

META CLASSIFIERS (01)

ENSEMBLE CLASSIFIERS.

WEAK  
CLASSIFIERS

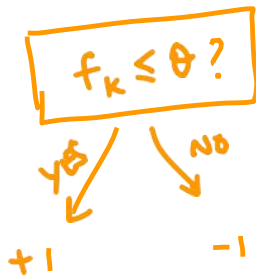
[better than  
random]



STRONG  
CLASSIFIERS

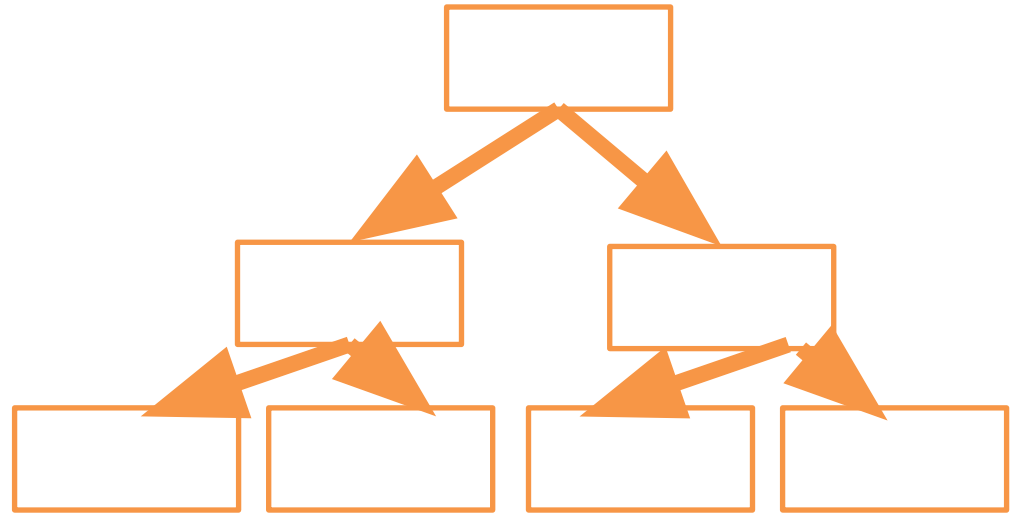
## Weak classifiers

DECISION  
STUMP



high bias, low variance

Overfit decision tree



...



.....



low bias, high variance

$$x_1, x_2, \dots, x_n \sim \mathcal{N}(\mu, 1)$$

$$\hat{\mu}_1 = x_1$$

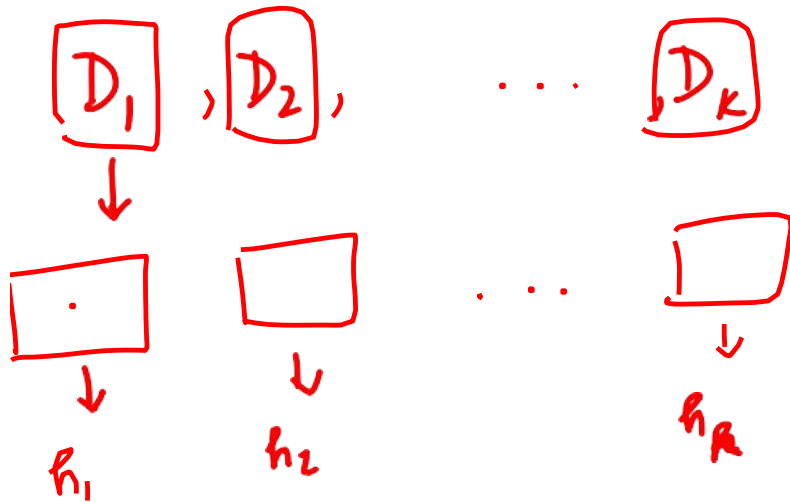
$$\hat{\mu}_2 = x_2$$

, ...

$$\hat{\mu}_n = x_n$$

$$\hat{\mu}_{ML} = \frac{1}{n} \sum x_i$$

Overfit  
decision  
trees



$$h_i: \mathbb{R}^d \rightarrow \{\pm 1\}$$

$$h^*(x) = \text{majority}(h_1(x), \dots, h_k(x))$$

## BAGGING - Bootstrap Aggregation.

Chance that a point appears in a dataset

$$1 - \underbrace{\left(1 - \frac{1}{n}\right)\left(1 - \frac{1}{n}\right) \dots \left(1 - \frac{1}{n}\right)}_{\left(1 - \frac{1}{n}\right)^n}$$
$$1 - \frac{1}{e} \quad (\text{as } n \rightarrow \infty)$$
$$\approx 66\%$$

$$D = \{(x_1, y_1) \dots (x_n, y_n)\}$$

- Create datasets  $D_1, \dots, D_k$  from  $D$  by  
    "Sampling with replacement".
- Run weak classifier on  $D_1, \dots, D_k$  to get  
     $h_1, \dots, h_k$
- Aggregate  $h_1, \dots, h_k$  using majority.

FEATURE BAGGING

→ Bag the features in addition to data points

Feature bagged decision trees → RANDOM FOREST

---

BOOTSTRAP - Sampling with Replacement ?

AGGREGATION - Majority. ?

# BOOSTING



ADA-BOOST

[ Freund & Schapire  
1995  
Nobel Prize ]

Distribution

$D$  over

$(x \times y)$   
 $\swarrow \mathbb{R}^d \quad \swarrow \{+1, -1\}$



unknown but fixed.

$x_1, \dots, x_n$  are iid from  $D$ .

$$h: \mathbb{R}^d_x \rightarrow \{+1, -1\}_y$$

Measure performance using

$$P_{\substack{(x,y) \sim D}} (h(x) \neq y)$$

Misclassification  
probability.

A weak learner is one which outputs a classifier  
Strong

$h$  for which

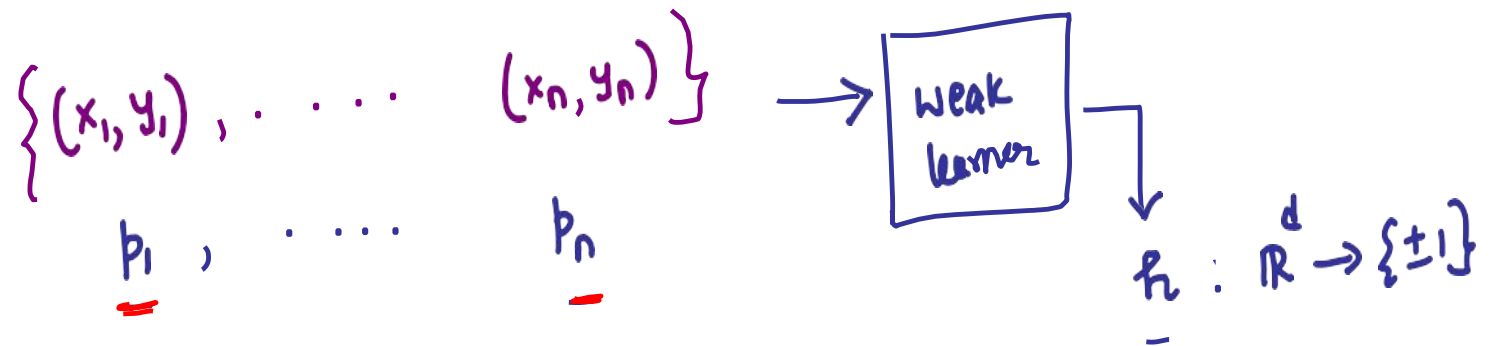
$$P_{\substack{(x,y) \sim D}} (h(x) = y) \geq \frac{1 - \epsilon}{2 + \gamma}$$

$$\gamma > 0$$

for any unknown but fixed  
distribution  $D$ .

# BOOSTING

Weak learner  $\rightarrow$  Strong learner.



$$\sum_i p_i = 1$$

$$\sum_{i=1}^n \underline{p_i} \underbrace{1(h(x_i) \neq y_i)} \leq \underbrace{\frac{1}{2} - \gamma}$$

$$\gamma > 0$$

Strong learner classifies all data points correctly.



# BOOSTING [ADABOOST]

$$S = \{ (x_1, y_1) \dots (x_n, y_n) \}$$

$$x_i \in \mathbb{R}^d$$
$$y_i \in \{\pm 1\}$$

Initialize

$$D_1(i) = \frac{1}{n} \quad \forall i$$

iteration

for  $t = 1, \dots, T$

$h_t$  = [Input  $S, D_t$  to the weak learner to get  $h_t$ ]

$$D_{t+1}(i) = \frac{D_t(i) \cdot e^{-\alpha_t y_i R_t(x_i)}}{\sum_{j=1}^n D_t(j) e^{-\alpha_t y_j R_t(x_j)}}$$

$\forall i$

$\alpha_t > 0$   
 $\forall t$

end.

$$\underline{H(x)} = \sum_{t=1}^T \alpha_t \boxed{R_t(x)}$$

$$\boxed{C(x) = \text{sign}(H(x))}$$

## ANALYSIS OF BOOSTING

$$D_{t+1}(i) = \frac{D_t(i) e^{-\alpha_t y_i h_t(x_i)}}{Z_t}$$

→ update rule.

$$Z_t = \sum_{j=1}^n D_t(j) e^{-\alpha_t y_j h_t(x_j)}$$

$$D_{t+1}(i) = \left( \frac{D_{t-1}(i) e^{-\alpha_{t-1} y_i h_{t-1}(x_i)}}{Z_{t-1}} \right) \cdot \frac{e^{-\alpha_t y_i h_t(x_i)}}{Z_t}$$

$$D_{t+1}(i) = \frac{\frac{1}{n} e^{-\sum_{l=1}^t \alpha_l \underline{y_i} h_l(x_i)}}{\prod_{l=1}^t z_l}$$

$$D_{T+1}(i) = \frac{\frac{1}{n} e^{-y_i H(x_i)}}{\prod_{l=1}^T z_l} \longrightarrow \sum_{l=1}^T \alpha_l h_l(x_i)$$

$$\underbrace{\sum_{i=1}^n D_{t+1}(i)}_{=1} = \frac{\sum_{i=1}^n \frac{1}{n} e^{-y_i H(x_i)}}{\prod_{l=1}^T z_l}$$

$$\prod_{l=1}^T z_l = \frac{1}{n} \sum_{i=1}^n e^{-y_i H(x_i)} \quad \text{--- (A)}$$

Recall,  $C(x) = \text{Sign}(H(x))$

$$\mathbb{1}(C(x_i) \neq y_i) = \mathbb{1}(\underline{H(x_i) y_i} < 0) \leq e^{-H(x_i) y_i} \quad \text{--- (B)}$$

$$\Rightarrow \prod_{\ell=1}^T z_{\ell} = \frac{1}{n} \sum_{i=1}^n e^{-y_i H(x_i)} \quad (\text{from (A)})$$

$$\geq \frac{1}{n} \sum_{i=1}^n \mathbb{1}(y_i H(x_i) < 0) \quad (\text{from (B)})$$

$$\geq \frac{1}{n} \underbrace{\sum_{i=1}^n \mathbb{1}(c(x) \neq y_i)}_{\text{error}(c)}$$

$$\boxed{\text{error}(c) \leq \prod_{\ell=1}^T z_{\ell} \leq \square}$$

$$\prod_{l=1}^T z_l = \prod_{l=1}^T \left[ \sum_{j=1}^n D_l(j) e^{-\alpha_l y_j h_l(x_j)} \right]$$

$$= \prod_{l=1}^T \left[ \frac{\sum_{j=1}^n e^{-\alpha_l} D_l(j) \mathbb{1}(y_j = h_l(x_j))}{\sum_{j=1}^n e^{\alpha_l} D_l(j) \mathbb{1}(y_j \neq h_l(x_j))} \right]$$

$$= \prod_{l=1}^T e^{-\alpha_l} \left( 1 - \text{error}(h_l) \right) + e^{\alpha_l} \left( \text{error}(h_l) \right)$$

Choose  $d_L$  st

$e^{-d_L} (1 - \text{error}(h_L)) + e^{d_L} (\text{error}(h_L))$  is  
as small as possible.

$$d_L^* = \ln \sqrt{\frac{1 - \text{er}(h_L)}{\text{er}(h_L)}}$$

↑

$$= \prod_{L=1}^T 2 \sqrt{\text{error}(h_L) (1 - \text{error}(h_L))}$$

$$\leq \prod_{L=1}^T 2 \sqrt{\left(\frac{1}{2} - \gamma\right) \left(\frac{1}{2} + \gamma\right)}$$

[show this]

$$= \prod_{L=1}^T 2 \sqrt{\frac{1 - 4\gamma^2}{4}}$$

$$e^{-d_L} (1 - \theta) + e^{d_L} \theta$$

$$(1 - \theta) e^{-d_L} (-1) + e^{d_L} \theta = 0$$

$$-1 + \theta e^{-d_L} + \theta e^{d_L} = 0$$



$$= \prod_{\ell=1}^T \sqrt{1 - 4\gamma^2}$$

$$\leq \prod_{\ell=1}^T \left( e^{-4\gamma^2} \right)^{1/2} = \prod_{\ell=1}^T \underline{e^{-2\gamma^2}}$$

$$= \underline{e^{-2T\gamma^2}} \quad \text{--- (2)}$$

$$\frac{1}{n} \sum_i \mathbb{1}(cx_i \neq y_i)$$

↓

$$\underline{\text{error}(c)}$$

$$\leq \prod_{\ell=1}^T z_{\ell} \leq \underbrace{e^{-2T\gamma^2}}_{\uparrow} \leq \frac{1}{2n}$$

$$\text{if } T \geq \frac{1}{2\gamma^2} \ln(2n)$$

$$\Rightarrow \underline{\text{error}(c) = 0}$$

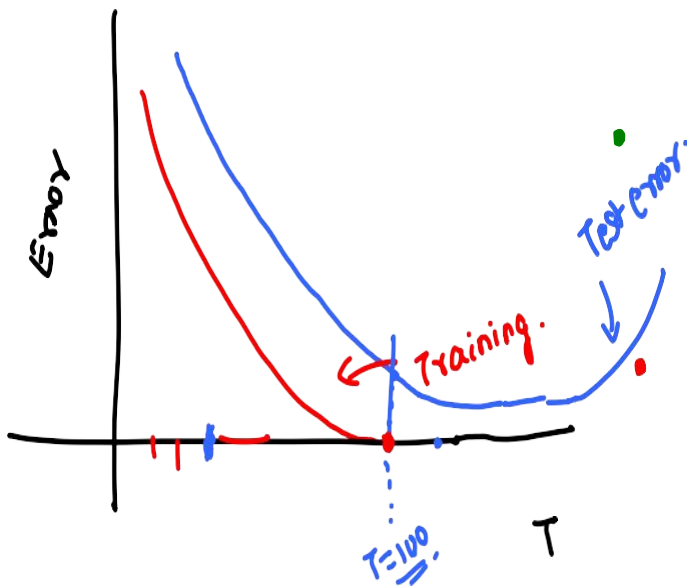
⇓  
Strong learner!

## Boosting-Summary

- Every bag depends on previous bag's performance

Next bag "focusses" on previous bag's mistakes

Can show training error goes to 0 as number of rounds increases



Cannot run in parallel – hence sometimes bagging is preferred in practice

Usually weak learners are decision stumps  $\rightarrow$  high bias, low variance  
Boosting reduces bias without affecting variance a lot

Why are there so many binary classification algorithms?

$$x \in \mathbb{R}^d$$

$$y \in \{\pm 1\}$$

$$\min_w \sum_{i=1}^n \mathbb{1}(w^T x_i \cdot y_i < 0)$$

NP-Hard problem



Linear Regression for classification

$$\min_h \sum_{i=1}^n (h(x) - y_i)^2$$

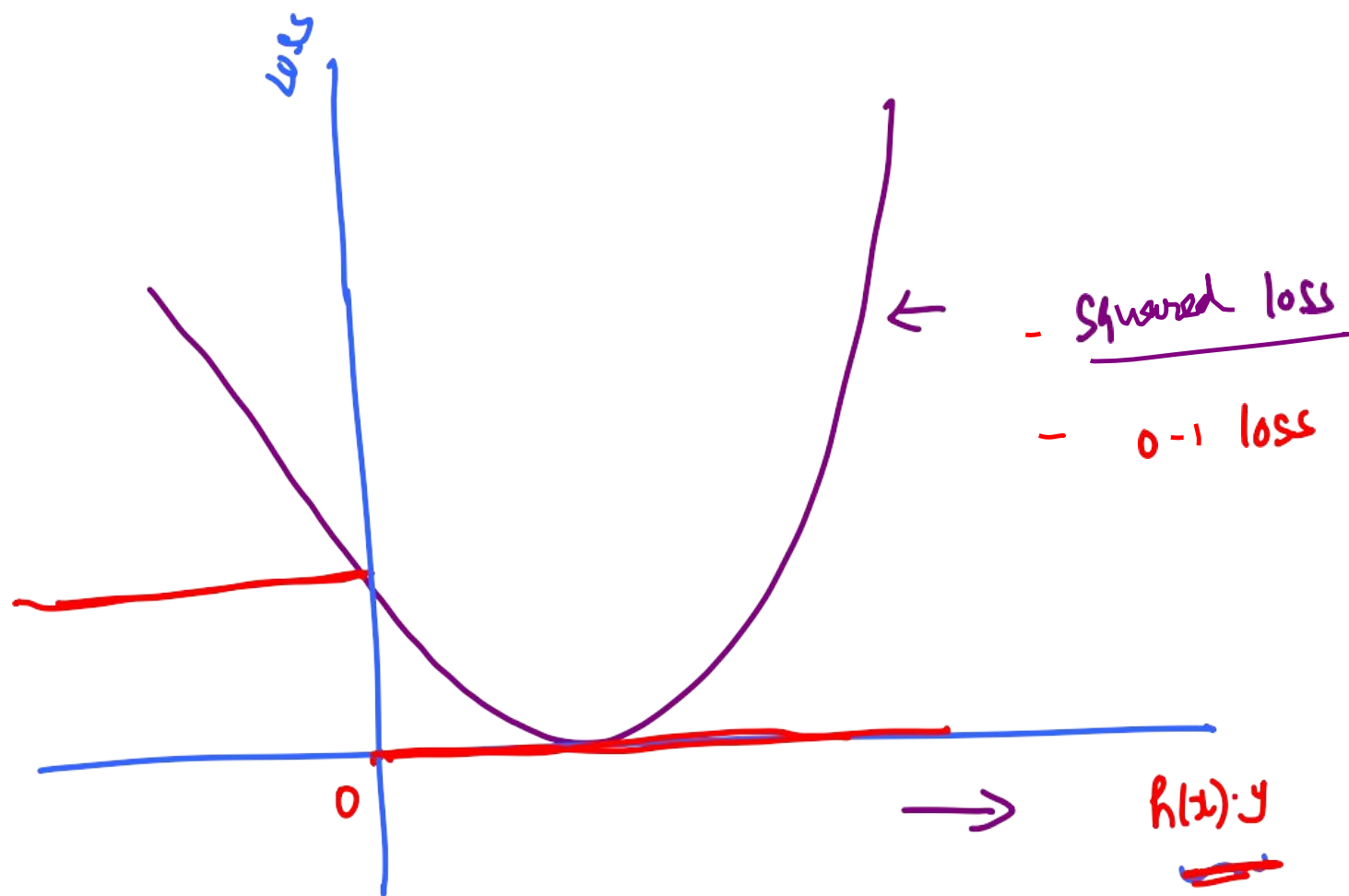
Final classifier:  $\text{Sign}(h(x))$

Loss per point for h

$$(h(x) - y)^2$$

$$= \begin{cases} (h(x) - 1)^2 & \text{if } y = 1 \\ (h(x) + 1)^2 & \text{if } y = -1 \end{cases}$$

$$\begin{aligned} \text{if } y = 1 & \quad \frac{(h(x))^2 + 1 - 2h(x)}{2} = (h(x))^2 + 1 - 2h(x)y \\ y = -1 & \quad \frac{(h(x))^2 + 1 + 2h(x)}{2} = (h(x))^2 + 1 - 2h(x) \cdot y \\ & \quad = (h(x) \cdot y - 1)^2 \end{aligned}$$



## SUPPORT VECTOR MACHINES

$$\min_{w, \xi} \quad \frac{1}{2} \|w\|^2 + C \sum_{i=1}^n \xi_i$$

$$\begin{aligned} \underline{w^T x_i y_i + \xi_i} &\geq 1 \quad \forall i \\ \xi_i &\geq 0 \quad \forall i \end{aligned}$$

Equivalently

$$\begin{aligned} \xi_i &\geq 1 - w^T x_i y_i \quad \forall i \\ \xi_i &\geq 0 \quad \forall i \end{aligned}$$

Equivalently

$$\xi_i \geq \max(1 - w^T x_i y_i, 0)$$

•  $\min_{w, \xi}$

$$\frac{1}{2} \|w\|^2 + C \sum_{i=1}^n \xi_i$$

$$\xi_i \geq \max(1 - w^T x_i y_i, 0)$$

Equivalently,

•  $\min_{w, \xi}$

$$\frac{1}{2} \|w\|^2 + C \sum_{i=1}^n \max(1 - w^T x_i y_i, 0)$$

**Model  
Regularization  
term**

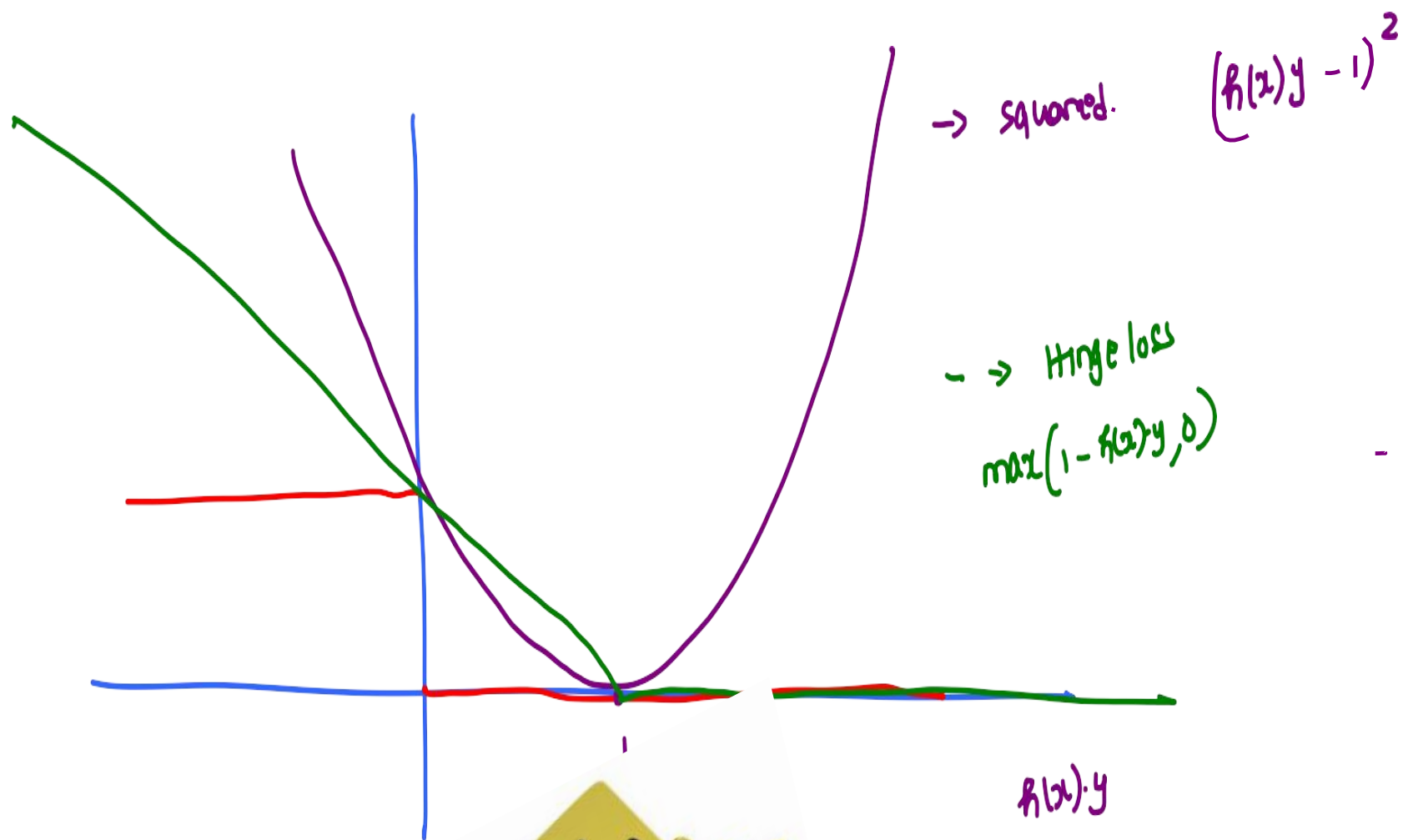
**Data  
Loss  
term**

So what exactly is the loss?

$$\ell(w, (x, y)) = \max(1 - w^T x y, 0)$$

**HINGE LOSS**





**Hinge**



Courtesy: Google images

# Logistic regression

$$\max_{\omega} \prod_{i=1}^n \underbrace{\left( g(\omega^T x_i) \right)^{z_i} \left( 1 - g(\omega^T x_i) \right)^{(1-z_i)}}_{}$$

$$\begin{aligned} z_i &= 1 & \text{if } y_i &= 1 \\ z_i &= 0 & \text{if } y_i &= -1 \end{aligned}$$

$$g(\theta) = \frac{1}{1 + e^{-\theta}}$$

$$\max_w \prod_{i=1}^n \underbrace{\left( \underline{g(\omega^T x_i)} \right)^{z_i} \left( \underline{1 - g(\omega^T x_i)} \right)^{(1-z_i)}}$$

$$\max_w \sum_{i=1}^n z_i \log(g(\omega^T x_i)) + (1-z_i) \log(1 - g(\omega^T x_i))$$

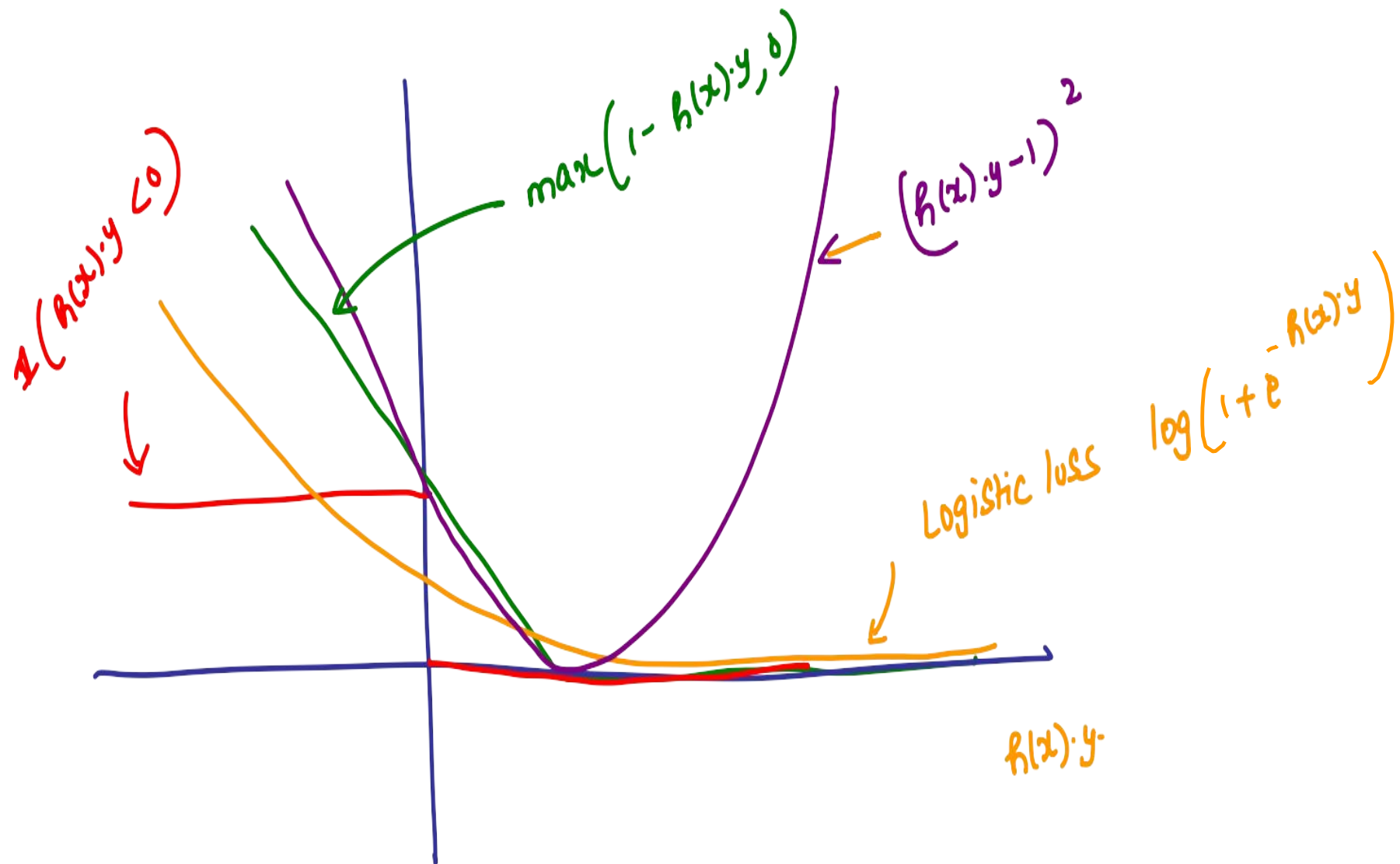
$$\equiv \min_w \sum_{i=1}^n \left[ -z_i \log(g(\omega^T x_i)) + (z_i - 1) \log(1 - g(\omega^T x_i)) \right]$$

Loss for a single point when  $z_i = 1$  ( $y_i = 1$ )

$$\begin{aligned} -\log(g(w^T x_i)) &= -\log\left(\frac{1}{1 + e^{-w^T x_i}}\right) \\ &= \log(1 + e^{-w^T x_i}) = \boxed{\log(1 + e^{-y_i w^T x_i})} \end{aligned}$$

Loss for a single point when  $z_i = 0$  ( $y_i = -1$ )

$$\begin{aligned} &= -\log(1 - g(w^T x_i)) = -\log\left(1 - \frac{1}{1 + e^{-w^T x_i}}\right) \\ &= -\log\left(\frac{e^{-w^T x_i}}{1 + e^{-w^T x_i}}\right) = \log(1 + e^{w^T x_i}) = \boxed{\log(1 + e^{-y_i w^T x_i})} \end{aligned}$$

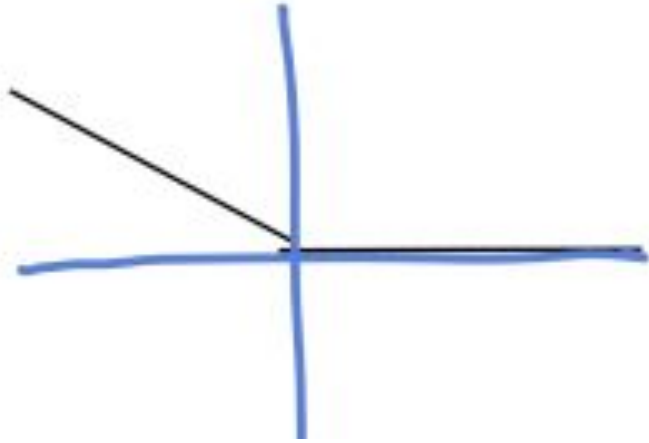


## Perceptron

$$W_{t+1} = W_t + x_t y_t$$

consider the hinge loss

$$\begin{aligned} \text{loss}(w, (x, y)) \\ = \max(0, 1 - w^T x y) \end{aligned}$$



$$\text{loss}(w, (x, y)) = \max(0, 1 - \vec{w}^T x y)$$

$$\frac{\partial \text{loss}}{\partial w} = \begin{cases} -x \cdot y & \text{if } (\vec{w}^T x) y < 0 \\ = 0 & \text{if } (\vec{w}^T x) y > 0 \\ = [-1, 0] \cdot x \cdot y & \text{if } (\vec{w}^T x) y = 0 \end{cases}$$

sub-gradient

↳ choose  $-xy$  if mistake.

When mistake

$$W_{t+1} = W_t - \eta_t (-x_t y_t)$$

↘ = 1

- Perceptron can be viewed as  
S.G.D with Hinge loss with step size  
= 1



## BOOSTING

Recall

$$\left( \prod_{t=1}^T z_t \right) = \frac{1}{n} \sum_{i=1}^n e^{- \left( \sum_{t=1}^T \alpha_t h_t(x_i) \right) y_i}$$
$$:= \frac{1}{n} \sum_{i=1}^n e^{- H_T(x_i) y_i}$$

at round  $T$ , we so far have

$$\sum_{t=1}^{T-1} \alpha_t h_t(x)$$

We need to add one more  
classifier in round  $T$  to get

$$\sum_{i=1}^{T-1} \alpha_t h_t(x) + \underline{\alpha_T h_T(x)}$$

↙ greedy choice  
to minimize

$$\frac{1}{n} \sum_{i=1}^n e^{-\left[ \sum_{t=1}^{T-1} \alpha_t h_t(x_i) + \alpha_T h_T(x) \right] y_i}$$

---

Define  $\text{loss}(h, y) = e^{-y h(x)}$

- At every round, one chooses

$\alpha_t$  to minimize  $Z_t$ .

- Thus AdaBoost can be viewed  
as "greedy / co-ordinate descent"  
on exponential loss

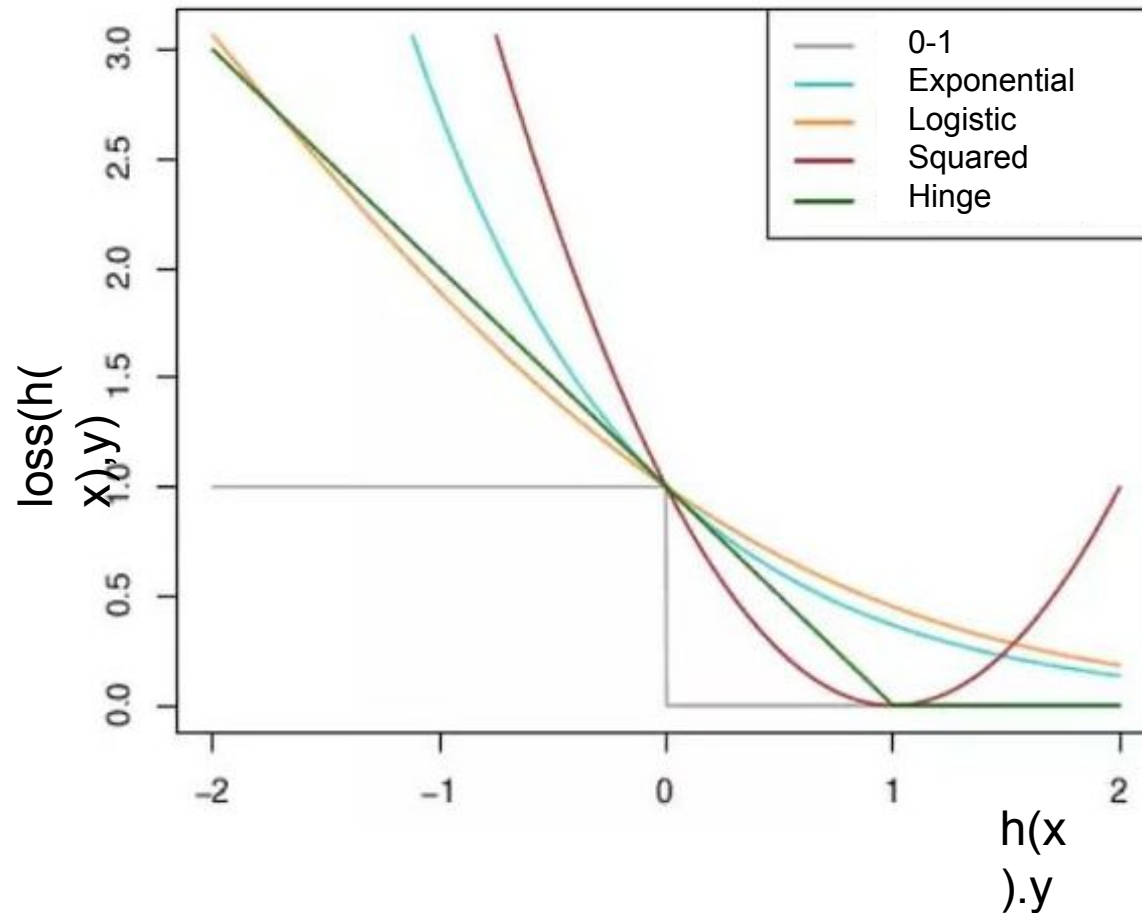
## What about Perceptron and Boosting?

- one can argue perceptron update rule is equivalent to doing a SSGD on hinge loss with stepsize = 1  
↓  
Stochastic sub-gradient descent.

$$W_{t+1} = W_t + \underline{x_i y_i}$$

- Boosting (Adaboost) can be viewed as "greedy / co-ordinate wise" descent of exponential loss  
↑  
 $e^{-f(x) \cdot y}$   
exponential loss

## SUMMARY



-> The 0-1 loss is NP-hard to optimize even for linear classifiers

-> Different algorithms get around this by using a “surrogate” loss function

-> Surrogates are usually **“convex” surrogates** that are easy to optimize