

NNs

representation learning
 excellent perf
 non-convex / costly
 hard to interpret
 an "art"

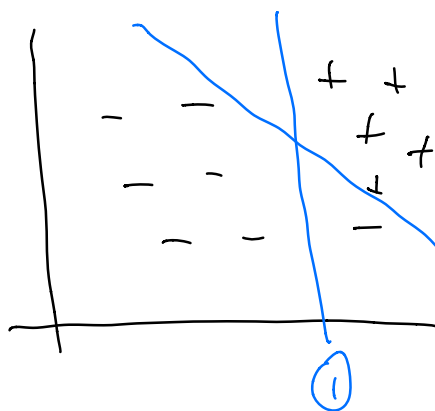
SVMs

basis engineering
 very good performance
 convex / easy to train
 simple, interpretable
 coherent theory

Setting Binary classification discr.

$$\hat{y} = \begin{cases} +1 & \text{if } \{ \underline{w}^T \underline{x} + w_0 \} > 0 \\ -1 & \text{otherwise} \end{cases}$$

For now assume separable data

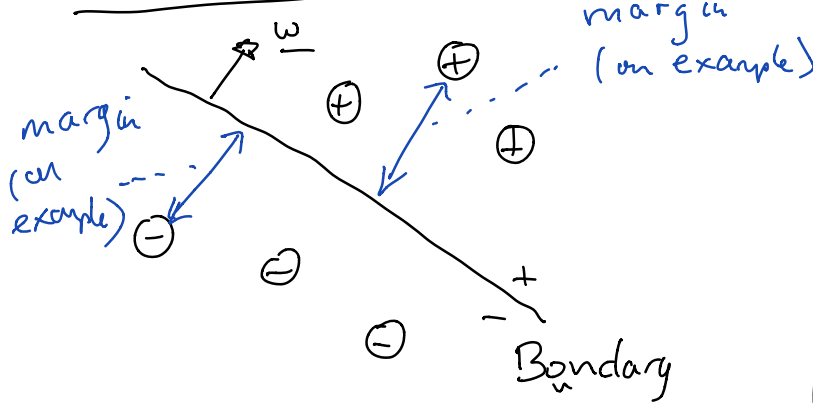


Prefer ① to ②

↳ generalize better

② (small perturbation to data will not matter)

Max margin

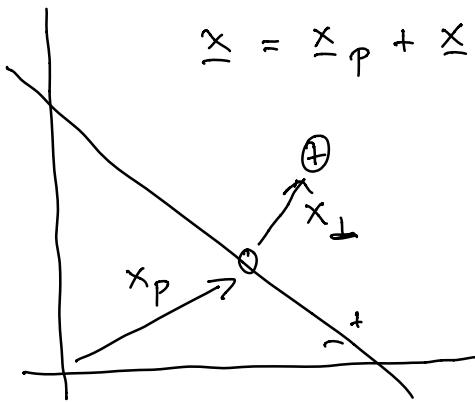


margin on correctly classified example is the absolute, normalized, orthogonal distance to boundary

"margin on data" = min margin on correct examples

Goal: find separator that maximizes margin

Geometry



$$\underline{x} = \underline{x}_p + \underline{x}_{\perp} = \underline{x}_p + r \frac{\underline{w}}{\|\underline{w}\|_2}$$

$$\underline{w}^T \underline{x} = \underline{w}^T \underline{x}_p + r \frac{\underline{w}^T \underline{w}}{\|\underline{w}\|_2}$$

$$= -w_0 + r \|\underline{w}\|_2$$

Boundary
 $\underline{w}^T \underline{x} + w_0 = 0$

$$r = \frac{\underline{w}^T \underline{x} + w_0}{\|\underline{w}\|_2}$$

Generally (pos ex. + neg ex.):

$$\text{margin}(x_n, y_n) = y_n \left(\frac{\underline{w}^T \underline{x}_n + w_0}{\|\underline{w}\|_2} \right) \quad (> 0) \quad (*)$$

Note ① Margin is invariant to multiplying (\underline{w}, w_0) by scalar $\beta > 0$

[But note that $y_n(\underline{w}^T x_n + w_0)$ increases!]

② (*) is negative if used as a misclassified example

Hard max-margin formulation

$$\textcircled{1} \max_{\underline{w}, w_0} \left[\min_n y_n \left(\frac{\underline{w}^T x_n + w_0}{\|\underline{w}\|_2} \right) \right]$$

↳ will find a separator!

↳ looks "ugly"

② By invariance to scaling by $\beta > 0$, w.l.o.g. to impose $y_n(\underline{w}^T x_n + w_0) \geq 1$

$$\max_{\underline{w}, w_0} \frac{1}{\|\underline{w}\|_2} \min_n y_n(\underline{w}^T x_n + w_0)$$

$$\text{s.t. } y_n(\underline{w}^T x_n + w_0) \geq 1 \text{ for all } n$$

Equivalent { one or more constraints will bind }

$$\max_{\underline{w}, w_0} \frac{1}{\|\underline{w}\|_2}$$

$$\text{s.t. } y_n(\underline{w}^T \underline{x}_n + w_0) \geq 1 \text{ for all } n$$

Hard-margin formulation

$$\min_{\underline{w}, w_0} \|\underline{w}\|_2^2$$

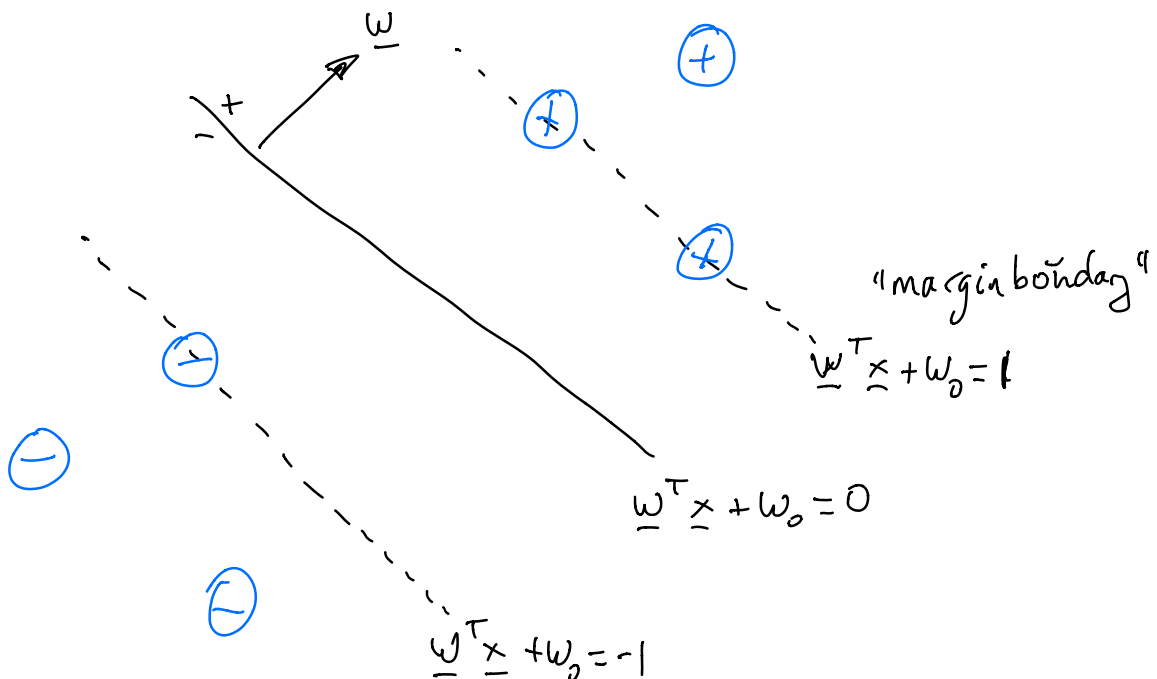
$$\text{s.t. } y_n(\underline{w}^T \underline{x}_n + w_0) \geq 1, \text{ all } n$$

[Note, can also write $\min \frac{1}{2} \|\underline{w}\|_2^2$!]

Nice! Convex (quadratic objective, linear constraints)

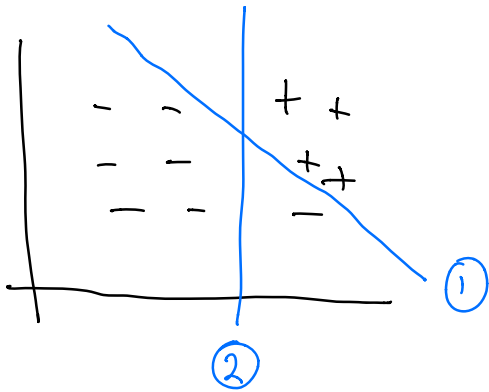
$$\left[\text{margin} = \frac{1}{\|\underline{w}\|_2} \right]$$

\cup



Soft-margin formulation

- Regularization
- Non-separable data

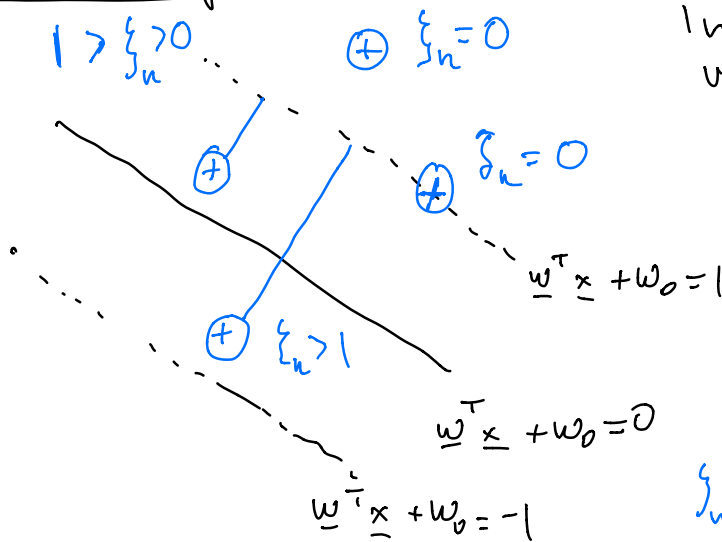


May prefer ② over ①, even though ② is not separating

(and data may not be separable!)

Relaxed formulation

(x_i)



Introduce new variable $\xi_n \geq 0$, each n

"How much is x_n on the (wrong side) of the margin boundary"

- $\xi_n = 0$ correct class
- $0 < \xi_n < 1$ correct class but smaller margin than $1/\|\underline{w}\|_2$
- $\xi_n > 1$ incorrect class

Soft-margin formulation

For some $C > 0$ (regularization parameter)

$$\min_{\underline{w}, w_0, \xi} \frac{1}{2} \underline{w}^T \underline{w} + C \sum_n \xi_n$$

s.t. $y_n (\underline{w}^T \underline{x}_n + w_0) \geq 1 - \xi_n$, all n
 $\xi_n \geq 0$, all n

⊠ "Pretend" margin $\frac{1}{\|\underline{w}\|_2}$ (ignoring examples live "in the margin")

⊠ Allows misclassified points

⊠ Large C ($C = \infty$ is hard margin)

↳ less regularization

(try to get correct classification all all n)

⊠ Smaller C (larger)

↳ Better $\frac{1}{\|\underline{w}\|_2}$ ("margin")

more mistakes

Equivalently

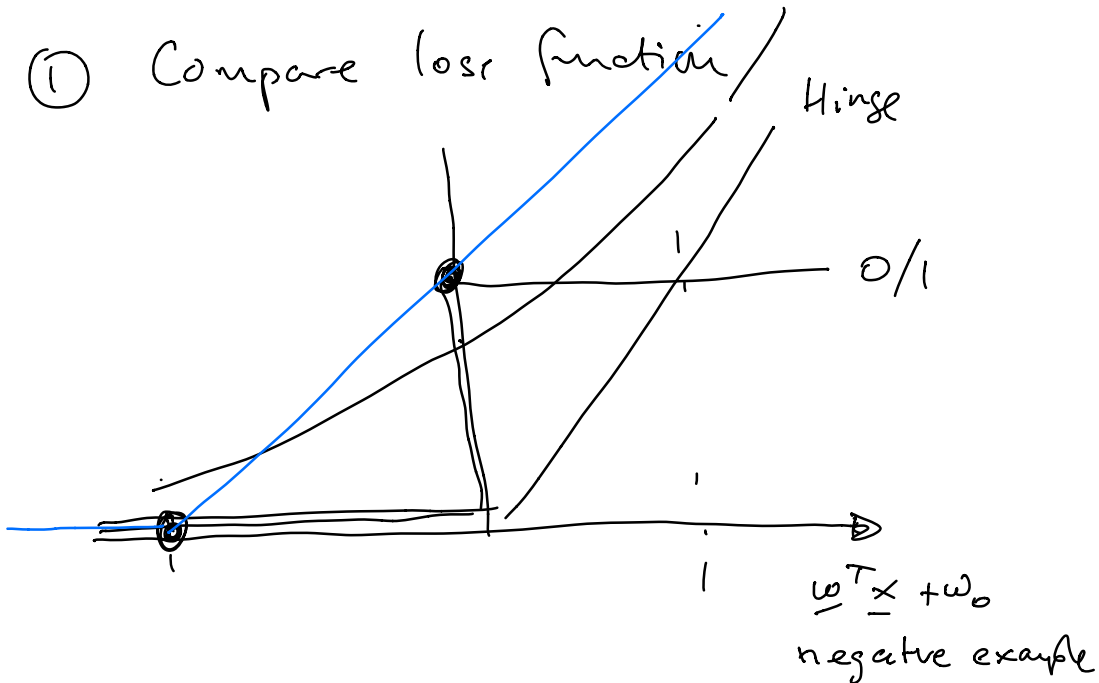
$$\min_{\underline{w}, w_0} \frac{1}{2} \underline{w}^T \underline{w} + C \sum_n \max(0, 1 - y_n (\underline{w}^T \underline{x}_n + w_0))$$

Convex! \cup + \leftarrow \nearrow \Rightarrow SGD
 + diff (almost everywhere)

(Notes)

SVM Logistic

① Compare loss function / Hinge



Why is it better than logistic?

② Quadratic + linear constraint
↳ Duality math; works with
basis functions nicely

② Large margins help to generalise
(think about small perturbation)