(*) At prediction time
(∅ Naive Bayes):

(*) $P(y=1|x) = \dfrac{P(x|y=1) \cdot P(y=1)}{P(x|y=1) \cdot P(y=1) + P(x|y=0) \cdot P(y=0)}$

(*) Lets say there is a word NIPS$^{(j>6017)}$ in your dictionary but it does not appear in your email:

$$P(X_{6017}=1|y=1) = \dfrac{0}{\#\{y=1\}}.$$

→ Statistically it is a bad idea to say that something will not definately happen.

(*) This mean that during prediction:

① $P(y|x) = \dfrac{\prod\limits_{i=1}^{10000} P(x_j|y) \cdot P(y)}{\prod\limits_{i=1}^{10000} P(x_j|y=1)\cdot P(y=1) + \prod\limits_{i=1}^{10000} P(x_j|y=0)\cdot P(y=0)}$

(*) → Due to non-presence of ith word - we have not seen before in our dataset → the product highlighted in red will become zero during predict. ith word might be in your testing instance but during training it was not present so product becomes zero.

Ⓨ How to Improve this error of zero?
   ↝ Laplace smoothing

Example (Stanford football team performance)          won

| | Wake Forest | 0 |
| 9/12 | Arizona | 0 |
| 10/17 | Caltech | 0 |
| 11/21 | Oklahoma | 1 |
| 12/31 | | |

→ What should be there winning chance in fourth game?

→ $\dfrac{\text{# of wins}}{\text{# of wins} + \text{# of lost}} = \dfrac{0}{0+3} = 0$ (Absolute certainty of Losing).

→ Again statistically it is had to assume that with absolute certanity Slandford won't win next gumes

→ Here comes in $\boxed{\text{Laplace smoothing}}$ which means:

→ Add One to both ~~numerator and~~ ~~one to denominator:~~ events.

$$\frac{\# \text{Wins} + \boxed{1}}{\# \text{Wins} + \boxed{1} + \#\text{of Loses} + \boxed{1}} = \frac{1}{6} \quad (\text{More reasonable}).$$

(*) $\boxed{\text{More henrally}}$

$$X \in \{1, \dots k\}$$

Estimate: (Laplace Smoothing)

$$P(X=j) = \frac{\sum_{j=1}^{m} 1\{x^{(i)} = j\} + 1}{m + 1}$$

(X) Laplace Smoothed Naive Bayes

(a) $\emptyset_{y=} = P(y) \dfrac{\sum_{i=1}^{m} 1\{y^{(i)} = 1\} }{m}$

(a) $\emptyset_{j|y=1} = \dfrac{\sum_{i=1}^{m} \{x_j = 1, y^{(i)} = 1\} + 1}{\sum_{i=1}^{m} \{y^{(i)} = 1\} + 2}$

→ Laplace Smoothed CD:o Similar for $\emptyset_{j|y=0}$.

⑦ What to do when features are multinomial?

$$X_i \in \{1, \ldots k\}$$

| Size | <400feet | 400< <800 | 800-1200 | >1200 |
|------|----------|-----------|----------|-------|
| $X_i$ | 1 | 2 | 3 | 4 |

→ Made Buckets of this feature.

$$P(X_g | y) = \frac{1}{t} \prod_{j=1}^{m} \cdot P(X_j | y)$$

→ Multinomial mean that each feature $X_j$ can now assume more than one value rather than only binary value.

→ multinomial ⟶ (more than one classes)

→ Better Representation of Naive Bayes

(✗) So Far = (Multivariate Bernolli event model)

$$X = \begin{bmatrix} 0 \\ 0 \\ 0 \\ 1 \\ \vdots \\ 1 \\ \vdots \\ 1 \\ 0 \end{bmatrix}$$
a ← 1
aadiohe 2
$\vdots$
buy ← 800  → $X_i \in \{0,1\}$
drugs ← 1600   ↳ Disregard
now ← 6200      the count
                of wood

↳ "Drugs! Buy drugs now!"

→ New representation:
(Multinomial Event Model)

$$X = \begin{bmatrix} 1600 \\ 800 \\ 1600 \\ 6200 \end{bmatrix}$$
$\in R^n$ (where $n_i$ is the length of the sentence, in this case four for the sentence "Drugs! Buy drugs now!"

but $X_j \in \{1, \dots 10000\}$.
different
$n_i$ = length of email $i$. (Varies for each instance).

(✗)

① Lets build a generative model for Multinomial Event model:-

$$P(x, y) = P(x|y) \cdot P(y)$$

assume

$$\simeq \underbrace{\prod_{j=1}^{n} P(x_j|y)}_{\text{Multinomial}} \cdot P(y)$$

→ Depends on single instance (will vary).

↓

Multinomial

→ **Parameters**

$$\boxed{\phi_y = P(y=1)}, \boxed{\phi_{k|y=0}} = P(x_j = k | y = 0)$$

↓

② Chance of a word $\frac{k|c}{\text{if } y = 0}$ being $\underline{k|c}$, if $y = 0$

$$\phi_{k|y=1} = P(x_j = k | y = 1).$$ Of all the emails how many ↑ times word K appears

→ MLE:- (Parameters)

$$\phi_{k|y=0} = \frac{\sum_{i=1}^{m} \mathbb{1}\{y^{(i)} = 0\} \sum_{j=1}^{n_i} \mathbb{1}\{x_j^{(i)} = k\} + 1}{\boxed{\sum_{i=1}^{m} \mathbb{1}\{y^{(i)} = 0\} \cdot n^{(i)}} + \boxed{10005}}$$

Total number of words in non-spam emails ←

① Naive Bayes in quick to implement and is computationally efficient.

② Gaussian and Naive Bayes are quick to implement

③ # Support Vector Machines

⊙



Fig-1

→ Turn key Property → Support Vector Machine do not have many hyperparameter
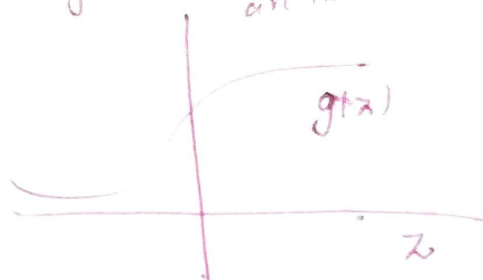
→ Optimal margin classifier (Separable case)



Fig

Ⓧ • Functional Margin: (How confident you are about an instance)

Logistic Regression
a) motivation example

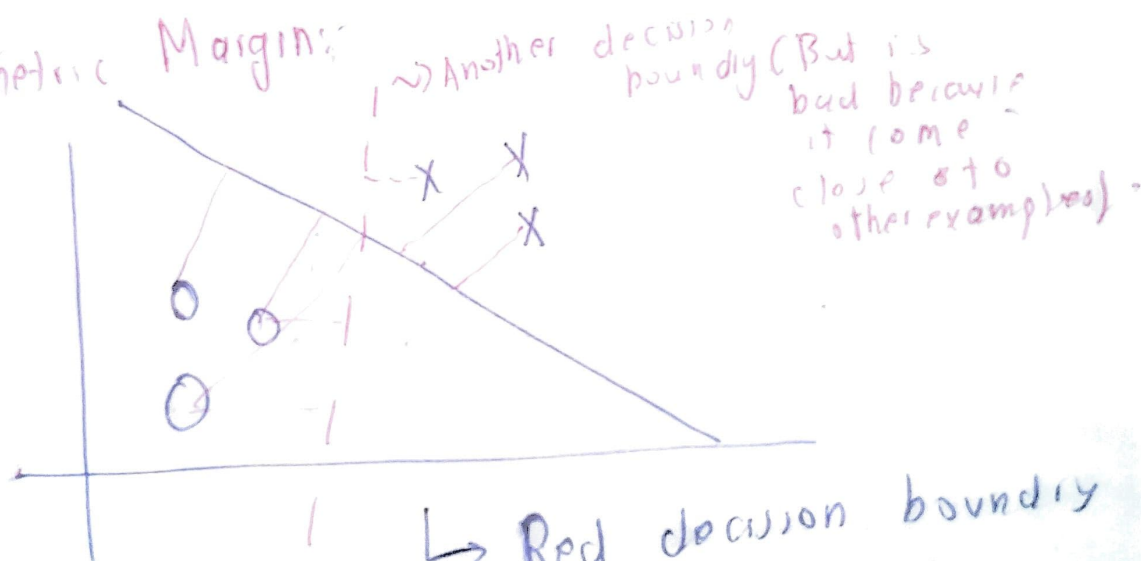$$h_\theta(x) = g(\theta^T x)$$

"1" if $\theta^T x > 0$

"0" Otherwise

ⓐ If $y^{(i)} = 1$, hope that $\theta^T x^{(i)} >> 0$
↓
Much greater

ⓑ if $y^{(i)} = 0$, hope that $\theta^T x^{(i)} << 0$
↓
Much less

Ⓧ Geometric Margin: ~ Another decision boundry (But is bad because it come close to other examples).

→ Red decision boundry has a lower geometric

→ Blue decision boundry has a higher decision boundry (Better!!!)

Intro - Some changes in SVM
Notation:

→ Labels → $y \in \{-1, +1\}$ → Not 0,1

Have h output values

in $\{-1, +1\}$ → Not probability

$$\rightarrow g(z) = \begin{cases} 1 & \text{if } z \geq 0 \\ \\ -1 & \text{other w.i.s.} \end{cases}$$

o Previous

$$h_\theta(x) = g(\theta^T x)$$

$\nwarrow \mathbb{R}^{n+1}, \; x_0 = 1.$

for
SVM
↓

Please
not
the differences

$$h_{w,b}(x) = g(w^T x + b)$$

$\mathbb{R}^n \quad \mathbb{R}$

$$\sum_{j=1}^{n} w_j x_j + b$$

( 

**⊛ Functional margin of (w,b) →hyperplane**

ⓡ Hyperplane defined by (w,b) w.r.t $(x^{(i)}, y^{(i)})$

$$\hat{\gamma}^{(i)} = y^{(i)}(w^T x^{(i)} + b)$$

If $y^{(i)} = 1$, want $w^T x^{(i)} + b > > 0$

If $y^{(i)} = -1$, want $w^T x^{(i)} + b < < 0$

⟶ $\hat{\gamma}^{(i)} > > 0$ (Combining these two)

(Want such a functional margin).

ⓢ if $\hat{\gamma}^{(i)} > 0$

that $h(x^{(i)}) = y^{(i)}$

→ Functional Margin w.r.t traing set:

$$\hat{\gamma} = \min_{i=1\dots m} \hat{\gamma}^{(i)}$$
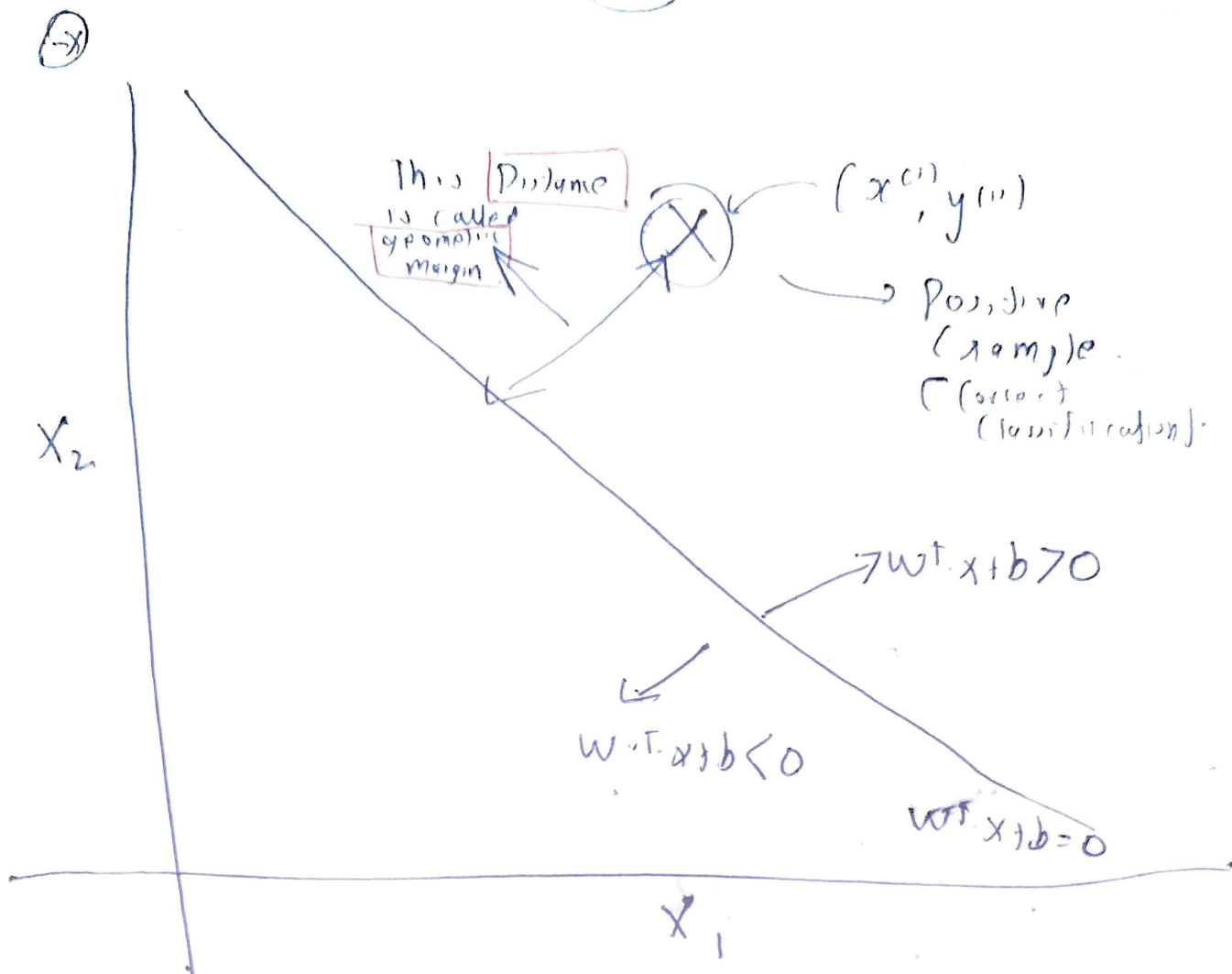
(Worst -Case Notion)

You can cheat this defination of functional margin by increasing the w w &b magnitude indefinately (Scaling up) $\rightarrow$ Would not change the decision boundry.

$\quad\llcorner\rightarrow$ You could normalize the length of your parameters.

$\quad\quad\llcorner\rightarrow$ Replace $(w,b) \rightarrow \left( \dfrac{w}{\|w\|}, \dfrac{b}{\|b\|} \right)$ (Normalization)

$\quad\quad\quad\Downarrow$

Boundry remain but cheating involving making the functional margin very high is resolved.

$\circledast$ Geometric margin

This [Distance] is called geometric margin

$(x^{(i)}, y^{(i)})$

→ Positive (sample) (correct classification).

→ $w^T x + b > 0$

$w^T x + b < 0$

$w^T x + b = 0$

$X_2$

$X_1$

↝ Formally Geometric margin of hyperplane $(w, b)$ w.r.t $(x^{(i)}, y^{(i)})$.

there is a proof to this!!

$$\gamma^{(i)} = \frac{y^i (w^T x^{(i)} + b)}{\| w \|}$$

↝ Relationship between Functional & geometric margin

$$\gamma^{(i)} = \frac{\hat{\gamma}^{(i)}}{\| w \|}$$

→ Functional margin

Geometric margin ←

(*) Geometric w.r.t training set

$$\gamma = \min_i \gamma^{(i)}$$

$$\hat{\gamma} = \text{functional margin}$$

$$\gamma = \text{geometric margin}$$

(*) ⊙ Optimal margin classifier

• (Choose w,b to maximize $\gamma$)
  (Geometric margin).



$\gamma$  $\gamma$

→ Maximize
   the distance.

How? (Mathematically)

$$\max_{\gamma, w, b} \gamma$$

$$\text{s.t} \quad \frac{y^{(i)}(w^T x^{(i)} + b)}{\|w\|} \geqslant \gamma$$

$i = 1 \ldots m$

Ⓧ this form is not a convex optimization problem therefore not solvable.

Ⓧ However it can be reformed:

$$\min_{w,b} \|W\|^2$$

$$s.t \quad y^{(i)}\left(w^T \cdot x^{(i)} + b\right) \geq 1$$

→ Same as before but rewritten and problem is convex optimization in this case.