

03/23/15

SUPPORT VECTOR MACHINES (SVM)

- ① Invented by Vapnik & Chervonovnik in '63.
- "Modern" form by Cortes & Vapnik in '93 & '95.

→ Intuition: Don't solve a harder problem as an intermediate step towards your goal.

If you care about low-loss (misclassification rate), then optimize for that.

→ Goal: Classification

① Generative approach: estimate $P(X|Y)$ \Rightarrow Bayes Rule $\Rightarrow P(Y|\bar{X}) \Rightarrow \hat{y}_{MAP}$
(NB)

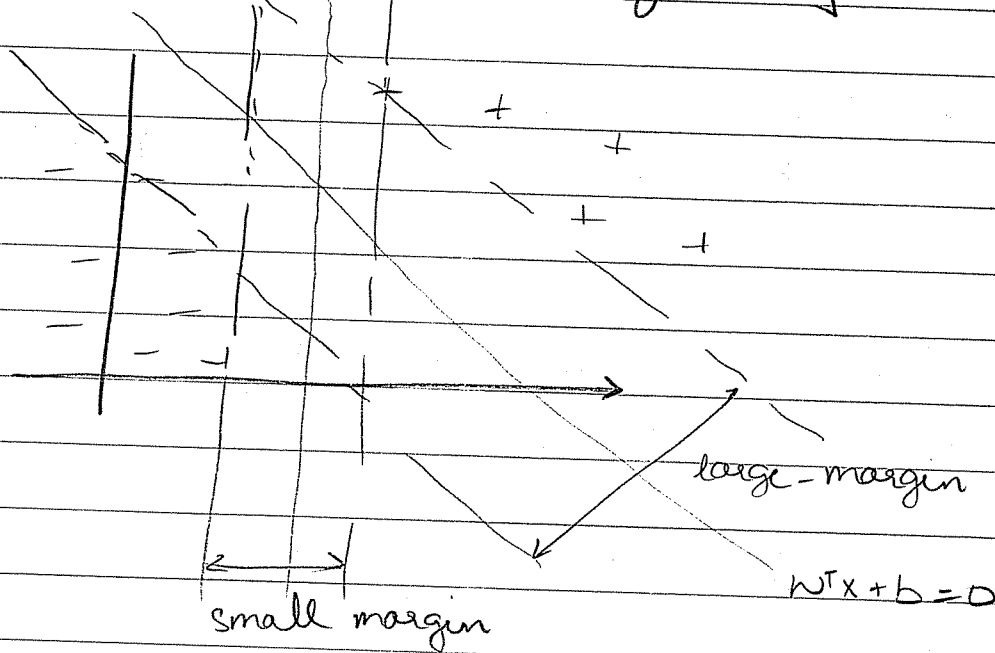
② Discriminative approach: estimate $P(Y|\bar{X}) \Rightarrow \hat{y}_{MAP} = \operatorname{argmax} P(Y|\bar{X})$
(LR)

③ Discriminative #2: estimate $\text{Score}(x) = w^T x$ (or $w^T x + b$)
(SVM)
 $\hat{y} = \operatorname{sign}(w^T x + b)$

← Notation Change: $\rightarrow y \in \{+1, -1\}$ (rather than $\{0, 1\}$)

→ \bar{w} no longer has bias weight included
So $\bar{w}^T \bar{x} + b \geq 0$ rather than $\bar{w}^T \bar{x} \geq 0$
 $\bar{w} \in \mathbb{R}^d$ $\bar{w} \in \mathbb{R}^{d+1}$
 \bar{x}

② SVM [Linear SVM for now]



→ Many different (\vec{w}, b) work well (actually perfectly) on training data. Which one should we choose?

→ Ans: The one with the largest margin!

Why?

- Trust me, I'm Vapnik
- Math will be easy & elegant.
- There will be bounds on generalization!
- Many different interpretations lead to this; Has to be correct!

Mathematically
(or pseudo-mathematically)
for now

max
 \vec{w}, b

Margin

s.t

Correct Classification
on training data

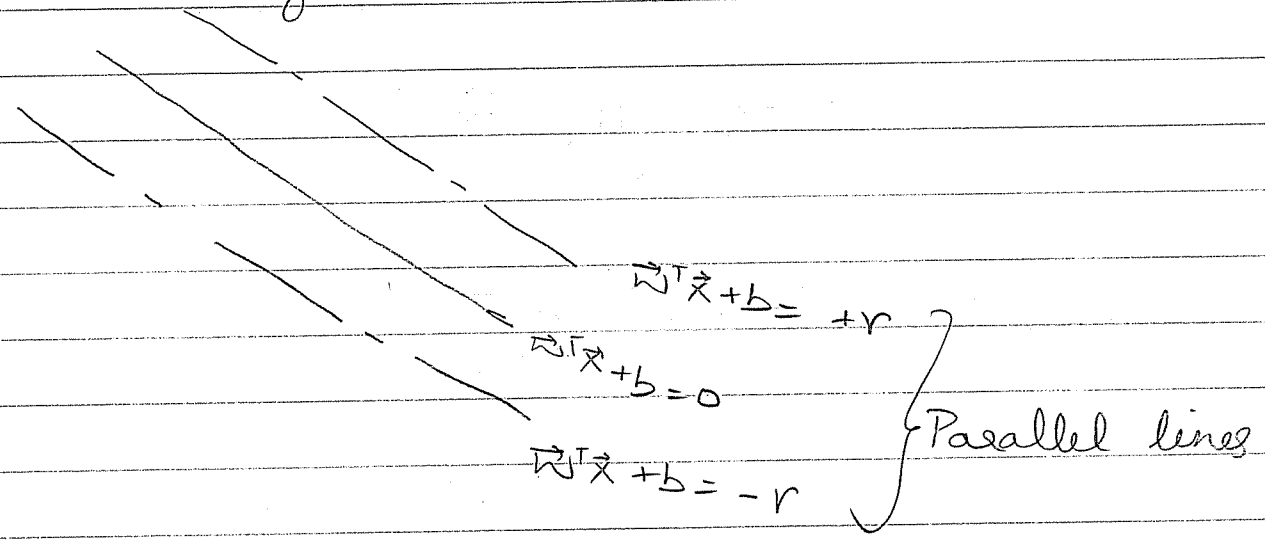
What is correct classification?

$$\Rightarrow \left\{ \begin{array}{l} \vec{w}^T \vec{x}_i + b \geq 0 \quad \forall y_i = +1 \\ \vec{w}^T \vec{x}_i + b < 0 \quad \forall y_i = -1 \end{array} \right\}$$

together

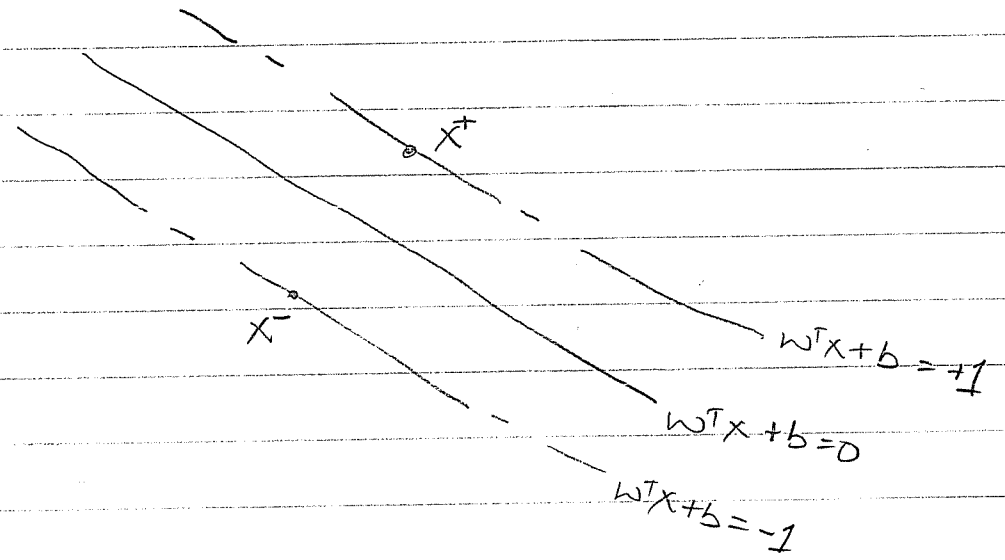
$$y_i (\vec{w}^T \vec{x}_i + b) \geq 0 \quad \forall i$$

What is Margin?



→ Notice that we can always scale both w & b to get arbitrary $\pm r$. So lets fix $r = +1$

Now



Let x^+ be a pt on $w^T x + b = +1$

" x^- be a pt on $w^T x + b = -1$

Let line from x^+ to x^- be orthogonal to our hyperplane $w^T x + b = 0$

$$\Rightarrow \|x^+ - x^-\| \equiv \text{Margin}$$

We know $(x^+ - x^-) \perp$ hyperplane

$$\Rightarrow (x^+ - x^-) \parallel \vec{w}$$

parallel

$$\Rightarrow (x^+ - x^-) = \lambda \vec{w}$$

scalar

Let's left & right multiply by \vec{w}^T

$$w^T (x^+ - x^-) = \lambda w^T w$$

$$\underbrace{w^T x^+}_{\downarrow} - \underbrace{w^T x^-}_{\downarrow} = \lambda w^T w$$

$$(1-b) - (-1-b) = \lambda w^T w$$

$$\Rightarrow \lambda = \frac{2}{w^T w}$$

$$\Rightarrow \text{Margin} = \|x^+ - x^-\|$$

$$= \frac{2}{w^T w} \|\vec{w}\| = \frac{2}{\|\vec{w}\|}$$

3

Finally, Linear SVM with a "hard margin"

$$\max_{\vec{w}, b} \frac{2}{\|w\|} \quad \left. \vphantom{\max} \right\} \text{Not very nice to optimize}$$

$$\text{s.t. } y_i (\vec{w}^T \vec{x}_i + b) \geq 0$$

\approx changes

$$\left\{ \begin{array}{l} \min_{\vec{w}, b} \frac{1}{2} \|w\|^2 \equiv \frac{1}{2} w^T w \\ \text{s.t. } y_i (\vec{w}^T \vec{x}_i + b) \geq (+1) \quad \forall i \end{array} \right. \quad \text{why?}$$

Very well studied optimization problem

→ Quadratic Program (QP)

→ $w^T w \equiv$ quadratic in \vec{w} & b

→ Convex

$$\rightarrow \text{Hessian} = \frac{1}{2} \begin{bmatrix} I_{d+1} & 0 \\ 0 & 0 \end{bmatrix} \succeq 0 \quad (\text{PSD})$$

→ $d+1$ variables, N constraints

→ Standard QP solvers will solve this optimally for you.

③ Soft-Margin Linear SVM

→ What if data is not linearly separable?

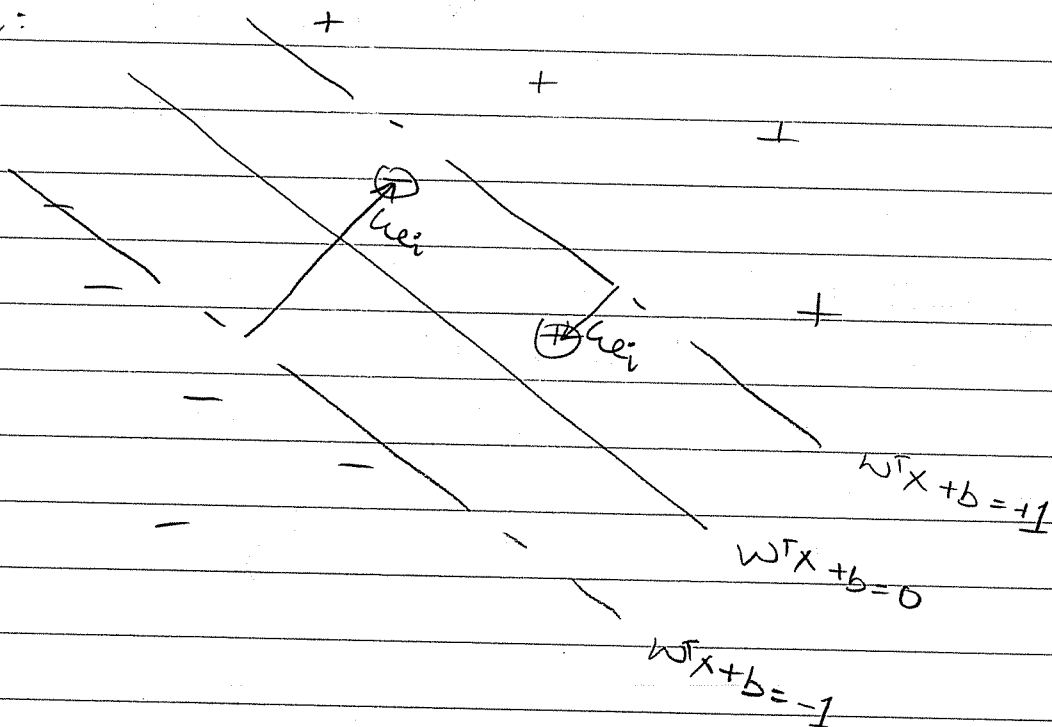
→ QP will be infeasible \Rightarrow no solution

→ How can we allow for some mistakes?

→ Idea: allow some "slack"

allow every data-point to violate constraint by some "slack" ϵ_i

Picture:



$$\min_{\vec{w}, b, \epsilon_i} \frac{1}{2} \vec{w}^T \vec{w} + C \sum_{i=1}^N \epsilon_i$$

penalize the violations / slacks

$$\vec{w}^T \vec{x}_i + b \geq 1 - \epsilon_i \quad \forall i$$

\rightarrow violation in margin

$$\epsilon_i \geq 0$$

\rightarrow can't set slack to negative to game the objective

4

→ Still a convex QP

→ $d+1 + N$ variables

\vec{w} b $\{\epsilon_i\}$

→ $2N$ constraints

$C \equiv$ trade-off parameter

$C=0 \Rightarrow w^* = \vec{0} \quad b^* = 0$ (Opt soln)

$C = \infty \Rightarrow$ Hard-Margin SVM
(might be impossible for linearly non-separable data)



④ Hinge-Loss

Hinge Function is defined as:

$h(z) = \max\{0, 1-z\}$

In SVM QP

$\left\{ \begin{array}{l} \epsilon_i \geq 1 - y_i(w^T x_i + b) \\ \epsilon_i \geq 0 \end{array} \right\}$

$\Rightarrow \epsilon_i \geq \max\{0, 1 - y_i(w^T x_i + b)\}$

Almost like $h(y_i(w^T x_i + b))$ but inequality
let's fix that.

Let $s_i \equiv$ score by SVM of i th data-point
 $= w^T x_i + b$

$$\xi_i \geq \max\{0, 1 - y_i s_i\}$$

Claim: Let (w^*, b^*, ξ_i^*) be QP optimum,
 $\xi_i^* = \max\{0, 1 - y_i s_i^*\}$

Proof: By contradiction. Assume equality doesn't hold

$$\Rightarrow \xi_i^* > \max\{0, 1 - y_i (w_i^{*T} x_i + b^*)\}$$

Define $\xi_i^{\text{new}} = \max\{0, 1 - y_i (w_i^{*T} x_i + b^*)\}$

by hinge loss
construction

$$\left\{ \begin{array}{l} \xi_i^{\text{new}} \geq 0 \\ \xi_i^{\text{new}} \geq 1 - y_i (w_i^{*T} x_i + b^*) \end{array} \right\}$$

$\Rightarrow (w^*, b^*, \xi_i^{\text{new}})$ is a FEASIBLE soln
to SVM QP

Also $\{w^*, b^*, \xi_i^{\text{new}}\}$ reduces objective by
 $C \sum_{i=1}^n (\xi_i^* - \xi_i^{\text{new}}) > 0$

$\Rightarrow (w^*, b^*, \xi_i^*)$ is NOT OPT.

\Rightarrow contradiction

\Rightarrow QED

$$\text{So } \xi_{ei}^* = \max \{0, 1 - y_i (\tilde{w}^T x_i + b^*)\}$$

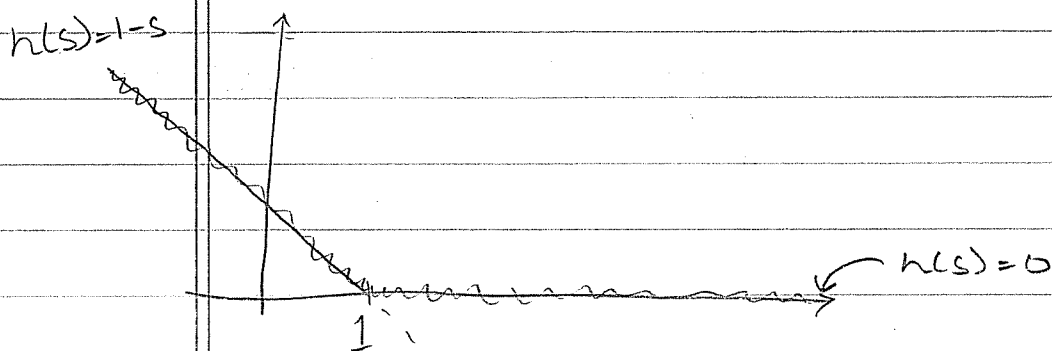
⇒ SVM QP can be re-written as

Very cool!
Beginning to
look like
things we
recognize
Loss + Regularization

$$\min_{w, b} \left[\frac{1}{2} w^T w + C \sum_{i=1}^N h(y_i (\tilde{w}^T x_i + b)) \right]$$

Unconstrained QP where

$$h(z) = \max \{0, 1 - z\} \text{ - Hinge-Loss!}$$



⑤ Comparison between LR & SVM

→ SVM: $\min_w \text{hinge-loss} + \frac{1}{2} w^T w$

→ LR: $\max_w \log P(y|x, w) - \lambda w^T w$

Can we compare the first terms?

→ Consider just 1 training pt ($\tilde{x}_1, y_1 = +1$)

[Datasets are never this small, but just for illustration]

LR

$$P(Y=1 | X, w) = \frac{1}{1 + e^{-w^T X + b}}$$

$$\vec{w}, b \sim N(\vec{0}, \sigma^2 I)$$

$$\vec{w}_{MAP}, b_{MAP} = \underset{w, b}{\operatorname{argmax}} \log P(Y=1 | \vec{x}, w, b) + \log P(w, b)$$

$$= \underset{w, b}{\operatorname{argmax}} \log \left[\frac{1}{1 + e^{-w^T x + b}} \right] - \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(w^T x + b)^2}{2\sigma^2}}$$

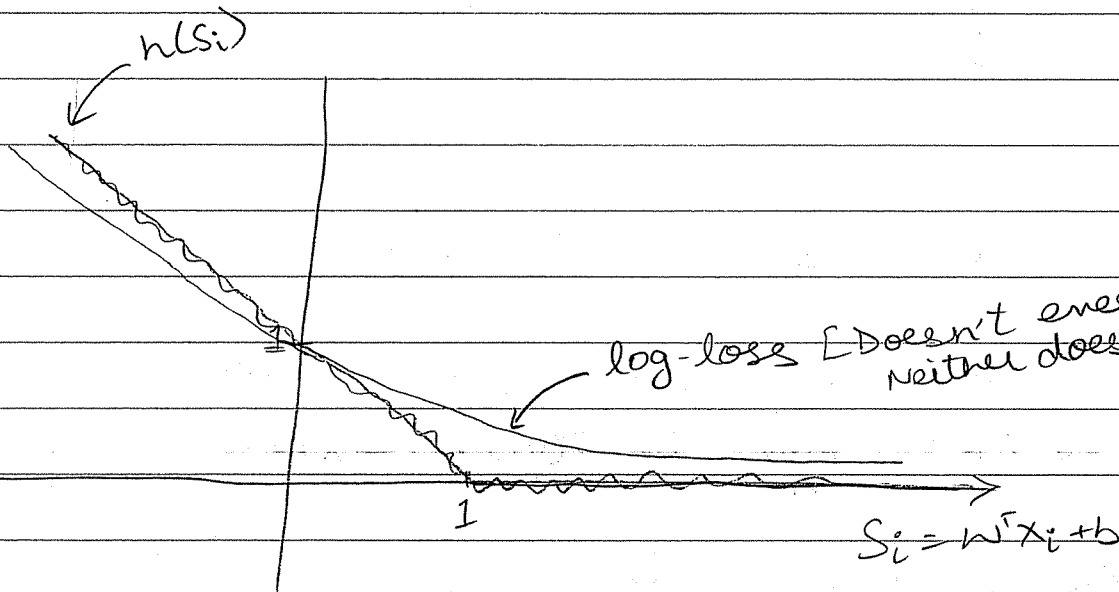
$$= \underset{w, b}{\operatorname{argmin}} \frac{1}{2} w^T w + \lambda \underbrace{\log(1 + e^{-w^T x + b})}_{\text{log-loss}}$$

SVM

max Margin
s.t. correct Class

=

$$\underset{w, b}{\operatorname{min}} \frac{1}{2} w^T w + \text{hinge}(w^T x + b)$$



become zero]
neither does the gradient.