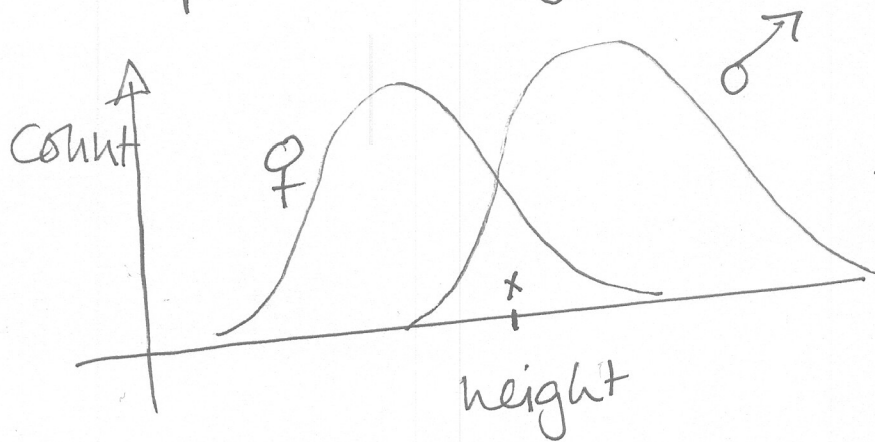# Classification:

- Input Features, output one bit
- (more complicated models later)

## example:

input:  height          output: gender



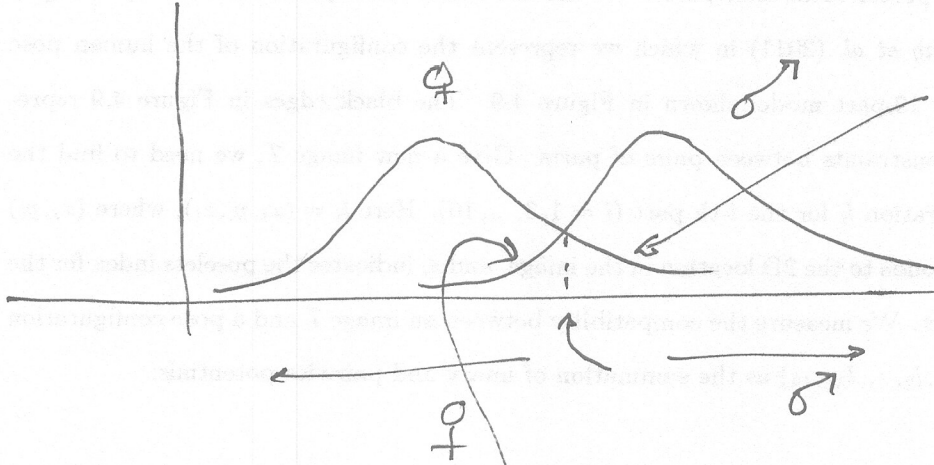← histograms of height w/ gender

- For the moment, assume that ♀ ♂ are evenly dist.

  - how would one classify?
    - choose a height threshold
    - mistakes are inevitable
    - we need to choose the least expensive.

# Notice guaranteed error



we will get these ♀'s wrong.

we will get these ♂'s wrong

## strategy

$$h > t \quad ♂$$
$$\text{otherwise} \quad ♀$$

Cases:
- males, females equally common
  $$t \text{ at } x$$

- males very common $♀$'s rare
  $$t \ll x$$

- $♂$ rare, $♀$ common
  $$t \gg x$$

Q: <u>reasonable way to set $t$</u>

<u>Choose $t$ to produce the minimum
expected-cost of errors</u>

- two types of error
  $$( ♂ \rightarrow ♀ )$$
  $$( ♀ \rightarrow ♂ )$$

we assume reward for <u>right</u> answer is 0

we have a feature $x$.

if we say ♀, we get $\begin{cases} 0 & \text{right} \\ L(\male \to \female) & \text{wrong} \end{cases}$

if we do this many times,
we get 0 with frequency

$$p(\female | x)$$

and $L(\male \to \female)$ with freq $p(\male | x)$

So expected loss of ♀ is

$$0 \cdot p(\female | x) + L(\male \to \female) p(\male | x)$$

Similarly, expected loss of ♂ is

$$0 \, p(\male | x) + L(\female \to \male) p(\female | x)$$

so in principle we have a rule ④

at $x$,
$$\begin{cases} \text{cost of } ♀ \text{ is: } L(♂→♀) \, P(♀|x) \\[1em] \text{cost of } ♂ \text{ is: } L(♀→♂) \, P(♂|x) \end{cases}$$

. choose the least expensive, say that.

. But where do we get $P(♂|x)$, $P(♀|x)$?

1) $P(♀|x) = 1 - P(♂|x)$     (only 2 options)

2) ~~$P(♀$~~   $P(♂|x) = \dfrac{P(x|♂) \, P(♂)}{P(x)}$

$$= \frac{P(x|♂) \, P(♂)}{\left[ P(x|♂) \, P(♂) + P(x|♀) P(♀) \right]}$$

we could read this ←     ↑    ↑ prior
off the histograms.

$\boxed{1 - P(♂)}$

We can now build one useful form of classifier.

- measure $p(x|♀)$, $P(♀) = \pi$, $p(x|♂)$

  (say, histogram)

- say

  ~~♂~~ ♀ if $L(♀→♂) \cdot p(♀|x) \overset{\gtrless}{\underset{<}{=}} L(♂→♀) \cdot p(♂|x$

  ~~♂~~ ♂

  doesn't matter

now, consider

$$L(♀→♂) \, p(♀|x) = L(♂→♀) \, p(♂|x)$$

all that matters is $\quad R = \dfrac{L(♀→♂)}{L(♂→♀)}$

So we care about

$$R \cdot p(♀|x) = p(♂|x)$$

i.e. $\quad \dfrac{R \cdot p(x|♀)\,\Pi}{p(x)} = \dfrac{p(x|♂)\,(1-\Pi)}{p(x)}$

i.e $\quad \underbrace{\dfrac{p(x|♀)}{p(x|♂)}} = \dfrac{(1-\Pi)}{\Pi} \cdot \dfrac{1}{R}$

$$\longrightarrow \quad \underline{\text{likelihood ratio}}$$

i.e. if

$$\dfrac{p(x|♀)}{p(x|♂)} \quad \begin{array}{c} > \\ = \\ < \end{array} \quad g(\Pi, R) \quad \begin{array}{l} \text{say } ♀ \\ \text{doesn't matter} \\ \text{say } ♂ \end{array}$$

Now equations give a way to
determine g. but we can manage
without.

Rule : $\dfrac{p(x \mid q)}{p(x \mid 0^q)} > t$ , say $q$

## Plot errors w/ varying t

## Receiver operating curve.



plot for different
t.

$t = c$

$t = b$

P(true
detn)

$t = a$

P (false detection)

Model: we are trying to detect $\vec{\sigma}$'s in a population of $\varphi$'s (or vice versa).

$$P \text{ (false detect)} = \frac{\# \text{ of } \varphi\text{'s we called } \vec{\sigma}\text{'s}}{\# \text{ of times we classified}}$$

$$P \text{ (true det)} = \frac{\# \text{ of } \vec{\sigma}\text{'s we called } \vec{\sigma}\text{'s}}{\# \text{ of } \vec{\sigma}\text{'s in population}}$$

We evaluate this on __test__ data.

Training:

- form histograms $p(x|\varphi)$ etc.
- can be hard to do with high dimensional $x$.

Major Problem:
- a 1-D histogram w/ n cells in each
  - dir has n cells
  - 2D $n^2$
  - 3D $n^3$
  - d D $n^d$

We cannot build such histograms:

Strategies:
- Simplify the model
- model $p(q \mid x)$ directly
- Search for decision boundary directly.

# Simplify model

## model

$$p(x_1, x_2 \cdots x_n | q)$$
$$= p(x_1 | q) \, p(x_2 | q) \, p(x_3 | q) \cdots p(x_n | q)$$

This model is usually wrong.
- But it's convenient
- and works surprisingly well

## Naive Bayes:

- method: 
  - one histogram each in each direction.
  - form likelihood ratio
  - test against threshold.

B) Model $p(q/x)$ directly.

- parametric models, perhaps later

-

C) Find <u>decision boundary</u>:

- By continuity reasoning
  (Nearest neighbours)
- By search.

<u>Nearest neighbours</u>:

<u>alg</u>:

- find $x_i \in$ examples such that $\|x_i - x\|^2$ is smallest
- class of $x$ is class of $x_i$

## $k - l$ nearest neighbors:

alg:
- find the $k$ $x_n \in$ examples that are closest to $x$.
- find the most common class in these neighbors
- if there are $l$ in this class, classify $x$ with that class, otherwise, don't know.
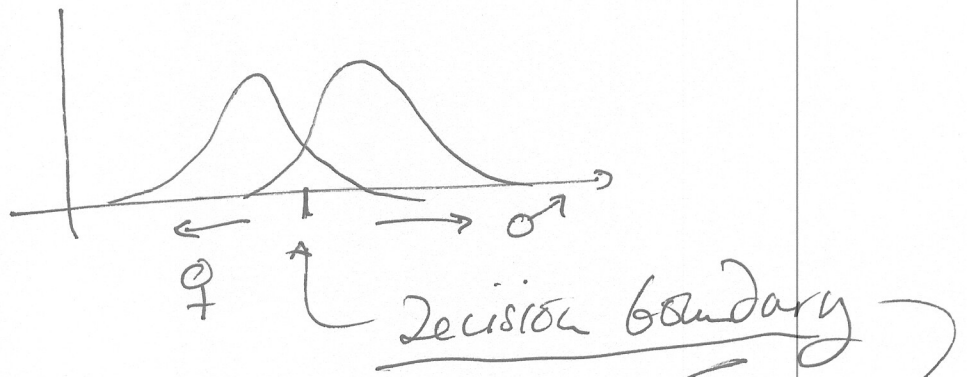
## Properties:

- with enough examples, error rate is no more than $2 \times$ best possible
- we must worry about scaling
  - dimension

- 
- Algorithmicly complex.

Ⓑ model $P(\mathcal{G}|x)$ directly

(we'll talk about this later)

Ⓒ <u>Find decision boundary directly</u>

. recall


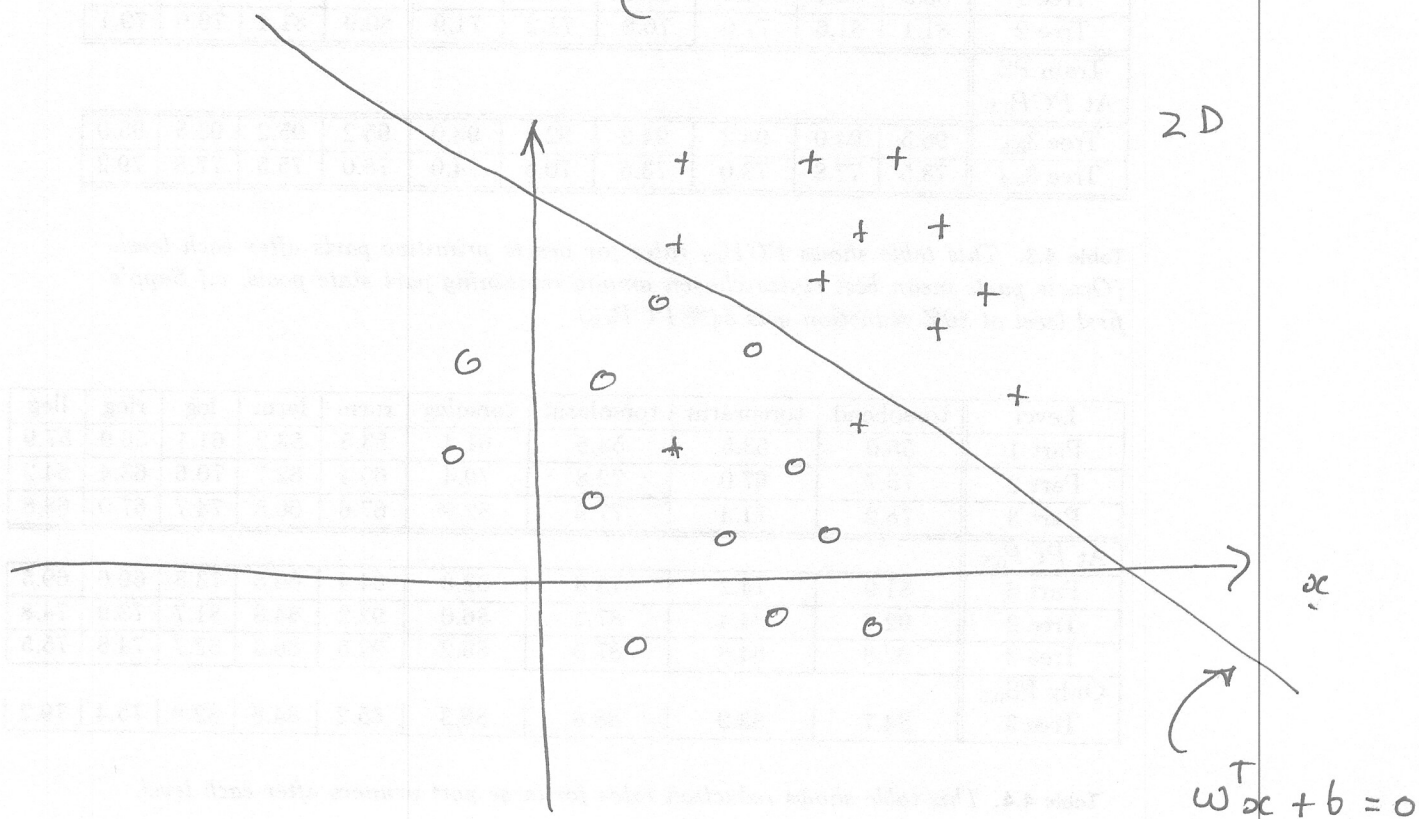
. eg in 2D



. Hard to search for a curve in 2D
or more-D

# Easy, highly successful strategy
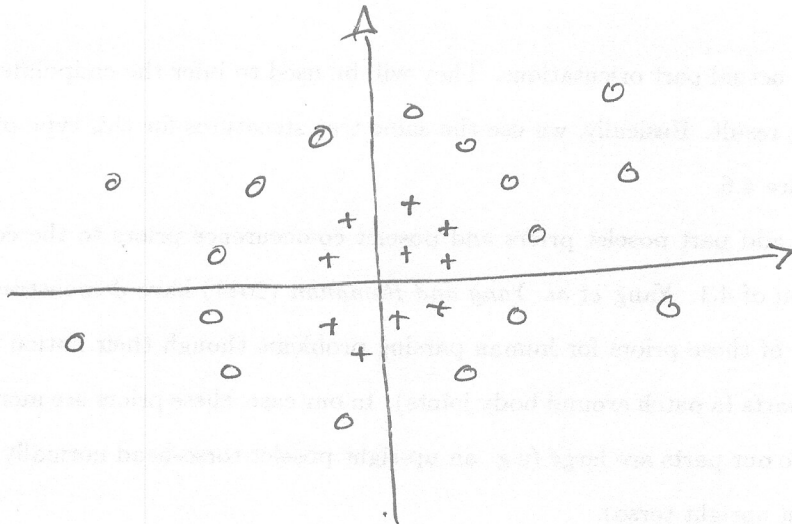
- decision boundary is a flat (line, plane, hyperplane)

- ie.

$$\left(\omega^T x + b\right) \begin{cases} > 0 & \text{class 1} \\ = 0 & \\ < 0 & \text{class 2} \end{cases}$$
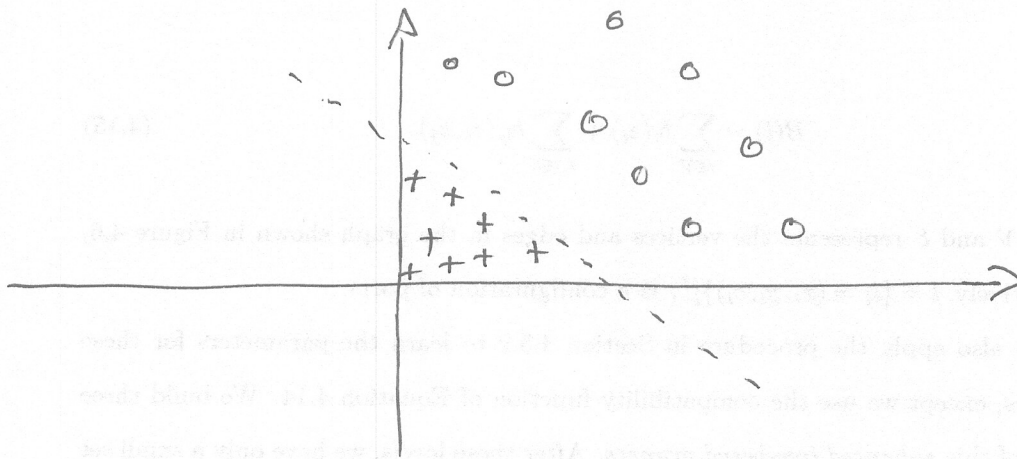
2D

$\omega^T x + b = 0$
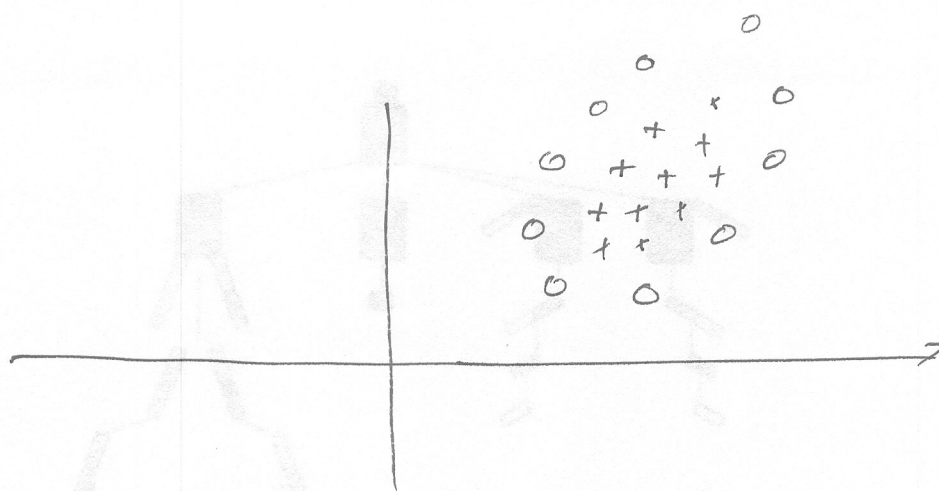
clearly, this won't work always

eg



but: map these points by

$$(x, y) \longrightarrow (x^2, y^2)$$

but

$$(x, y) \longrightarrow (x^2, xy, y^2, x, y)$$

( recall general ellipse is
$$ax^2 + bxy + cy^2 + dx + ey + f = 0 $$ )
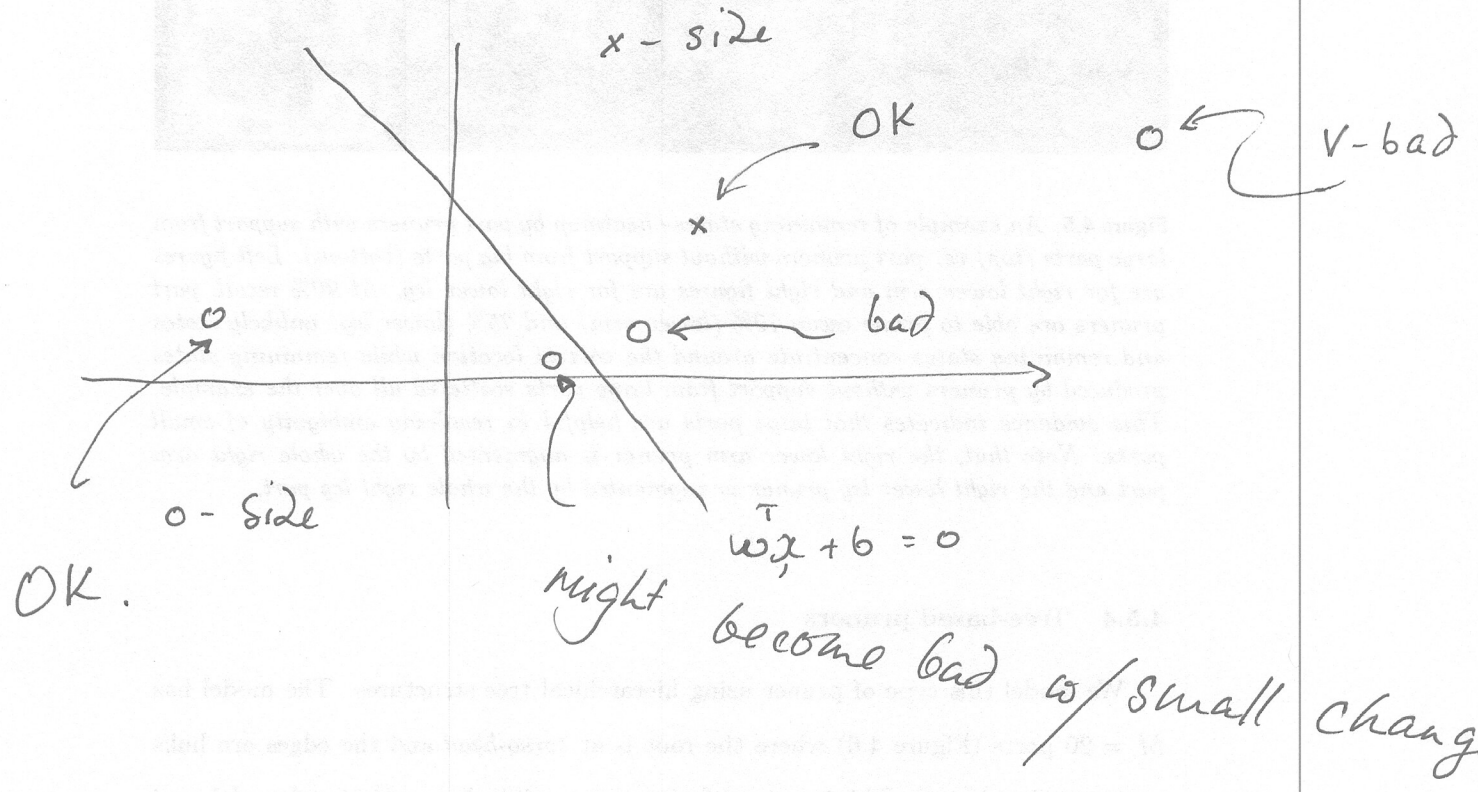
and we get linear boundary.

## General principle here :

- with enough features a linear classifier will behave well.
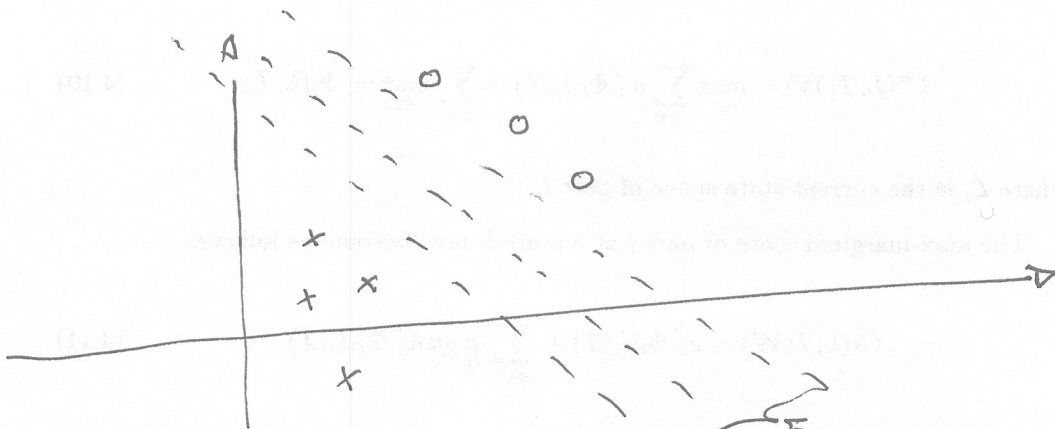
# How to choose a linear classifier

Q: what is $\omega, b$ ?

A:      minimize     <u>loss</u>     of using classifie

$x$ - side

OK

O ← V-bad

bad

O - side

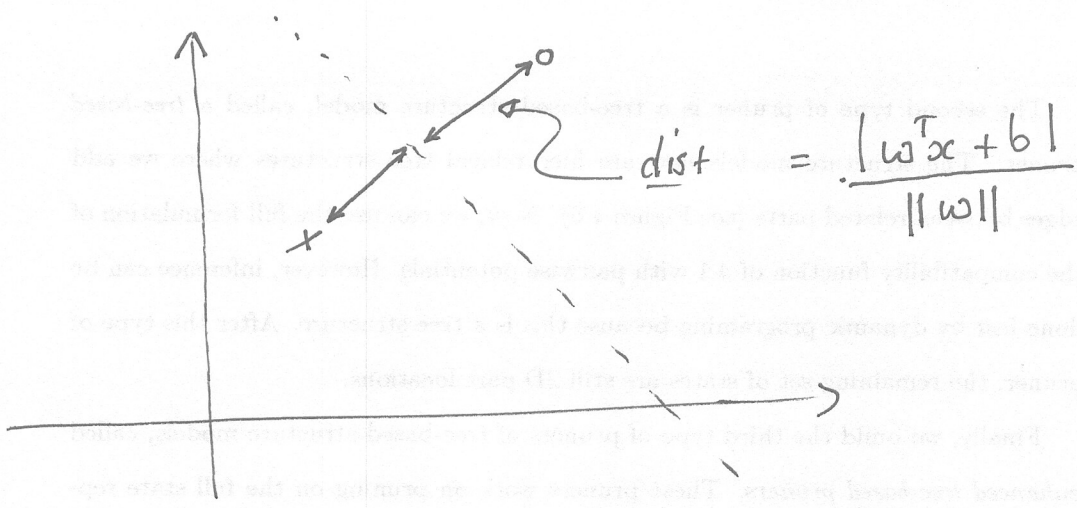$\omega^T x + b = 0$

OK.

might become bad w/ small change.

# We need to deal w/ easy cases



each of these lines looks OK — but which do we choose?

- ## Good choice
  - closest examples should be as far away as possible
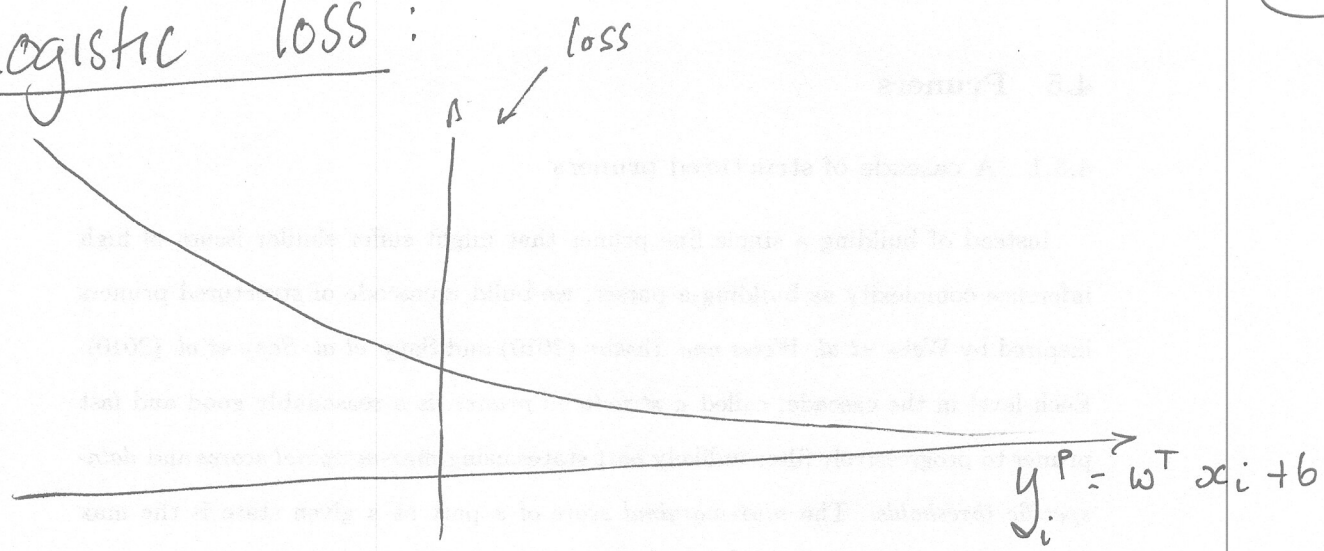
$$dist = \frac{|w^T x + b|}{\|w\|}$$

- hence, if loss is zero, we would like to minimize $\|w\|^2$

- if loss is non zero, small $\|w\|^2$ is a good idea

- Hence minimize

$$Loss + \theta \|w\|^2$$
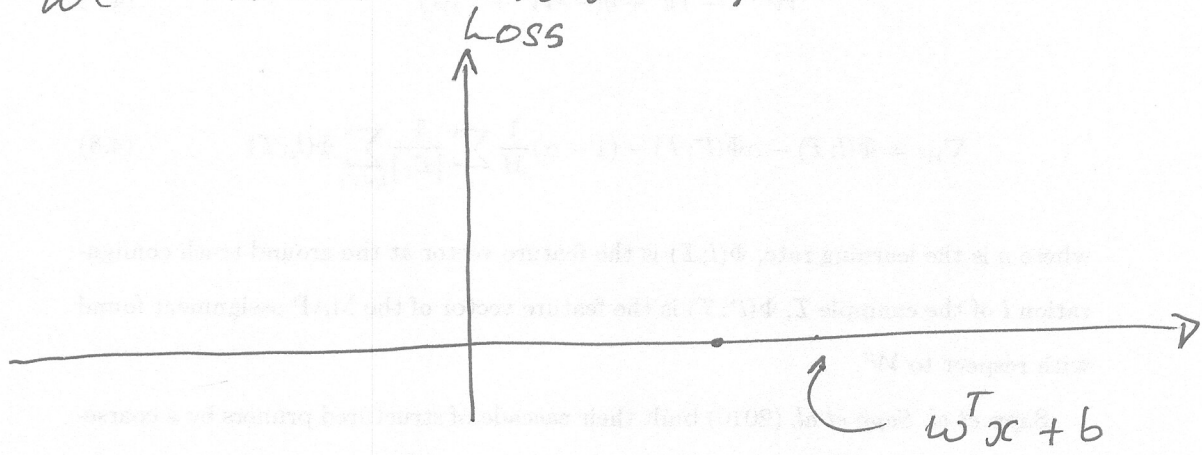
↑ weighting parameter choose later

# Logistic loss:

loss

$$y_i^p = w^T x_i + b$$

$$\log\left[1 + \exp\left(-y_i \cdot y_i^p\right)\right]$$

- leads to

$$\min \Theta \|w\|^2 + \sum_i \log\left[1 + \exp\left(-y_i \cdot y_i^p\right)\right]$$

# LOSS:

- consider an example of class 1
- we want $w^T x + b > 0$



- for $w^T x + b \gg 0$, $\text{LOSS} = 0$
- for $w^T x + b < 0$, loss is big
- for $w^T x + b \lll 0$, bigger

- loss should <u>not</u> grow too fast, otherwise one example dominates

- loss should be non-zero for small +ve $w^T x + b$

The graph above shows the hinge loss curve with axis labeled $y_i^p = w^T x_i + b$.

## Hinge loss :

- example has label $y_i \in \{1, -1\}$

- we predict $y_i^p = w^T x_i + b$

- loss is

$$\max\left\{0, \ 1 - y_i \cdot y_i^p\right\}$$

- plotted above for $y_i = 1$.

- leads to

$$\min \ \theta \|w\|^2 + \sum_i \max\left\{0, 1 - y_i y_i^p\right\}$$

which is hard.