# Lecture 24: Boosting

## Loss Functions

### Regression:

① squared error loss

$$L(y, f(x)) = (y - f(x))^2$$

where $r = y - f(x)$

② absolute loss

$$L(y, f(x)) = |y - f(x)|$$

where $r = y - f(x)$

## Classification: (binary)

Two parameterizations for $Y$

① $Y \in \{0, 1\}$  $\xrightarrow{\frac{(Y+1)}{2}}$  ② $Y \in \{-1, 1\}$

$\xrightarrow{2Y-1}$

losses:  0-1 loss

$$L(y, f(x)) = \mathbb{I}(y \neq f(x)) = \begin{cases} 0 & y = f(x) \\ 1 & y \neq f(x) \end{cases}$$

If I use $-1/1$ encoding then

$$y \in \{-1, 1\}, \quad f(x) \in \{-1, 1\}$$

correct classification = signs of $y$, $f(x)$ matching

incorrect " = " don't match

Also note for any classifier $f$ there is some fn $h$ so that

$$f(x) = \text{sign}(h(x)) \qquad \text{← like a disc. fn } \delta$$

idea: $h(x) >> 0$ if class $1$
$h(x) << 0$ if class $-1$

example: linear classifier
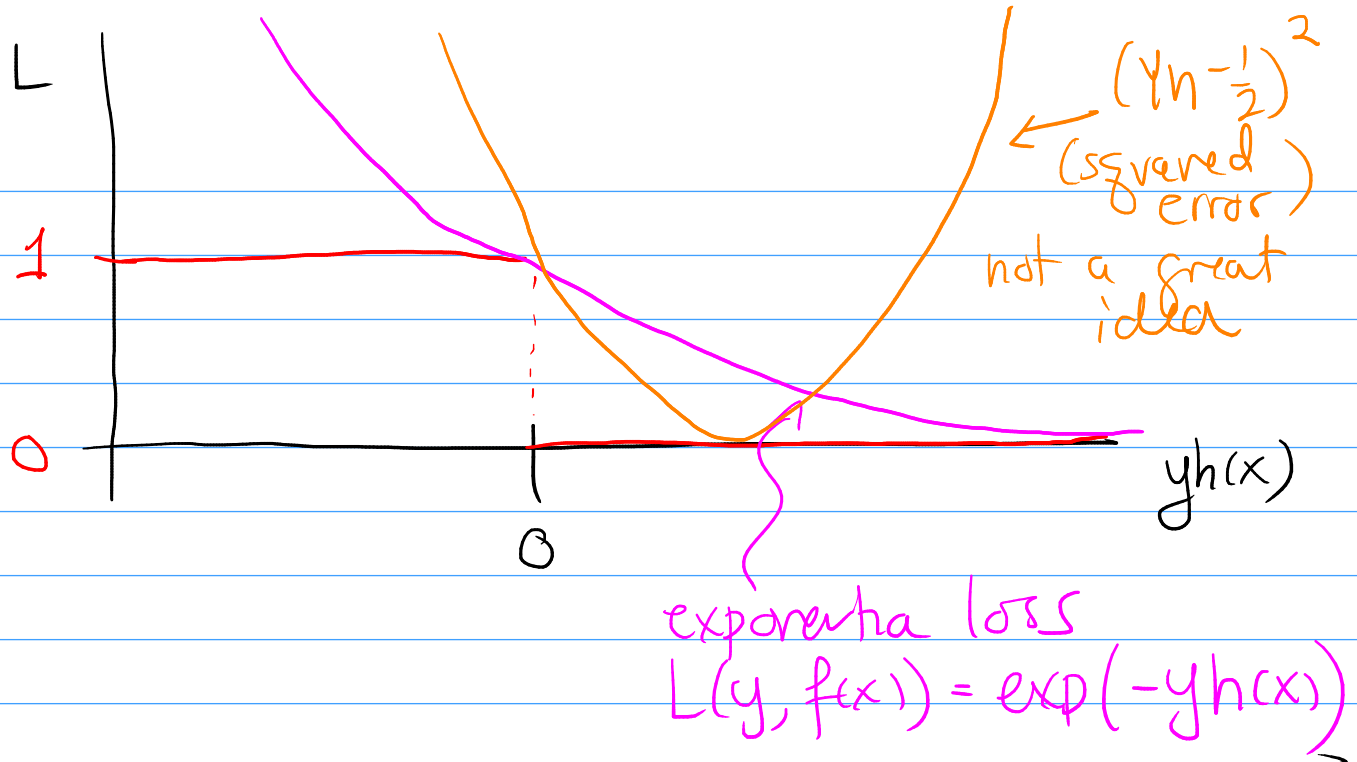
$$f(x) = \text{sign}(w^T x)$$

margin: $yh(x)$

$yh(x) > 0$  correct classification

$yh(x) < 0$  incorrect "

$|yh(x)|$ = amount of correct/incorrect

We can write $0-1$ loss as a fn of the margin

$$L(y, f(x)) = \mathbb{1}(yh(x) < 0)$$

$\left(yh - \frac{1}{2}\right)^2$
(squared error)

not a great idea

exponential loss
$L(y, f(x)) = \exp(-yh(x))$

ensemble

## Boosted Methods

Orig. motivated as a method to combine a series of weak classifiers into a stronger one.

Weak classifier: one that is not much better than random chance

→ typically a "stump" tree w/ 1 split

## Boosting:

① sequentially learn weak classifiers on a series of modified traing data

$\hat{f}_1, \hat{f}_2, \hat{f}_3, \ldots$

→ use a weighted loss

each classifier will be modified (weighted) to focus on traing data where prev. class. was bad.

②
$$\hat{f}(x) = \text{sign}\left(\sum_{b=1}^{B} \alpha_b \hat{f_b}(x)\right)$$

$h(x)$ — over the sum $\sum_{b=1}^{B} \alpha_b \hat{f_b}(x)$

$\hat{f}(x)$ is $\pm 1$

↳ weighted combn of individual $\hat{f_b}$

weight $\alpha_b$ is higher for better classifiers and lower for worse.

## Ada Boost
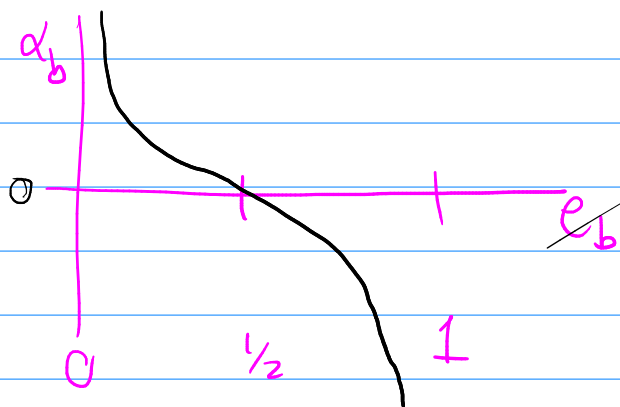
① $w_n = \frac{1}{N}$ = weight for $n^{th}$ traing pt.

② For $b = 1, \ldots, B$

ⓐ fit $\hat{f_b}$ using weighted loss $(w_n)$

ⓑ compute error for $\hat{f_b}$

$$e_b = \text{weighted mis. class. err.}$$

ⓒ $\alpha_b = \frac{1}{2} \log\left(\frac{1-e_b}{e_b}\right)$



ⓓ update $w_n$

$$w_n \longleftarrow \exp(\pm \alpha_b) w_n$$

If $\hat{f_b}$ got $n^{th}$ point correct $e^{-\alpha_b} w_n$

incorrect $e^{\alpha_b} w_n$

(2) $\hat{f}(x) = \text{Sign}\left(\sum_{b=1}^{B} \alpha_b \hat{f}_b(x)\right)$

[as prev.]

---

What is boosting doing?

## General form:

$$\hat{f}(x) = \text{Sign}(h(x))$$
$$h(x) = \sum_{b=1}^{B} \alpha_b \hat{f}_b(x)$$

## Generalized Additive method

$$h(x) = \sum_{b=1}^{B} \alpha_b \beta_b(x; \gamma_b)$$

$\gamma_b$ = params for basis fu.

$\beta_b$ is same collection of "basis" fns

want to find $\alpha_b, \gamma_b$ to fit data well

$$\underset{\{\alpha_b, \gamma_b\}_{b=1}^{B}}{\text{argmin}} \text{Loss}(h(x))$$

Lots of params! A problem.

Soln: Use a greedy <sup>forward</sup> "stagewise"
add. modeling approach where
do this for one pair $\alpha_b, \gamma_b$ at a
time.

For $b = 1, \ldots, B$

(a) $\alpha_b, \gamma_b = \underset{\alpha, \gamma}{\text{argmin}} \, \text{Loss}\left(\hat{f}_{b-1}(x) + \alpha \beta(x, \gamma)\right)$

(b) $f_b(x) = f_{b-1}(x) + \alpha_b \beta(x, \gamma_b)$

Punchline: Ada-Boost is basically
   Forward stagewise add. modeling
where $\beta$s are "stumps" and we
use an Exp. loss to measure error.