

Regression & Classification

Losses: MSE, L1

0/1, hinge, logistic reg.

Models:

linear models,
basis models,
neural networks

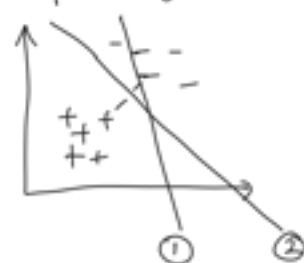
Idea: choose loss, model that make sense for the problem,
but still computationally tractable.

Broader theme: we want models that allow us to generalize.
losses

Today: loss that (a) has some nice computational properties when
paired with a broad class of models
(b) has excellent generalization

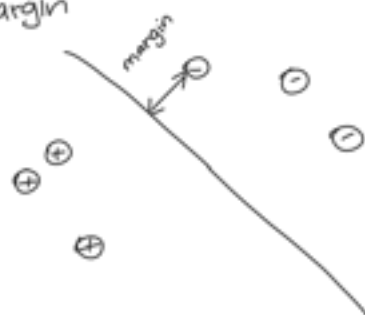
Setting: Binary classification: $\hat{y} = \begin{cases} 1 & \text{if } f(x, w) > 0 \\ -1 & \text{else} \end{cases}$

for now, assume separability; also linear case $f(x, w) = w^T x + w_0$



perceptron (hinge loss)
→ no preference

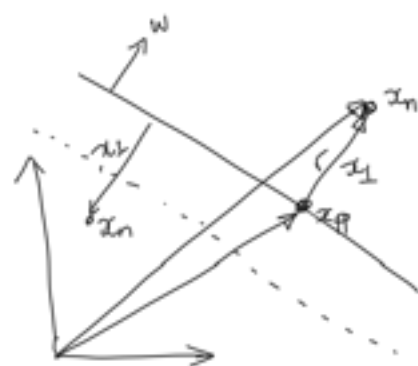
↳ max margin



how far does a pt. need
to move before its
classification changes?

goal: maximize the minimum
margin (later add bases...)

Some geometry:



recall: $w^T x + w_0 = 0$
defines the boundary

$w^T x + w_0 = c$ is parallel
to the boundary

$$x_n = x_{\parallel} + x_{\perp}$$

$$x_{\perp} = r \frac{w}{\|w\|_2}$$

margin unit vector in the
direction of w

$$x_{\parallel}: w^T x_{\parallel} + w_0 = 0$$

multiply my original eqn by w^T :

$$w^T x_n = w^T x_{\parallel} + w^T x_{\perp}$$

$$w^T x_n = -w_0 + r \frac{w^T w}{\|w\|_2} = -w_0 + r \frac{\|w\|_2^2}{\|w\|_2} = -w_0 + r \|w\|_2$$

solve for r

$$w^T x_n + w_0 = r \|w\|_2$$

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

signed
margin



unsigned (or just "margin") for some x_n : $y_n \left(\frac{w^T x_n + w_0}{\|w\|_2} \right)$

observe: correctly classified pts will have $\text{margin}(x_n) \geq 0$;
misclassified pts. will have negative margin.

margin(X)

$$\text{margin}(x_n) = y_n \left(\frac{w^T x_n + w_0}{\|w\|_2} \right)$$

$$\text{margin}(X) = \min_{x_n} \text{margin}(x_n)$$

now, we can write our objective:

$$\max_{w, w_0} \text{margin}(X) = \max_{w, w_0} \min_{x_n} y_n \left(\frac{w^T x_n + w_0}{\|w\|_2} \right)$$

anything misclassified \rightarrow we'll try to correct it!
even if correctly classified \rightarrow still try to make margin large!



great idea! looks UGLY!!

let's first notice an over-parametrization in this problem..

$$\begin{aligned} w^T x + w_0 &= 0 \text{ as boundary} \\ \beta w^T x + \beta w_0 &= 0 \\ \underbrace{\beta w^T}_{w'^T} x + \underbrace{\beta w_0}_{w'_0} &= 0 \\ w'^T x + w'_0 &= 0 \end{aligned}$$

$$\beta > 0$$

shows up in the margin also:

$$y_n \frac{1}{\beta \|w\|_2} (\beta w^T x + \beta w_0) \Leftrightarrow \text{same margin value!}$$

choose a β that will achieve this.

Let's decide on a scaling of w, w_0 s.t. $y_n (w^T x_n + w_0) \geq 1$

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

s.t. $y_n (w^T x_n + w_0) \geq 1$ my constraint

note: optimal sol'n will touch the constraint

$$\max_{w, w_0} \frac{1}{\|w\|_2} \text{ s.t. } y_n (w^T x_n + w_0) \geq 1$$

$$\min_{w, w_0} \|w\|_2 \text{ s.t. } y_n (w^T x_n + w_0) \geq 1$$

$$\min_{w, w_0} w^T w \text{ s.t. } y_n (w^T x_n + w_0) \geq 1$$

$$\begin{aligned} \min \|w\|_2 &= \min \|w\|_2^2 \\ &= \min w^T w \end{aligned}$$

very cool because this is a quadratic program with linear constraints.

→ convex!
→ very efficient solvers.

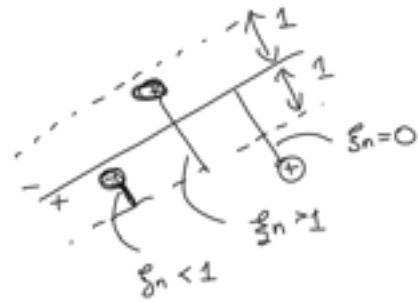
but what if ~~these~~ data are not separable?



in general, tension btw margin (X) size and correctly classified data/misclassification rate

defined only on correct data

move from "hard margin" to "soft margin"



$\xi_n \geq 0$ measure how much margin violation exists for pt. x_n



$$\min_{w, w_0} \frac{1}{2} \|w\|_2^2$$

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

$$\min_{w, w_0} \frac{1}{2} \|w\|_2^2 + C \sum \xi_n$$

$$\text{s.t. } y_n (w^T x_n + w_0) \geq 1 - \xi_n$$

$$\xi_n \geq 0$$

pretend "margin = 1, and allow some slack"

larger $C \rightarrow$ enforce constraints more, smaller margin/reg. less
small $C \rightarrow$ too many errors

usually, we set C w/ cross-validation

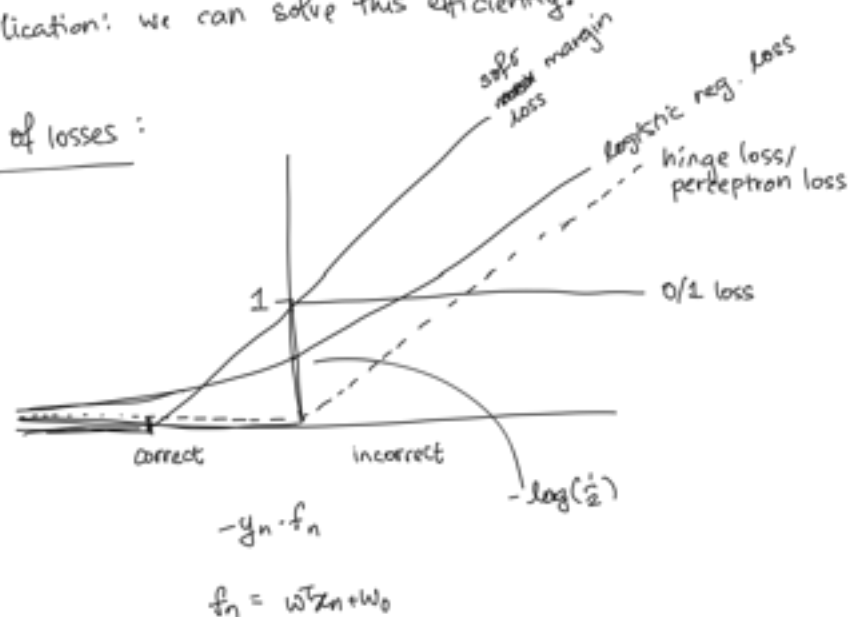
question: is this still convex? yes!

\Rightarrow sum of convex functions is still convex

$$\min_{w, w_0} \underbrace{\frac{1}{2} \|w\|_2^2}_{\text{quadratic}} + C \sum \underbrace{\max(0, 1 - y_n (w^T x_n + w_0))}_{\text{hinge loss}}$$

\Rightarrow implication: we can solve this efficiently!

Review of losses:



end of our supervised learning unit.

Concept Question: close to the end of ch 1

models: nonparametric
parametric: linear basis
NN
practicals: random forests

losses: \rightarrow dimensionality & neighbors

When should you choose each model/loss:

discuss: data? $\frac{\text{semi-sup vs. not}}{\text{linear vs. nonlin.}}$ \rightarrow size: computation, expressiveness

convex vs. non \rightarrow opt? \rightarrow access to knowledge? \rightarrow interpretability
basis vs. nn
type of fit: linear vs. nonlinear

Max margin