

2<sup>nd</sup> Method to solve for the separating hyperplane in the realizable case:

We know that  $y_i \langle \underline{w}, \underline{x}_i \rangle \geq 0$  for correct classification.

So, if any example is incorrectly classified, it would give

$$\rightarrow \underbrace{y_i \langle \underline{w}, \underline{x}_i \rangle}_{\text{LHS}} < 0 \Rightarrow \text{incorrect classification}$$

We design an iterative method

$$\underline{w} \leftarrow \underline{w}^{(0)}$$

We update  $\underline{w}$  in steps (iterations)

Scan the examples  $(x_i, y_i)_m$

If the  $i^{\text{th}}$  example is incorrectly classified,

$$\underbrace{y_i \langle \underline{w}^{(t)}, \underline{x}_i \rangle}_{\text{LHS}} < 0$$

We apply an update step proposed by Rosenblatt

$$\underline{w}^{(t+1)} \leftarrow \underline{w}^{(t)} + \underbrace{y_i \underline{x}_i}_{\text{LHS}}$$

This update step will increase the value of  $y_i \langle \underline{w}^{(t)}, \underline{x}_i \rangle$

$$\begin{aligned} \text{LHS. } y_i \langle \underline{w}^{(t+1)}, \underline{x}_i \rangle &= y_i \langle \underline{w}^{(t)} + y_i \underline{x}_i, \underline{x}_i \rangle \\ &= \underbrace{y_i \langle \underline{w}^{(t)}, \underline{x}_i \rangle}_{\text{LHS. after the update}} + \underbrace{y_i y_i \underline{x}_i^T \underline{x}_i}_{\text{always +ve}} \\ &= \underline{y_i \langle \underline{w}^{(t)}, \underline{x}_i \rangle} + \underline{y_i^2 \|\underline{x}_i\|^2} \end{aligned}$$

Perceptron algorithm  
Linear hyperplanes.

This is the Perceptron Algorithm.

$$\underline{w}^{(0)} \leftarrow (0, \dots, 0)$$

for  $t = 0, 1, 2, \dots$

$i := 1$

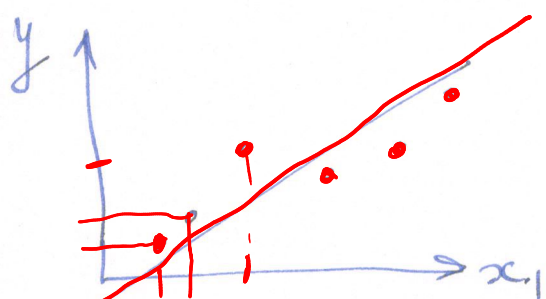
if  $\exists i$  such that  $y_i \langle \underline{w}^{(t)}, \underline{x}_i \rangle \leq 0$

$$\underline{w}^{(t+1)} \leftarrow \underline{w}^{(t)} + \underline{y_i x_i}$$

else return  $\underline{w}^{(t)}$

The algorithm is guaranteed to converge  
but only for the Realizable case.

## Linear Regression



$\underline{x} \equiv [x_1]$  in the 1D case.

$$y = w_0 + w_1 x_1$$

$y$  depends linearly on  $x$ .

$(x_i, y_i)_{i=1, \dots, n}$

A linear hypothesis can be written as

$$\underline{h(x)} \equiv \langle \underline{w}, \underline{x} \rangle$$

Augmented vectors  
 $\underline{h_w(x)}$  homogeneous linear map.

## Loss function

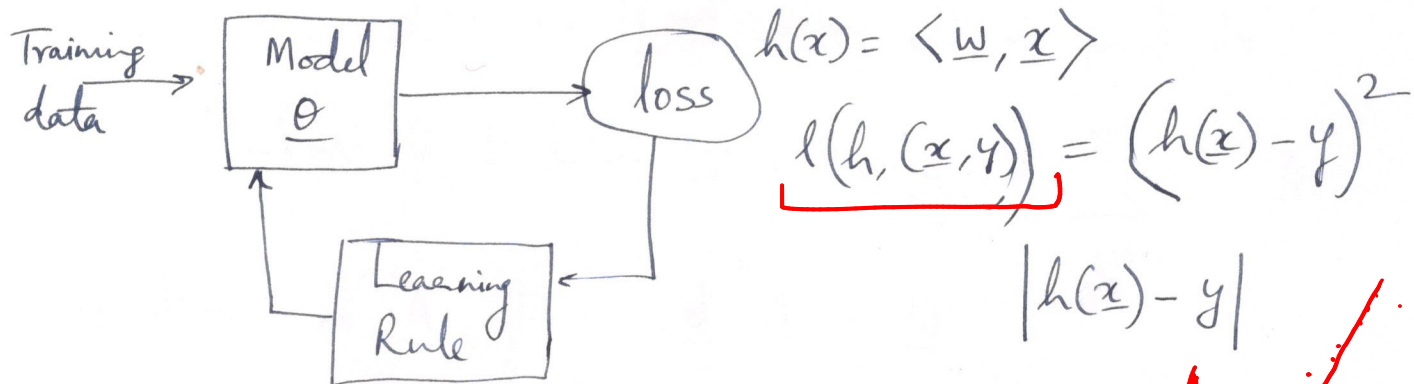
$$\ell(\underline{h}, (\underline{x}, y)) = ?$$

$$(h(\underline{x}) - y)^2$$

squared error loss

5 x 4 loss  
x 6 loss  $\rightarrow$  reward.

absolute error loss.



## Empirical Risk Minimization (ERM Rule)

Emp. Risk is expected loss over the  $m$  training examples.

Risk.

$$L_S(h) = \left( \frac{1}{m} \right) \sum_{i=1}^m (h(\underline{x}_i) - y_i)^2$$

ERisk.

$S$ : Sample  
(Training Set)

Non realizable

$\mathcal{D}$ : True Risk  
Data generating distribution

Objective: To solve the ERM problem

$$\underline{w}^* = \arg \min_{\underline{w}} L_S(h) = \arg \min_{\underline{w}} \frac{1}{m} \sum_{i=1}^m (\langle \underline{w}, \underline{x}_i \rangle - y_i)^2$$

Differentiating w.r.t  $\underline{w}$  and equating to zero.

$$\frac{2}{m} \sum_{i=1}^m (\langle \underline{w}, \underline{x}_i \rangle - y_i) \underline{x}_i = 0$$

$$\sum_{i=1}^m (\underline{x}_i^T \underline{w} - y_i) \underline{x}_i = 0$$

$$\sum_{i=1}^m \underline{x}_i \underline{x}_i^T \underline{w} = \sum_{i=1}^m y_i \underline{x}_i$$

$\underline{x}_i$   $\underline{x}_i^T$

$\begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix}$   $\begin{bmatrix} x_1 & x_2 & x_3 \end{bmatrix}$

$3 \times 1$   $1 \times 3$   $= 3 \times 3$  outer product matrix

$A$   $b$



Writing in the form  $A \underline{w} = \underline{b}$

$$A = \sum_{i=1}^m \underline{x}_i \underline{x}_i^T \quad \text{outer product of vectors.}$$

$$\underline{w} = \underline{A}^+ \underline{b}$$

if  $A$  is invertible

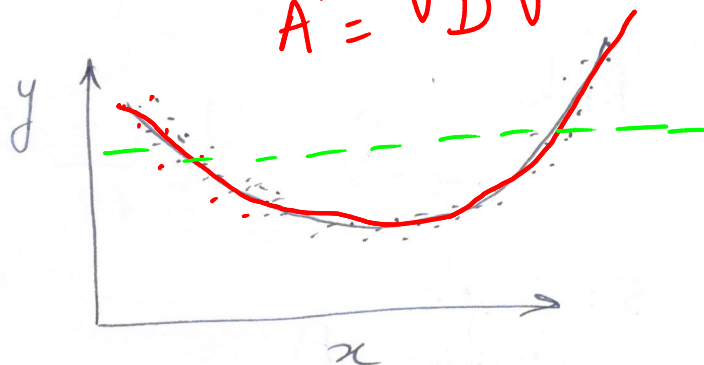
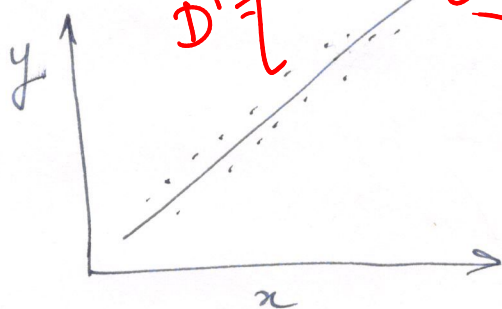
else

where  $A^+$  is the pseudo inverse.

$$A^+ = V D^+ V^T$$

$$A = V D V^T$$

$$D^+ = \begin{bmatrix} 1/a_1 & 1/a_2 & \dots & 0 \end{bmatrix}$$



A linear function is not a good fit.

## Polynomial Regression

Assuming that  $x$  is a scalar

$H =$  linear fn.  
Quadratic  
Cubic  
quartic  $n??$

$$y = \underline{p(x)} = a_0 + a_1 x + a_2 x^2 + a_3 x^3 + \dots + a_n x^n$$

$y$  depends on  $x$  through a polynomial function.

$$\begin{bmatrix} a_0 \\ a_1 \\ \vdots \\ a_n \end{bmatrix}$$

Unknown.

We use a trick

$$x \rightarrow \underline{x} = \begin{bmatrix} 1 \\ x \\ x^2 \\ \vdots \\ x^n \end{bmatrix}$$

$$\underline{\underline{w}} = \begin{bmatrix} a_0 \\ a_1 \\ a_2 \\ \vdots \\ a_n \end{bmatrix}$$

$$p(x) = \langle \underline{x}, \underline{w} \rangle$$

$$\mathbb{R} \rightarrow \mathbb{R}^{n+1}$$

After transformation of the input features, we can estimate the polynomial mapping of the input  $x \in \mathbb{R}$  to  $y$  by using linear regression.

Linear regression maps  $x$  to a scalar  $y$ .  
$$x \xrightarrow[\text{linear}]{\langle w, x \rangle} y \xrightarrow{\text{squash}} [0, 1]$$
  
 $y \in \mathbb{R}$       probability.

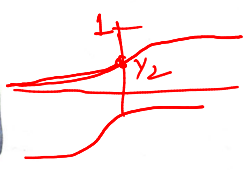
$P(y \in +1) = p_1$   
 $P(y \in -1) = 1 - p_1$

If we can restrict the output in the range  $[0, 1]$ , it can be interpreted as probability.

If we consider the task of binary classification, we can interpret the squashed output of linear regression as the probability of input  $x$  taking the class label  $y=1$ .  
target  $y \in \{0, 1\}$ .

$$\underline{x} \xrightarrow[\text{linear}]{\langle w, x \rangle} \underline{z} \xrightarrow[\text{squashing function}]{\phi = \frac{1}{1 + \exp(-z)}} \phi(z)$$
  
 $z \in \mathbb{R}$

squashing  
$$\phi(z) = \frac{1}{1 + \exp(-z)}$$
  
sigmoid  
 $z \in \mathbb{R}$   
 $\phi(z) \in [0, 1]$



$\phi(z)$  is interpreted as the probability that  $x$  takes label  $y=1$   
 $1 - \phi(z)$  " "  $y=0$

By squashing the linear regression output, we have a classifier which gives probability of the class label for binary classification task.  
 $\rightarrow \{1, 0\}$   
 $\{1, +1\}$

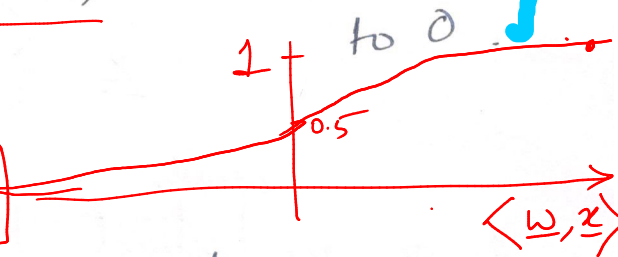
The Hypothesis class for a logistic regression model is given by

$$H_{\text{sig}} = \left\{ \underline{x} \mapsto \phi_{\text{sig}}(\langle \underline{w}, \underline{x} \rangle) : \underline{w} \in \mathbb{R}^d \right\}$$

If  $\langle \underline{w}, \underline{x} \rangle$  is large,  $\phi_{\text{sig}}(\langle \underline{w}, \underline{x} \rangle)$  is close to 1.  
small to 0

$$h_w(\underline{x}) \in [0, 1]$$

$$h_w(\underline{x}) = \phi(\langle \underline{w}, \underline{x} \rangle)$$



for  $y=1$   $h_w(\underline{x})$  large  
 GT What we expect for correct prediction

Loss formulation  
 $1 - h_w(\underline{x})$  should be small for correct  
 loss for  $y=1$

for  $y=0$   $h_w(\underline{x})$  small ✓ OK.  
 GT  $1 - h_w(\underline{x})$  is large OK.

$h_w(\underline{x})$  should be small for correct  
 for  $y=0$

Simplifying the loss for  $y=1$

$$1 - h_w(\underline{x}) = 1 - \frac{1}{1 + \exp(-\langle \underline{w}, \underline{x} \rangle)}$$

$$= \left( \frac{\exp(-\langle \underline{w}, \underline{x} \rangle)}{1 + \exp(-\langle \underline{w}, \underline{x} \rangle)} \right) \times \frac{\exp(\langle \underline{w}, \underline{x} \rangle)}{\exp(\langle \underline{w}, \underline{x} \rangle)}$$

$$= \frac{1}{1 + \exp(\langle \underline{w}, \underline{x} \rangle)}$$

Loss for  $y=0$   $h_w(\underline{x}) = \frac{1}{1 + \exp(-\langle \underline{w}, \underline{x} \rangle)}$



$y \in \{0, 1\}$

Loss for  $y=0$  /  $y=-1$

Loss for  $y=1$

$$\left( \frac{1}{1 + \exp(-\langle \underline{w}, \underline{x} \rangle)} \right)$$

$$\left( \frac{1}{1 + \exp(\langle \underline{w}, \underline{x} \rangle)} \right)$$

If we recode the target  $y \in \{-1, 1\}$  instead of  $y \in \{0, 1\}$ , we can unify the two losses for  $y=-1$  and  $y=1$

Unified loss

$$\frac{1}{1 + \exp(y \langle \underline{w}, \underline{x} \rangle)}$$

Captures both the loss terms.

ERM

The objective function to be minimized can be written as

$$\min_{\underline{w}} \frac{1}{1 + \exp(y \langle \underline{w}, \underline{x} \rangle)}$$

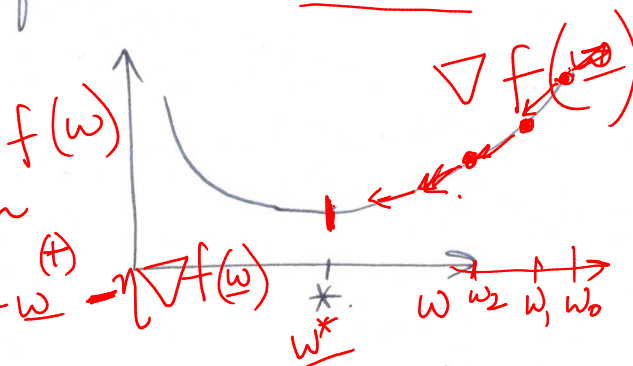
$$\text{or } \min_{\underline{w}} 1 + \exp(-y \langle \underline{w}, \underline{x} \rangle)$$

Log of the loss

$$\min_{\underline{w}} \log(1 + \exp(-y \langle \underline{w}, \underline{x} \rangle))$$

The ~~loss function~~ is convex.

Convex have only the global minimum



Use gradient descent



Half space classifiers

Realizable Case

0-1 loss

Linear Prog.

Linear Regression

Non Realizable Case

Squared Error loss

Pseudo Inverse

Logistic Regression

Non Realizable Case

Logistic Regression loss

Gradient Descent