

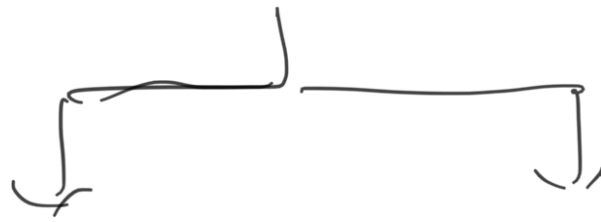
1. Bayes Classifier

- (i) Probability of Class 0, 1
- (ii) Conditional densities



Posterior and implement Bayes Classifier (Perfect information)

2. Data — $(x_i, y_i)_{i=1}^n$



Generalization



Estimates priors
and conditional
densities.

Discrimination



Posterior estimation
using ERM.

(i) K-NN

(ii) Naive Bayes

(i) Linear regression
for Classification

(ii) Logistic regression

Discriminative:

$$R(h) = E_{\substack{x, y \\ x \in \mathbb{R}^d, y \in \{0, 1\}}} [L(y, h(x))] \quad \leftarrow \text{Risk of a classifier.}$$

$$h^* = \min_h E_{x, y} [L(y, h(x))].$$

$$\mathcal{D} = \{(x_i, y_i)_{i=1}^n\}.$$

$$\hat{R}(h) = \underbrace{\sum_{i=1}^n L(y_i, h(x_i))}_n$$

$$(i) \quad \mathcal{H} = \{w^T x : x \in \mathbb{R}^d\}$$

$$h(x) = w^T x.$$

(iii) Perceptron bound on linear regression

$$L(y, h(x)) = \mathbb{1}_{\{y \neq h(x)\}},$$

↗

Indicator loss function.

$$L_{\text{lin}}(y, h(x)) = (y - h(x))^2$$

$$J(w) = \frac{1}{n} \sum_{i=1}^n (y_i - w^T x_i)^2.$$

$$w^* = \min_{w \in \mathbb{R}^d} J(w).$$

$$(a) \quad \nabla J(w) = 0 \quad \Rightarrow \quad w^* = (\bar{x}^T \bar{x})^{-1} \bar{x}^T y.$$

$$X = \begin{pmatrix} x_1^T \\ x_2^T \\ \vdots \\ x_n^T \end{pmatrix}$$

$$y = \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{pmatrix} \begin{matrix} \rightarrow +1 \text{ or } -1 \\ \rightarrow +1 \text{ or } -1 \\ \vdots \\ \rightarrow +1 \text{ or } -1 \end{matrix}$$

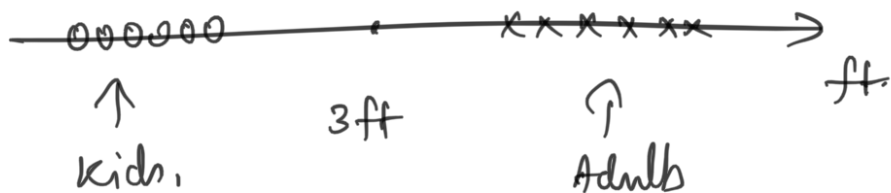
$$(b) \quad w_0 = \begin{pmatrix} 1 \\ \vdots \\ 1 \end{pmatrix}$$

(10)

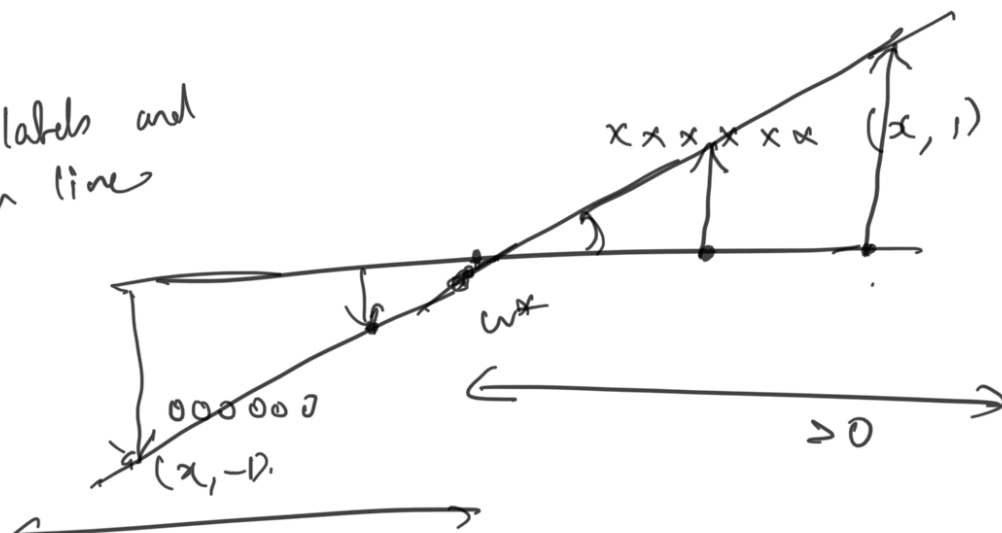
$$w_{n+1} = w_n - \alpha \nabla J(w_n)$$

- (i) Batch gradient
 - (ii) Stochastic gradient
 - (iii) mini batch gradient.
-

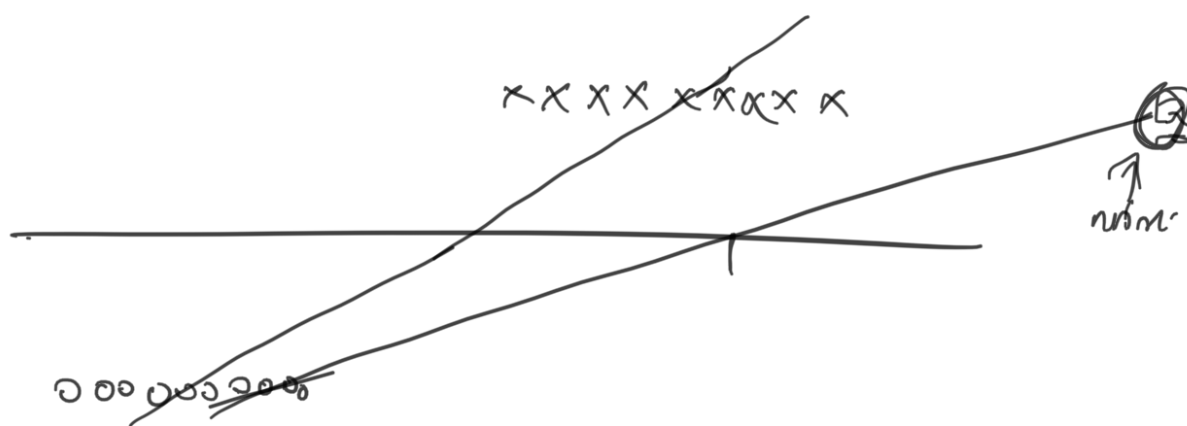
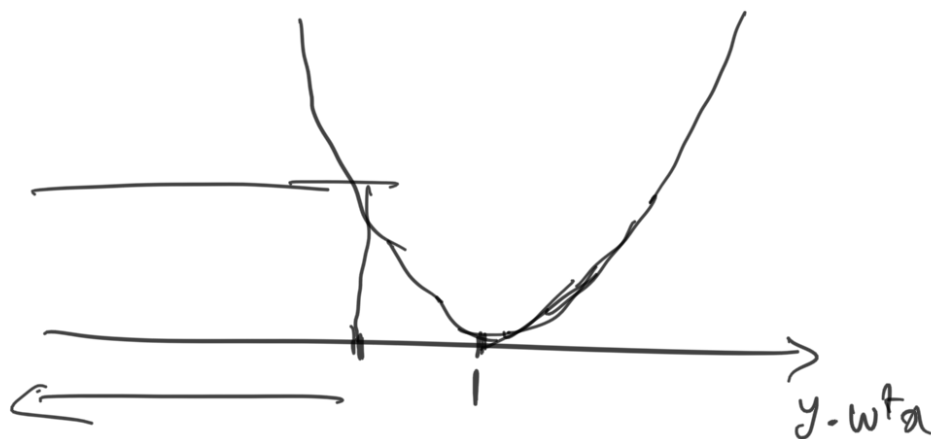
D.



with y labels and
find a line



$$h(x) = \begin{cases} +1 & \text{if } w^+ x > 0 \\ -1 & \text{if } w^+ x < 0 \end{cases}$$

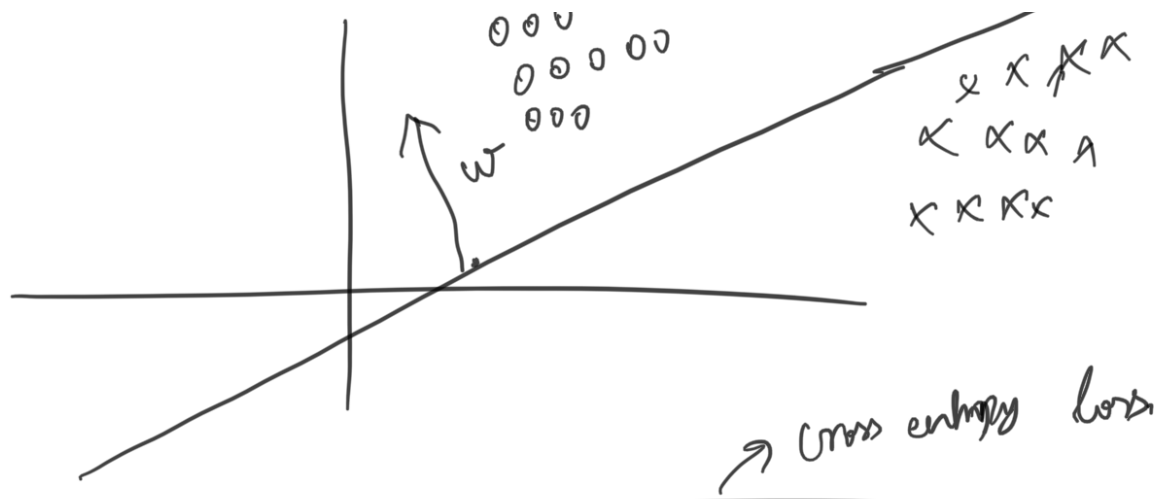


Disadv:

Linear regression for classification is not suitable when there is outlier in the present

Adv: Simplicity and closed form answer.

$\{ w^+ x \}$



$$L(y, h(x)) = y \log h(x) + (1-y) \log(1-h(x))$$

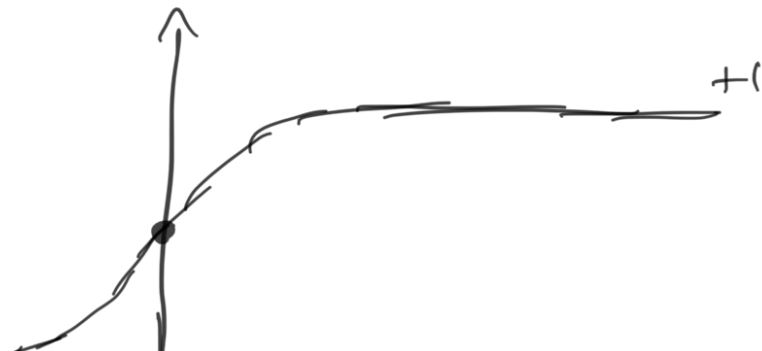
$$h(x) = \frac{1}{1 + e^{-(w^T x)}} = G(w^T x)$$

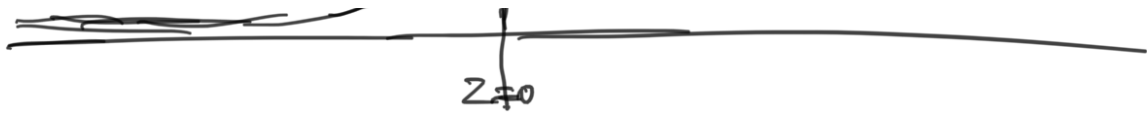
$$h: \mathbb{R}^d \rightarrow \mathbb{R}$$

↓ score function

$$h: \mathbb{R}^d \rightarrow [0, 1]$$

$G \rightarrow$ Sigmoid function $G(z) = \frac{1}{1+e^{-z}}$

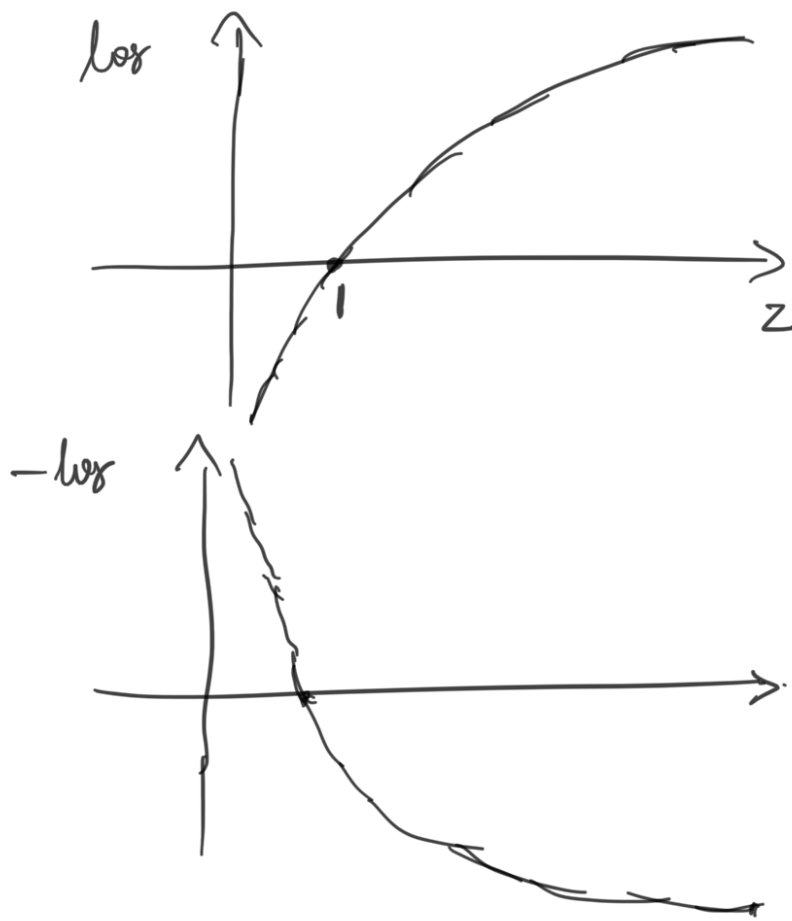




$$L(y, h(x)) = \underline{y \log h(x)} + \underline{(1-y) \log (1-h(x))}$$

$$h(x) = \frac{1}{1 + e^{-w^T x}}$$

$$\log h(x) = -\log(1 + e^{-w^T x})$$



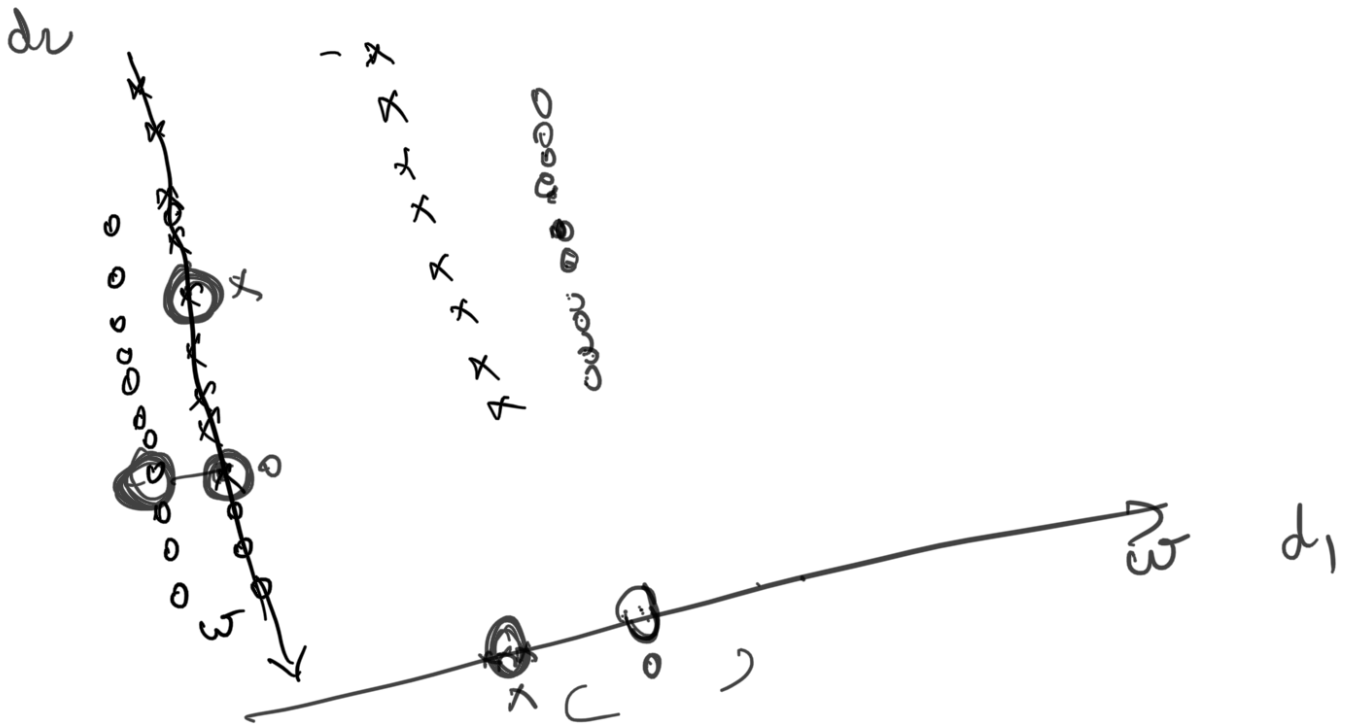
$x \in \text{class 1}$
 $w^T x \gg 0$
 $x \notin \text{class 1}$
 $w^T x \ll 0$

$$L(y, h(x)) =$$

$$-y \log(1 + e^{-w^T x}) + (1-y) \log \underline{e^{-w^T x}}$$

$$- \left[y \log (1 + e^{-w^T x}) + (1-y) \log (1 + e^{w^T x}) \right]$$

FLDA : Fisher Linear discriminant analysis.



$$C_0 = \{ i : y_i < 0 \}$$

$$C_1 = \{ i : y_i > 0 \}$$

$$\begin{aligned} & w^T x_i \\ & \underline{w^T x_i} \\ & w_0 x_{i0} + w_1 x_{i1} + \dots + w_n x_{in} \end{aligned}$$

$$m_0 = \sum_{i \in C_0} w^T x_i$$

$$m_1 = \sum_{i \in C_1} w^T x_i$$

$$\overline{|C_0|}$$

$$\overline{|C_1|}$$

$$s_0^2 = \frac{\sum_{i \in C_0} (w^T x_i - m_0)^2}{|C_0|} \rightarrow \text{Variance of class 0}$$

$$s_1^2 = \frac{\sum_{i \in C_1} (w^T x_i - m_1)^2}{|C_1|} \rightarrow \text{Variance of class 1}$$

$$J(w) = \frac{(m_1 - m_0)^2}{s_0^2 + s_1^2}$$

$$= \frac{w^T S_B w}{w^T S_W w}$$

$$S_B = \underbrace{(M_1 - M_0)}_{d \times 1} \underbrace{(M_1 - M_0)^T}_{1 \times d} \leftarrow d \times d \text{ Matrix between class}$$

$$M_1 = \frac{\sum_{i \in C_1} x_i}{|C_1|}$$

$$M_0 = \frac{\sum_{i \in C_0} x_i}{|C_0|}$$

$$M_1, M_0 \in \mathbb{R}^d$$

$$S_w = \sum_{i \in C_0} (x_i - M_0)(x_i - M_0)^T + \sum_{i \in C_1} (x_i - M_1)(x_i - M_1)^T$$

$$J(w) = \frac{w^T S_B w}{w^T S_w w}$$

$$w^* = \max_w J(w)$$

$$w^* \propto S_w^{-1} (M_1 - M_0)$$


Duda &
Hart
Pattern Classification
Chapter 3.

Logistic Regression:

$$h: \mathbb{R}^d \rightarrow \mathbb{R} \xrightarrow{G} [0, 1]$$

$$h(x) = w^T x$$

$$G(z) = \frac{1}{1 + e^{-z}}$$



$$P(y=1|x) = \frac{1}{1+e^{-w^T x}} = \sigma(w^T x)$$

$$\mathcal{D} = \{ (x_i, y_i)_{i=1}^n \}$$

Christoph
Bishop
chapter 3/4
Linear models

$$P(y=0|x) = 1 - \frac{1}{1+e^{-w^T x}} = 1 - \sigma(w^T x)$$

$$= \frac{e^{-w^T x}}{1+e^{-w^T x}} = \frac{1}{1+e^{w^T x}}$$

$\{0, 1\}$
 $\{-1, 1\}$

$$\underbrace{P(y=y_1|x_1)}_{\text{Perceptron}} = \frac{1}{1+e^{-y_1 w^T x_1}}$$

$$P(y=1|x_1) = \frac{1}{1+e^{-w^T x_1}}$$

$$P(y=-1|x_1) = \frac{1}{1+e^{w^T x_1}}$$

$$\prod_{i=1}^n P(y_i|x_i)$$

$$P(D) = \prod_{i=1}^n P(y_i | x_i)$$

$$J(w) = \prod_{i=1}^n \frac{1}{1 + e^{-y_i w^T x_i}}$$

$$w^* = \arg \max_w J(w) = \arg \max_w P(D)$$

$$\tilde{J}(w) = \log J(w) = \log \prod_{i=1}^n \frac{1}{1 + e^{-y_i w^T x_i}}$$

$$= \sum_{i=1}^n -\log(1 + e^{-y_i w^T x_i})$$

$$l(w) = \sum_{i=1}^n \log(1 + e^{-y_i w^T x_i}) = -\tilde{J}(w)$$

Logistic regression:

$$l(w) = \sum_{i=1}^n \log(1 + e^{-y_i w^T x_i})$$

$y_i \in \{-1, 1\}$

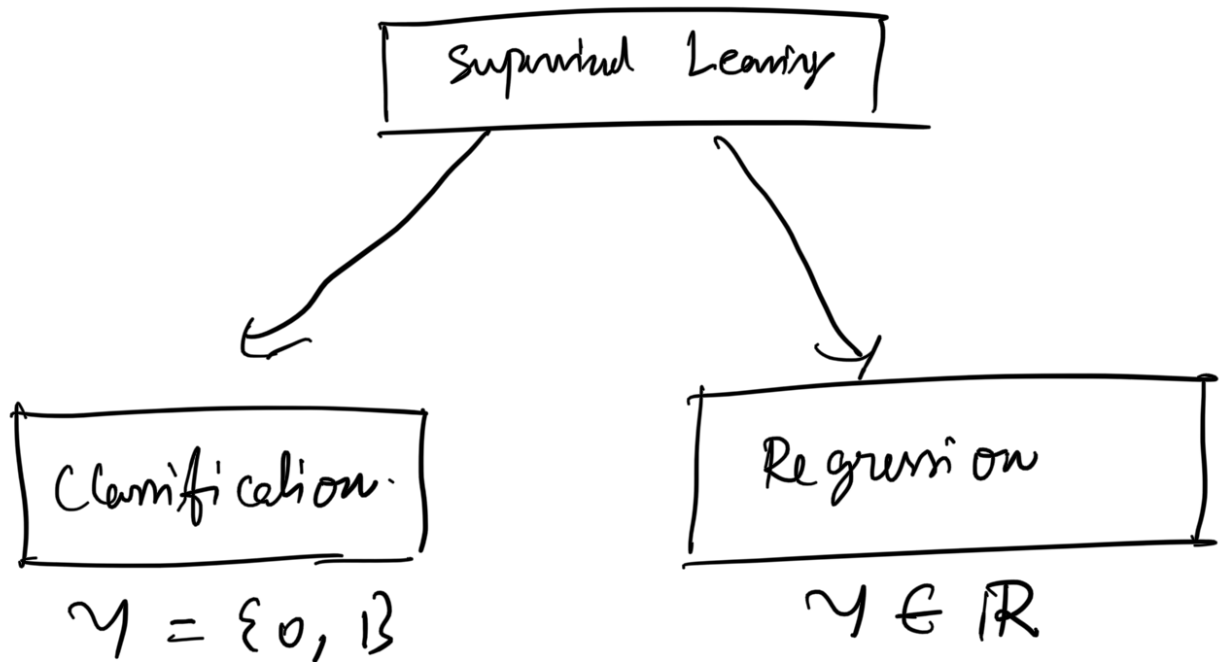
$$L(w) = -\sum_{i=1}^n y_i \log h(x_i) + (1 - y_i) \log(1 - h(x_i))$$

$$L(w) = -l(w)$$

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

Classifiers: (Discriminative)

1. Linear regression — closed form, iterative
2. Logistic regression — iterative
3. FLDA — closed form



(x_1, y_1)

(x_1, y_2)

$x_1, y_1, x_2, y_2, x_3, y_3,$

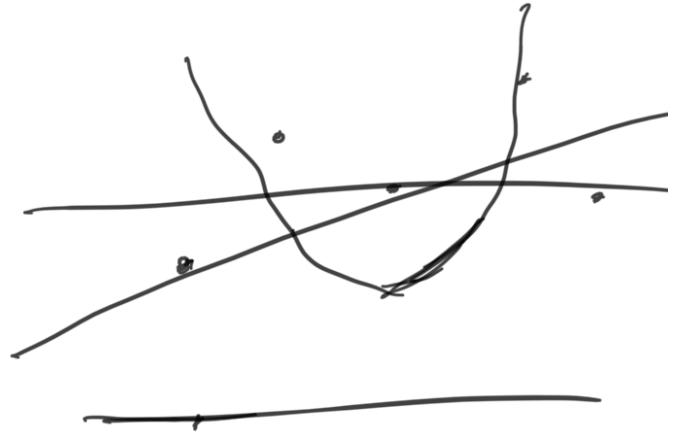
$x_4, y_4, x_5, y_5 \in \mathbb{R}^1$

$$(x_3, y_3)$$

$$(x_4, y_4)$$

$$(x_5, y_5)$$

$$(x, y)$$



$$\hat{y} = \underline{w_1} x + \underline{w_0} \rightarrow \text{Model } M=1$$

$$\hat{y} = w_2 x^2 + w_1 x + w_0 \rightarrow M=2$$

$$\hat{y} = w_3 x^3 + w_2 x^2 + w_1 x + w_0 \rightarrow M=3.$$

$$y = w_0 \rightarrow M=0$$

$$\sum_{i=1}^n (\hat{y}_i - y_i)^2 \Leftrightarrow J(w)$$