IS A CLASSIFICATION ALGORITHM

- $y \in \{-1, +1\}$
- $\Theta = (\Theta_0, \Theta_1, ..., \Theta_d) \in \mathbb{R}^{d+1}$
  - $b \in \mathbb{R}$    $\vec{w} \in \mathbb{R}^d$

SUPPOSE WE HAVE A DATASET



**HYPOTHESIS CLASS**
$$h_{w,b}(x) = sign(w^T x + b)$$

**HYPERPLANE**
$$r: w^T x + b = 0, \quad w \perp r \longrightarrow \begin{cases} (w^T x^{(i)} + b) \geq 0 & y = 1 \quad \bullet \\ (w^T x^{(i)} + b) < 0 & y = -1 \quad \times \end{cases}$$

**MARGIN**

WE CAN SEE IT AS THE DISTANCE BETWEEN $r$ AND A DATA POINT. WE WANT THE LARGEST VALUE POSSIBLE OF THE MARGIN.

→ **FUNCTIONAL**    $\hat{\gamma}^{(i)} = y^{(i)} \left[ w^T x^{(i)} + b \right] \rightarrow \hat{\gamma} = \min_{i=1,...,n} \hat{\gamma}^{(i)}$

PROBLEM, IF I WANT TO SCALE $w$ AND $b$ BY $t$, THEN $\hat{\gamma}^{(i)}$ GETS $t$-UPLED (NOT GOOD) (WHILE $h_{w,b}(x)$ REMAINS THE SAME). WE NEED THE GEOMETRIC

→ **GEOMETRIC**    $\gamma^{(i)} = y^{(i)} \left[ \frac{w^T x^{(i)}}{\|w\|} + \frac{b}{\|w\|} \right] \rightarrow \gamma = \min_{i=1,...,n} \gamma^{(i)}$    (IF $t$-SCALES, GETS $t$-CANCELED)

IN ORDER TO FIND THE OPTIMAL CLASSIFIER WE NEED TO APPLY A CONVEX OPTIMIZATION PROBLEM

OPERATION OP OP.

$$\max$$

OBJECTIVE

$$\gamma$$

VARIABLES W.R.T. OPTIMIZE → $\gamma, w, b$

$$\text{s.t.} \quad y^{(i)}\left(w^T x^{(i)} + b\right) \geqslant \gamma, \quad i = 1, \ldots, n$$

CONSTRAINTS

$$\|w\| = 1$$

⇓

**PRIMAL OPTIMIZATION PROBLEM**

$$\min_{w, b} \quad \frac{1}{2} \|w\|^2$$

$$\text{s.t.} \quad y^{(i)}\left(w^T x^{(i)} + b\right) \geqslant 1, \quad i = 1, \ldots, n$$

⇓

**DUAL OPTIMIZATION PROBLEM**
W/ LAGRANGE MULTIPLIERS

$$\max_{\alpha} \quad W(\alpha) = \sum_{i=1}^{n} \alpha_i - \frac{1}{2} \sum_{i=1}^{n} y^{(i)} y^{(j)} \alpha_i \alpha_j \langle x^{(i)}, x^{(j)} \rangle$$

$$\text{s.t.} \quad \alpha_i \geqslant 0, \quad i = 1, \ldots, n$$

$$\sum_{i=1}^{n} \alpha_i y^{(i)} = 0$$

SOLVING IT

α  —  OUTPUT FROM SOLUTION

WE CAN NOW MAKE A PREDICTION

$$\left( w^T x + b \right) = \left( \sum_{i=1}^{n} \alpha_i y^{(i)} x^{(i)} \right)^T x + b = \sum_{i=1}^{n} \alpha_i y^{(i)} \underbrace{\langle x^{(i)}, x \rangle}_{\text{KERNELIZED} \atop K(x^{(i)}, x^{(j)})} + b$$

KERNELIZED

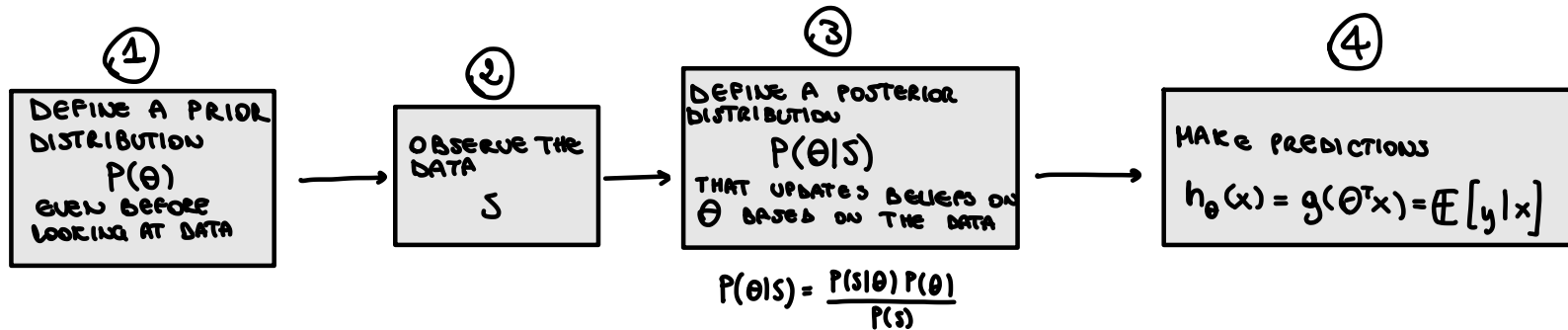$K(x^{(i)}, x^{(j)})$

<u>SUPPORT VECTORS α's</u>

α  ARE  MOSTLY  $\vec{0}$, THE MORE  THEY ARE DEFINED, THE  MORE  THE IMPACT ON THE LINEAR CLASSIFIER (SMALL MARGIN)

# BAYESIAN METHODS

- SO FAR WE WORKED W/ **FREQUENTIST APPROACH** → MLE, $\theta$ IS CONSTANT VALUED AND UNKNOWN

- **BAYESIAN APPROACH** → $\hat{\theta}$ $\underset{\theta}{\arg\max}$ $L(\theta, S)$     → $\theta$ IS A <u>RANDOM VARIABLE</u> AND UNKNOWN

  → $S = \left\{ x^{(i)}, y^{(i)} \right\}_{i=1}^{n}$ (TRAINING SET)

①
DEFINE A PRIOR DISTRIBUTION $P(\theta)$ GIVEN BEFORE LOOKING AT DATA

→

②
OBSERVE THE DATA $S$

→

③
DEFINE A POSTERIOR DISTRIBUTION $P(\theta|S)$ THAT UPDATES BELIEFS ON $\theta$ BASED ON THE DATA

$P(\theta|S) = \dfrac{P(S|\theta)\,P(\theta)}{P(S)}$

→

④
MAKE PREDICTIONS $h_\theta(x) = g(\theta^T x) = \mathbb{E}[y|x]$

POSTERIOR PREDICTIVE DISTRIBUTION (DISTRIBUTION, NOT PREDICTION)

AVERAGING         POSTERIOR DISTRIBUTION

$$P(y_*|x_*, S) = \int_\theta P(y_*|x_*, \theta)\, P(\theta|S)\, d\theta$$

WITH $y_*, x_*$ BEING UNKNOWN

$$= \mathbb{E}_{\theta \sim P(\theta|S)}\left[ P(y_*|x_*, \theta) \right]$$

WE CONSIDER EVERY VALUE OF $\theta$ GIVEN FROM THE DISTRIBUTIONS AS LINEAR COMBO OF THE ONES GIVEN IN THE POSTERIOR DISTRIBUT.

$$P(y_* \mid x_*, S) = \int_\theta P(y_*, \theta \mid x_*, S)\, d\theta$$

LET'S NOW APPLY THE CHAIN RULE

$$= \int_\theta P(\theta \mid x_*, S)\, P(y_* \mid \theta, x_*, S) = \int_\theta P(\theta \mid S)\, P(y_* \mid \theta, x_*)\, d\theta = \int_\theta P(y_* \mid x_*, \theta)\, P(\theta \mid S)\, d\theta$$

THE UNKNOWN $\theta$ IS INDEPENDENT FROM TEST EXAMPLES $X$ ||

## PROBLEMS

1. HOW DO WE CHOOSE A PRIOR?

2. HOW DO WE GO FROM PRIOR TO POSTERIOR?

$P(\theta)$ — PRIOR DISTRIBUTION

$P(S \mid \theta)$ — LIKELIHOOD (THIS DISTRIBUTION DEFINES OUR MODEL)

WE CAN GET TO THE POSTERIOR THANKS TO BAYES

$$P(\theta \mid S) = \frac{P(\theta)\, P(S \mid \theta)}{P(S)}$$

$$P(\theta \mid S) = \frac{P(\theta)\, P(S \mid \theta)}{\int_\theta P(\theta) P(S \mid \theta)\, d\theta}$$

EX MONTE CARLO APPROACH

$$\theta^{(0)}, \theta^{(1)}, \ldots, \theta^{(k)}$$

GIBBS SAMPLING : EXTRACTION OF FEATURES

$$\frac{1}{N} \sum_{i=1}^{n} P(y_* \mid x_*, \theta^{(i)}) \approx \mathbb{E}_{\theta \sim P(\theta \mid S)} \left[ P(y_* \mid x_*, \theta) \right]$$

USUALLY WE DON'T HAVE CLOSED-FORM EXPRESSIONS, BUT APPROXIMATIONS OF THE POSTERIOR DISTR. FUNCTION