1 \ 1 \ 1 \ \cdots \ 1]$, $y = 1$

$x = [1 \ 1 \ 1 \ 1 \ 1 \cdots 1]$, $y = 0$

✓ LAPLACE SMOOTHING.

## DECISION BOUNDARY ?

Predict $y_{test} = 1$ if $\dfrac{P\left(y_{test} = 1 / x_{test}\right)}{P\left(y_{test} = 0 / x_{test}\right)} \geqslant 1.$

$$\log\left(\frac{P(y_t = 1 / x_t)}{P(y_t = 0 / x_t)}\right) \geqslant 0$$

$$\log\left(\frac{P(x_t / y_t = 1) \cdot P(y_t = 1)}{P(x_t / y_t = 0) \cdot P(y_t = 0)}\right) \geqslant 0$$

$$\Rightarrow \quad \log\left(\left(\prod_{i=1}^{d}\left(\frac{\left(\hat{p}_i^1\right)^{f_i}\left(1-\hat{p}_i^1\right)^{(1-f_i)}}{\left(\hat{p}_i^0\right)^{f_i}\left(1-\hat{p}_i^0\right)^{(1-f_i)}}\right)\right)\cdot\left(\frac{\hat{p}}{1-\hat{p}}\right)\right) \geq 0$$

$$\sum_{i=1}^{d}\left[f_i \log\left(\frac{\hat{p}_i^1}{\hat{p}_i^0}\right) + (1-f_i)\log\left(\frac{1-\hat{p}_i^1}{1-\hat{p}_i^0}\right)\right] + \log\left(\frac{\hat{p}}{1-\hat{p}}\right) \geq 0$$

$$X_{test} = [f_1, \ f_2 \ldots \ f_d]$$

Predict 1 if

$$\underbrace{\sum_{i=1}^{d} f_i \ \log\left(\frac{\hat{p}_i^1 \left(1-\hat{p}_i^0\right)}{\hat{p}_i^0 \left(1-\hat{p}_i^1\right)}\right)}_{w_i} + \underbrace{\sum_{i=1}^{d} \log\left(\frac{1-\hat{p}_i^1}{1-\hat{p}_i^0}\right) + \log\left(\frac{\hat{p}^1}{1-\hat{p}}\right)}_{b} \geqslant 0$$

$$X_{test} = [f_1 \quad f_2 \ldots f_d]$$

$\vec{w}$

$\leftarrow \vec{w}^T x + b = 0$

$\leftarrow \vec{w}^T x = 0$

Data :  $\{(x_1, y_1), \cdots, (x_n, y_0)\}$

$x_i \in \boxed{\mathbb{R}^d} \leftarrow$

$y_i \in \{0, 1\}$

A Generative Story



$x/y=1 \sim$
$N(\mu_1, \Sigma)$

$x/y=0$
$\sim N(\mu_0, \Sigma)$

Note: In this model, Co-variances are same.

# Parameters:

$$1 + d + d + o(d^2)$$

$$\uparrow \quad \uparrow \quad \uparrow$$

$$p \quad M_1 \quad M_0$$

## Maximum-likelihood estimator

$$\hat{p} = \frac{1}{n} \sum_{i=1}^{n} y_i \qquad \text{(Fraction of points labelled 1)}$$

$$k \in \{0,1\} \quad \hat{A}_k = \frac{\sum_{i=1}^{n} \mathbb{1}(y_i = k) \cdot x_i}{\sum_{i=1}^{n} \mathbb{1}(y_i = k)}$$

$$\hat{\Sigma} = \frac{1}{n} \sum_{i=1}^{n} \left( x_i - \hat{M}_{y_i} \right) \left( x_i - \hat{M}_{y_i} \right)^\top \qquad \left[ \text{argue this} \right]$$

Prediction

As usual, by Bayes rule.

Predict 1 if $P\left( y_t = 1 \mid x_t \right) > P\left( y_t = 0 \mid x_t \right)$

$$P\left( x_t \mid y_t = 1 \right) \cdot P(y_t = 1) \;>\; P\left( x_t \mid y_t = 0 \right) \cdot P(y_t = 0)$$

$$\frac{P\left(x_t|_{y_t=1}\right) \cdot P(y_t=1)}{f\left(x_t; \hat{M_1}, \hat{\Sigma}\right) \cdot \hat{P}} > \frac{P\left(x_t|_{y_t=0}\right) \cdot P(y_t=0)}{f\left(x_t; \hat{M_0}, \hat{\Sigma}\right) \cdot (1-\hat{P})}$$

$$\Rightarrow \quad e^{-\frac{1}{2}\left(x_t - \hat{M_1}\right)^T \hat{\Sigma}^{-1} \left(x_t - \hat{M_1}\right)} \cdot \hat{P} > e^{-\frac{1}{2}\left(x_t - \hat{M_0}\right)^T \hat{\Sigma}^{-1} \left(x_t - \hat{M_0}\right)} \cdot (1-\hat{P})$$

$y_{test} =$
Predict 1

$\Rightarrow$

$$2\left(\hat{A}_1 - \hat{A}_0\right)^T \underbrace{\hat{\Sigma}^{-1}}_{} \; x_{test} \;\; + \;\; \underbrace{\hat{A}_0^T \hat{\Sigma}^{-1} \hat{A}_0 \;-\; \hat{A}_1 \hat{\Sigma}^{-1} \hat{A}_1 \;+ \log\left(\frac{1-\hat{p}}{\hat{p}}\right)}_{b} \;\geq 0$$

$$\underbrace{2\left(\hat{A}_1 - \hat{A}_0\right)^T \hat{\Sigma}^{-1}}_{W}$$

$$W^T x_{test} + b \geq 0 \quad \Rightarrow \quad \text{DECISION FUNCTION}$$
$$\text{IS LINEAR}$$

$$\text{only if } \Sigma \text{ is same for } x/y=0 \quad x/y=1$$

Same $\varepsilon$

$x_{test}$

$\hat{M}_1$

$\hat{N}_0$

$A_1 - A_0$

$(A - \lambda I)|_{x=0}$

Different $\Sigma$

Decision function
is Quadratic

$\hat{\mu}_1$

$\hat{\mu}_0$

▶ Gaussian Naive Bayes.

Question :

$\rightarrow$ Can we directly make linear assumption about $P(y/x)$

$$P\left(y=1/x\right) = 1 \quad \text{if} \quad \underline{w^T x} \geq 0$$
$$\phantom{P\left(y=1/x\right)} \qquad \qquad +b$$
$$\phantom{P\left(y=1/x\right)} = 0 \quad \text{otherwise.}$$

Dataset is
not allowed in
our model

LINEAR
SEPARABILITY
ASSUMPTION

Allowed.

$$P\left(y=1\mid x\right) \quad = \quad 1 \qquad \text{if} \qquad \vec{w}^{\scriptscriptstyle T} x \geqslant 0$$

$$= \quad 0 \qquad \text{o/w}$$

Linear
Separability
assumption.

Lin·Sep assumption :

$$\exists w \in \mathbb{R}^d \quad \text{s.t} \quad \underset{\downarrow}{\text{sign}}\left(\vec{w}^{\scriptscriptstyle T} x_i\right) = y_i \quad \forall i \in [n]$$

$$\{1, \ldots, n\}.$$

$$\text{sign}(z) = \begin{cases} +1 & \text{if } z \geqslant 0 \\ -1 & \text{o/w} \end{cases}$$

$$\text{DATA } \{(x_1, y_1) \dots (x_n, y_n)\}$$

$$\min_{w \in \mathbb{R}^d} \sum_{i=1}^{n} \mathbb{1}\left( \text{sign}(w^T x_i) \neq y_i \right) \rightarrow \text{NP-HARD}$$
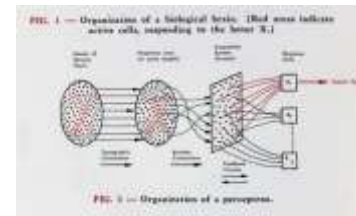
We know in general the problem of finding best w is NP-hard.

But is it still hard under linear separability assumption?

# PERCEPTRON



Frank Rosenblatt '50, Ph.D. '56, works on the "perceptron" –
what he described as the first machine "capable of having an original idea."

# PERCEPTRON - ALGORITHM

Input:   $\{ (x_i, y_i) \cdots \cdots , \qquad\qquad x_i \in \mathbb{R}^d$
$\qquad\qquad\qquad\qquad\qquad\qquad\qquad y_i \in \{ +1, -1 \}$

$$W^{0} = [0 \ 0 \cdots 0] \qquad 0 \in \mathbb{R}^d$$

<— Iteration

until   convergence

$\rightarrow$   Pick $(x_i, y_i)$   from   dataset

IF $\left( \text{Sign}(W^{t^T} x_i) = y_i \right)$

do   nothing

else

$\{ \ W^{t+1} = W^t + x_i y_i \in \{ +1, -1 \}$   $\mathbb{R}^d$   UPDATE RULE.

end.

end

## UPDATE RULE

$$W^{t+1} = W^t + x_i y_i$$

**Two types of mistake**

**Type 1**

Pred $\geqslant 1$
act $\rightarrow -1$

**Type 2**

Pred $\rightarrow -1$
act $\rightarrow 1$

**Type-1**

Pred $\rightarrow 1$
act $\rightarrow -1$

$$W^{t^T} x_i \geqslant 0 \quad \text{but} \quad y_i = -1$$

$$W^{t+1^T} x_i = \left( W^t + x_i y_i \right)^T x_i$$

$$= W^{t^T} x_i + y_i \| x_i \|^2$$
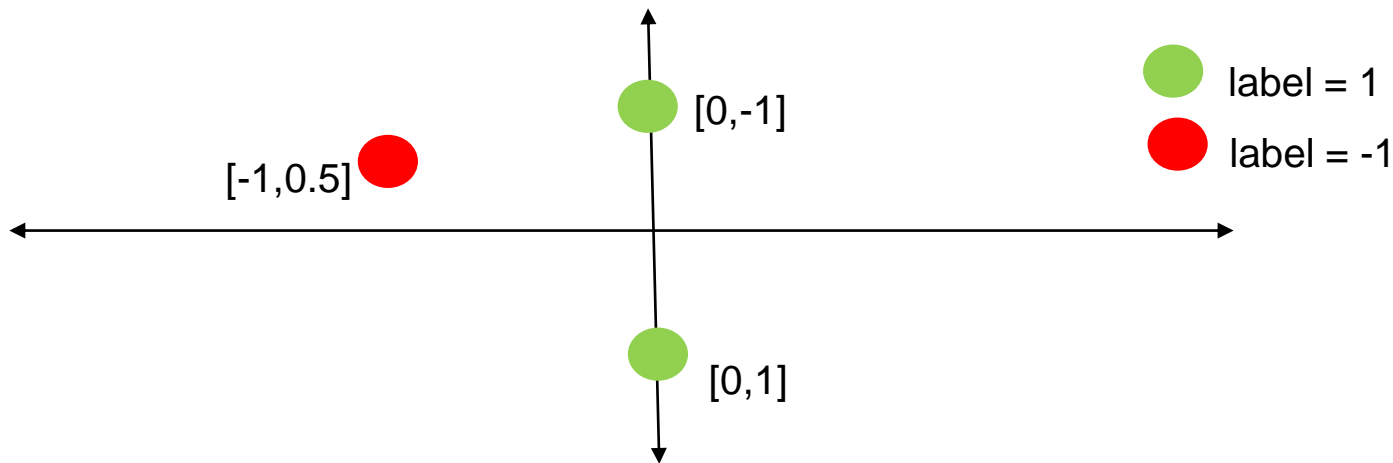
$\geqslant 0 \qquad < 0$

$\leq 0 \qquad > 0$

**Type-2**

**1**

**-1**

**Fixing one error might lead to more errors elsewhere**

In general, does perceptron work for linearly separable data?

Recall

$\underline{\text{Lin. Sep assumption}}$ :    $\exists w \in \mathbb{R}^d$   s.t   $\underline{\text{Sign}(w^T x_i)} = y_i$   $\forall i \in [n]$

$$\text{Sign}(z) = \begin{cases} +1 & \text{if } z \geq 0 \\ -1 & \text{olw} \end{cases}$$

$\{1, \ldots, n\}$

[0,-1]

[-1,0.5]

[0,1]

● label = 1

● label = -1

**Is this a linearly separable dataset?**

label = 1

label = -1

[0,-1]

[-1,0.5]

[0,1]
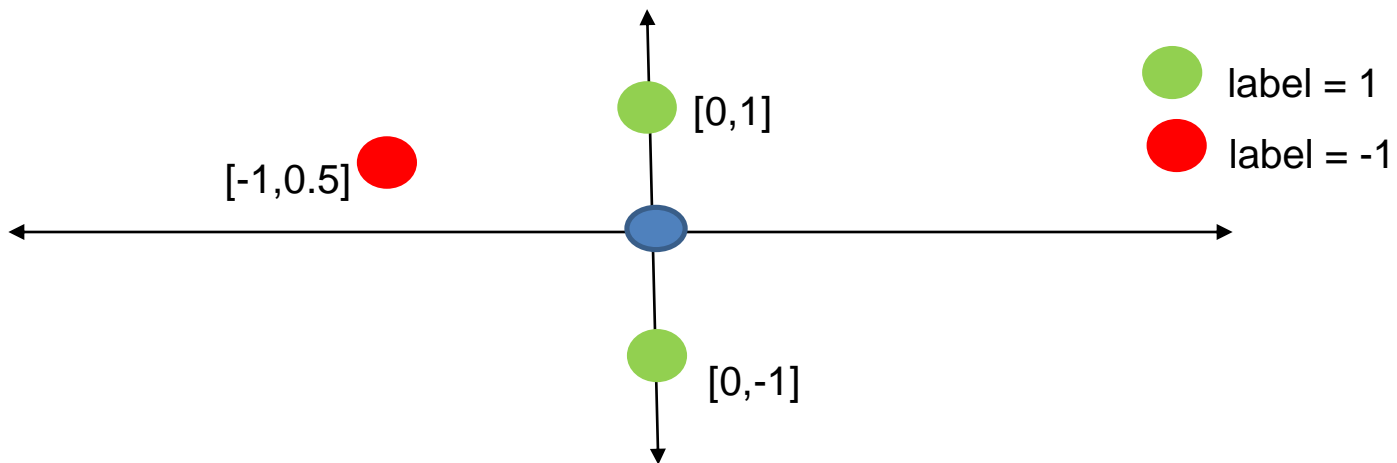
Any w in the positive x-axis linearly separates the data.
The dataset is linearly separable.

Let's see what perceptron learns from this data!

| | Predicted label | True label |
|---|---|---|
| [0 1] | 1 | 1 |
| [0 -1] | 1 | 1 |
| [-1 0.5] | 1 | -1 |

$w^0 = [0\ 0]$

| | label = 1 |
| | label = -1 |

[0,1]

[-1,0.5]

[0,-1]

$w^1 = [1\ -0.5]$

| | Predicted label | True label |
|---|---|---|
| [0 1] | -1 | 1 |
| [0 -1] | 1 | 1 |
| [-1 0.5] | -1 | -1 |

label = 1

label = -1

[0,1]

[-1,0.5]

[0,-1]

$w^2 = [1\ 0.5]$

| | Predicted label | True label |
|---|---|---|
| [0 1] | 1 | 1 |
| [0 -1] | -1 | 1 |
| [-1 0.5] | -1 | -1 |

label = 1

label = -1

[0,1]

[-1,0.5]

[0,-1]

**PERCEPTRON
DOES NOT CONVERGE**

$w^3$= [1 -0.5]

|  | Predicted label | True label |
|---|---|---|
| [0 1] | -1 | 1 |
| [0 -1] | 1 | 1 |
| [-1 0.5] | -1 | -1 |

- Optimal $w^* = \begin{bmatrix} c \\ 0 \end{bmatrix}$ $c > 0$ has datapoints that lie on the linear separator.

- If we assume this isn't the case, will perceptron converge?

# ASSUMPTIONS

- LINEAR    SEPERABILITY    with    $\gamma$ - MARGIN



$\{x : w^T x = \gamma\}$

$\{x : w^T x = 0\}$

$\{x : w^T x = -\gamma\}$

A Dataset $D = \{ (x_1, y_1) \cdots , (x_n, y_n)$ is $\underline{L \cdot S}$

with $\gamma$-margin if $\exists \; w^*$ s.t

$$(w^{*T} x_i) \; y_i \geq \gamma \qquad \forall i \qquad \text{for some} \qquad \gamma > 0$$

# PERCEPTRON

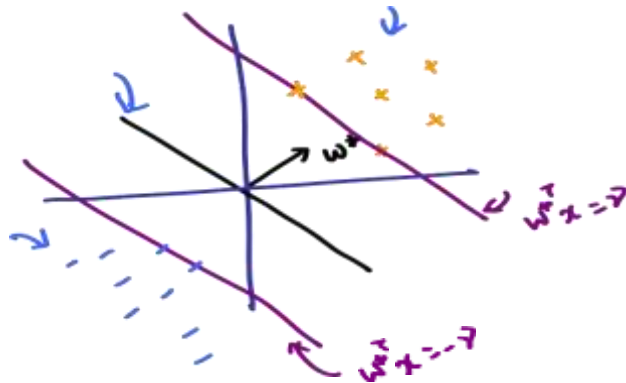$$w^{t+1} = w^t + x_i y_i$$

## ASSUMPTIONS

① • Linear separability with $\geq$ margin

A dataset $\{(x_1, y_1), \dots, (x_n, y_n)\}$ is L.S with $\geq$ margin
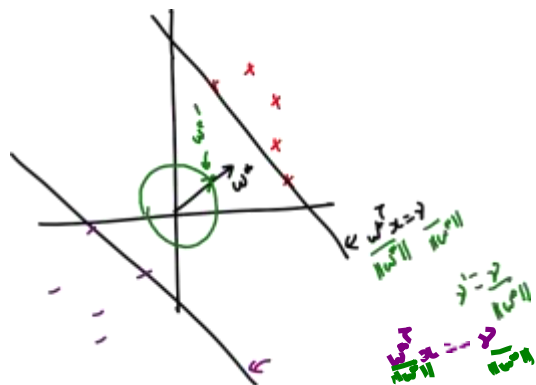
if $\exists\, w^*$ s.t $(w^{*T} x_i) y_i \geq \gamma$ $\forall i$ for some $\gamma > 0$

②   RADIUS   ASSUMPTION

$$\forall i \in D, \quad \| x_i \|_2 \leq R \quad \text{for some} \quad R > 0$$

③   Without loss of generality

$$\| w^* \| = 1$$

## ANALYSIS OF "mistakes" of Perceptron.

- observe that an update happens only when a "mistake" occurs.

- Say the current guess $= w_\ell$ and a mistake happens. w.r.t $(x,y)$

$$w_{\ell+1} = w_\ell + x \cdot y$$

$$\|w_{\ell+1}\|^2 = \|w_\ell + x \cdot y\|^2$$

$$= \left( w_\ell + x \cdot y \right)^T \left( w_\ell + x \cdot y \right)$$

$$\| w_{\ell+1} \|^2 = \| w_\ell \|^2 + 2 \underbrace{\left( w_\ell^T x \right) \cdot y}_{\leq 0} + \underbrace{\left( x^T x \right)}_{\|x\|^2 \leq 1} \cdot y^2$$

$$\underset{\substack{\leq 0 \\ \text{because} \\ \text{[mistake.]}}}{} \qquad \underset{\leq R^2}{}$$

Inductively

$$\leq \| w_\ell \|^2 + R^2 \leq \left( \| w_{\ell-1} \|^2 + R^2 \right) + R^2$$

$$\underbrace{\|w_0\|^2 + \ell R^2}$$

$$\boxed{\| w_{\ell+1} \|^2 \leq \ell \cdot R^2} \qquad \underline{\qquad} ①$$

$$w_{\ell+1}^T \, \omega^* \quad = \quad \left( w_\ell + x \cdot y \right)^T \omega^*$$

$$= \quad w_\ell^T \omega^* + \left( \omega^{*T} x \right) \, \underline{y}$$

$$\underbrace{\qquad\qquad}_{} \geq \gamma$$

$$\geq \left( w_{\ell-1}^T \omega^* + \underline{\gamma} \right) + \gamma$$

$$\vdots$$

$$\boxed{ w_{\ell+1}^T \, \omega^* \quad \geq \quad \ell \cdot \gamma }$$

$$\left(\frac{(x^Ty)}{\|y\|^2}\right) \cdot y$$

$$\|x\|^2 \geq \left\|\left(\frac{x^Ty}{\|y\|^2}\right)y\right\|^2$$

$$\geq \frac{(x^Ty)^2}{\|y\|^{4\,2}} \cdot \|y\|^2$$

$$\Rightarrow (x^Ty)^2 \leq \|x\|^2\|y\|^2$$

$$x^Ty \leq \|x\|\|y\|$$

(C.S)    Cauchy
        Schwartz    $\longrightarrow$
        inequality

From before

$$W_{\ell+1}^T \, w^* \geq \ell \cdot \gamma$$

$$\|W_{\ell+1}\|^2 \|w^*\|^2 \geq \left(W_{\ell+1}^T \, w^*\right)^2 \geq \ell^2 \gamma^2$$

$$\uparrow$$
$$\text{c.s}$$

$$\|W_{\ell+1}\|^2 \geq \ell^2 \gamma^2 \qquad - ②$$

Combining ① & ②.

$$\ell^2 \gamma^2 \le \| w_{\ell+1} \|^2 \le \ell R^2 \quad \longrightarrow ①$$
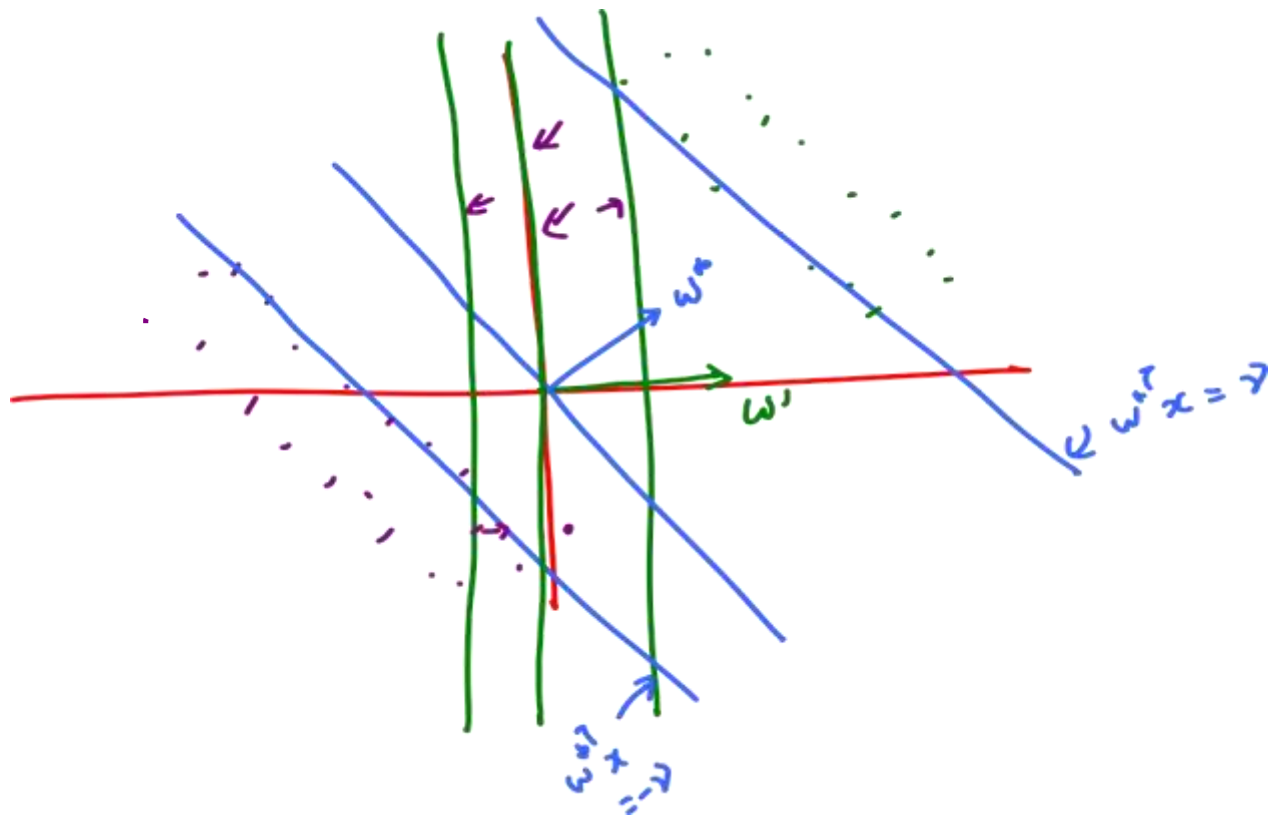
② $\longrightarrow$

$$\Rightarrow \quad \ell^2 \gamma^2 \le \ell R^2$$

$$\boxed{\ell \le \frac{R^2}{\gamma^2}}$$

RADIUS- MARGIN
BOUND.

$\Rightarrow$ # mistakes of Perceptron is bounded

$\Rightarrow$ Perceptron converges!

$w^*$

$w'$

$w^* x = 7$

$w^* x = 7$

Perceptron's # mistakes
depends on $w^*$.
But it might
output $w'$.