

## Loss Function

### Regression:

① Squared loss

$$L(f) = (Y - \underbrace{f(x)}_r)^2$$

- inner product
- differentiable
- not robust

② absolute loss

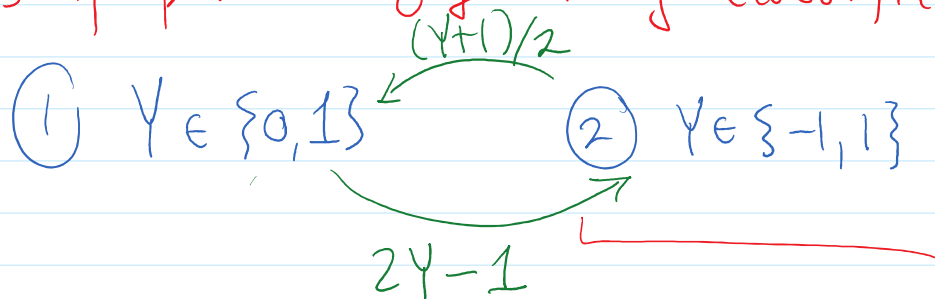
$$L(f) = |Y - \underbrace{f(x)}_r|$$

- not differentiable
- more robust

$$\boxed{r = Y - f(x)} \leftarrow$$

### Classification: (Binary)

Two ways of parameterizing binary classification



① 0-1 Loss

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

aside:  $Y \in \{-1, 1\}$  and  $f(x) \in \{-1, 1\}$

signs match = correct classification  
signs don't match = incorrect classification

for any  $f$  there is some  $h$  so that

$$f(x) = \text{sign}(h(x))$$

$$\text{sign}(a) = \begin{cases} 1 & a > 0 \\ -1 & a < 0 \\ 0 & a = 0 \end{cases}$$

real-valued function  $h$ ,  
 $\rightarrow h(x) > 0$  if  $Y = 1$

margin:

$Yh(x) > 0$  = correct classification

$\rightarrow h(x) < 0$  if  $Y = -1$

$< 0$  = incorrect classification

Similar to a residual in the regression case

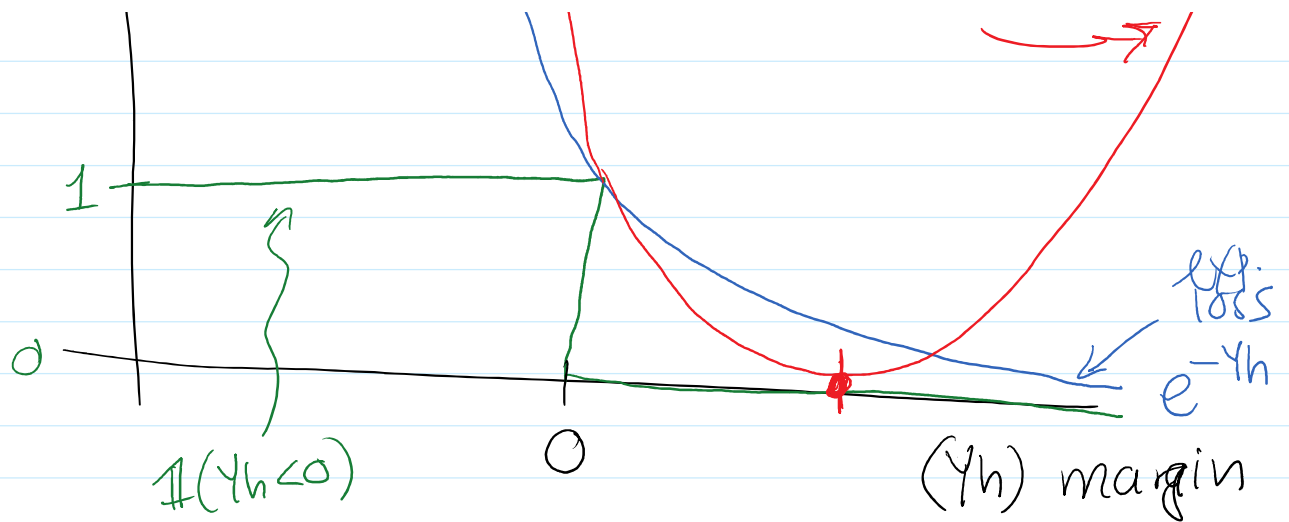
$|Yh(x)|$  = incorrect amount

We can write loss functions as a fn of margin.

0-1 Loss  $L(f) = \mathbb{1}(Yh(x) < 0)$   
 $= \mathbb{1}(f(x) \neq Y)$

$L(f)$

$(Yh)^2$



## ② Exponential Loss

$$L(f) = \exp(-Yh(x))$$

Boosted Methods : originally motivated as an ensemble method, consensus method

Idea: combine a series of weak methods to form a stronger one.

Weak Classifier Classification Setting (Binary)  
slightly better than random guessing

Boosting: ① sequentially learn weak class. to a series of modified training data

$\hat{f}_1, \hat{f}_2, \hat{f}_3, \hat{f}_4, \dots$

→ focus training on examples that were misclassified in prev. iter.

$t_1, t_2, t_3, t_4, \dots$

examples that were  
misclassified in prev. iter

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

weights:

higher for a good classifier,  
lower for a bad

weighted vote

## Alg: Ada Boost

(1)  $w_n = 1/N$  for each training example  $(x_n, y_n)$

(2) For  $b = 1, \dots, B$

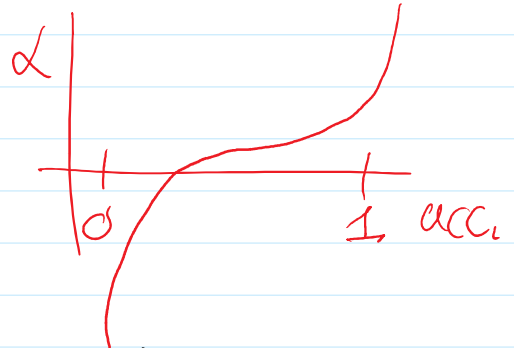
→ (a) fit  $\hat{f}_b$  on weighted  $(w_n)$  training data

weak method  
"a stump" - 1 split  
class tree

(b) compute err for  $\hat{f}_b$

$e_b = \text{weighted misclassification error}$

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



(d) re-weight training examples

(d) re-weight training examples

$w_n \leftarrow w_n$  if I classified correct  
or  
 $w_n \leftarrow e^{\alpha_b} w_n$  if " incorrect

(3) 
$$\hat{f}(x) = \text{Sign}\left(\sum_{b=1}^B \alpha_b \hat{f}_b(x)\right)$$

---

What is boosting really doing?

Why does it work?

Boosted method: additive model

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

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

Generalized additive method:

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

parameters for basis fn

$\alpha = \text{Linear function}$

$$| \quad \overset{b=1}{\underbrace{\quad \quad \quad}} \quad \beta = \text{basis function}$$

Ideally, find  $h$  by minimizing

$$\min_{\{\alpha_b, \gamma_b\}_{b=1}^B} \left( \sum_{b=1}^B \alpha_b \beta_b(x; \gamma_b) \right)$$

loss fn

Complicated problem. Soln: Stage-wise Additive Modeling  
Only optimize over one  $b$  at a time

For  $b=1, \dots, B$

(\*)  $\rightarrow$  (a)  $\alpha_b, \gamma_b = \underset{\alpha, \gamma}{\operatorname{argmin}} \sum_{n=1}^N L(\underbrace{f_{b-1}(x_n)} + \underbrace{\alpha \beta(x_n, \gamma_b)})$

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

---

Punchline: AdaBoost =  
Forward Stage-wise Additive Modeling  
w/  $L = \text{exponential Loss}$

---

Generalization! Can we do for Regression?

Yes!  $L = \text{squared loss}$

at each step, fit new model to residuals of prev. (in 2@)

More general losses?

Harder to do, Square loss or exp. loss lead to simple algs.

Soln:

Gradient Boosting: (for any loss)  
gradient descent to solve  $(*)$

Analogy:  $t = 1, 2, 3, \dots$   $\underline{\chi_t} \leftarrow \underline{\chi_{t-1}} + \underline{\alpha \nabla f(\chi_{t-1})}$

ensemble methods

RF: reduce variance, hard to overfit,  
easy to tune

Boosting: reduce bias, can overfit, needs tuning