

## Example 1) ML Classification task

Pattern  $\rightarrow$  Alg  $\rightarrow$  label

Notation :

①  $\mathcal{A}$  : Feature Space = Set of all feature vectors  
For e.g.  $\mathcal{A} = \mathbb{R}^d$

② Classifier  $h: \mathcal{A} \rightarrow \{1, \dots, M\}$

Simple case:  $M=2$ ,  $\{0, 1\}$  or  $\{-1, +1\}$

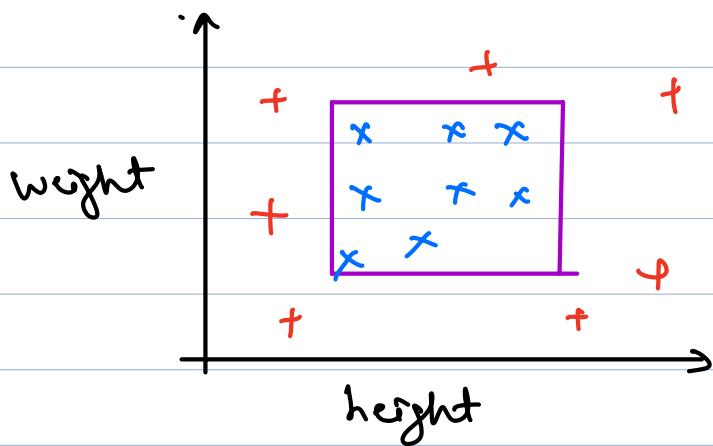
How to design classifiers & how to judge their performance?

Input:  $\{(x_i, y_i), i=1 \dots n\}$  Training dataset  
 $x_i$   $\nearrow$  feature vector  $y_i$   $\searrow$  class label

Using training data, learn an appropriate classifier  $h$

Test! Test & validate the  $h$  on "new" data

Example: Medium-build



Regression task:

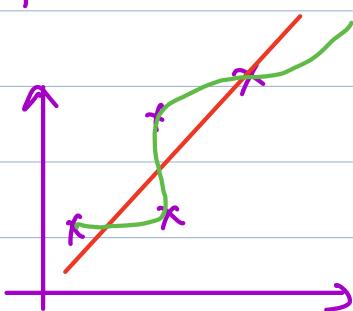
Training data:  $\{(x_i, y_i), i=1 \dots n\}$   
 $x_i \in \mathbb{R}^d, y_i \in \mathbb{R}$

Need is to learn a function  $h: X \rightarrow \mathbb{R}$

Examples: prediction of stock prices, etc

Curve-fitting:  $\{(x_1, y_1), \dots, (x_n, y_n)\}$

Find a function  $f$  s.t.  $y = f(x)$



Generalization:- Obtain a classifier/regression function

using a training dataset s.t. the error on "unseen" (test) data is low.

Assumption: There is a distribution underlying the data.

### Linear models

Ref: Bishop's book

Max-likelihood & least-squares

$$(x) \rightarrow y = \theta^T x + \epsilon, \text{ where } \epsilon \sim N(0, \frac{1}{\beta})$$

↑  
unknown  
indep & identically  
distributed  
precision parameter

Suppose  $\mathcal{D}_n = \{(x_i, y_i), i=1-n\}$  iid & satisfying (x).

$$L(\theta) = \prod_{i=1}^n \frac{\sqrt{\beta}}{\sqrt{2\pi}} \exp\left(-\frac{\beta}{2} (y_i - \theta^T x_i)^2\right)$$

$$l(\theta) = \log L(\theta) = \frac{n}{2} \log \beta - \frac{n}{2} \log 2\pi - \beta \left[ \frac{1}{2} \sum_{i=1}^n (y_i - \theta^T x_i)^2 \right]$$

To maximize  $l(\theta)$ , it is enough to find

$$\hat{\theta}_{ML} = \arg \min_{\theta} \frac{1}{2} \sum_{i=1}^n (y_i - \theta^T x_i)^2 \leftarrow \begin{matrix} \text{empirical risk} \\ \text{minimization} \end{matrix}$$

Linear

Regression: Given data  $\{(x_i, y_i); i=1 \dots n\}$

$$x_i \in \mathbb{R}^d, y \in \mathbb{R}$$

$$J(\theta) = \frac{1}{2} \sum_{i=1}^n (x_i^\top \theta - y_i)^2$$

Minimize  $J$ :

$$A = \begin{bmatrix} -x_1^\top - \\ -x_2^\top - \\ \vdots \\ -x_n^\top - \end{bmatrix}_{n \times d}$$

$$A\theta = \begin{bmatrix} x_1^\top \theta \\ \vdots \\ x_n^\top \theta \end{bmatrix} \quad y = \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix}$$

$$A\theta - y = \begin{bmatrix} x_1^\top \theta - y_1 \\ \vdots \\ x_n^\top \theta - y_n \end{bmatrix}$$

$$(A\theta - y)^\top (A\theta - y) = \sum_{i=1}^n (x_i^\top \theta - y_i)^2 = 2J(\theta)$$

$$\nabla J(\theta) = \frac{1}{2} \nabla (A\theta - y)^\top (A\theta - y) = A^\top (A\theta - y)$$

$$\text{So, } \nabla J(\theta) = 0 \Leftrightarrow (A^\top A)\theta = A^\top y$$

$$\nabla J(\theta) = 0 \Leftrightarrow (A^T A)\theta = A^T y$$

Can we write  $\theta = (A^T A)^{-1} A^T y$ ?

Yes, if  $A$  is full rank.

since full rank  $A = A^T A$  is invertible

(Why? argue  $\text{rank}(A) = \text{rank}(A^T A)$  by showing  
 $N(A) = N(A^T A)$ )

## Lecture-4: A popular ML algorithm for training

Want to solve  $\min_{\theta} f(\theta) = \frac{1}{m} \sum_{i=1}^m f_i(\theta)$

Common assumption:  $f_i$  smooth &  $f$  is convex/  
strongly-convex

Batch GD:

$$\Theta_{n+1} = \Theta_n - \alpha_n \left( \frac{1}{m} \sum_{i=1}^m \nabla f_i(\Theta_n) \right) \leftarrow \text{Nonscler algorithm}$$

Computationally expensive for large  $m$   
(e.g.  $m \rightarrow$  # training examples in a  
large dataset)

Alternative: SGD

Pick " $i_n$ " uniformly at random in  $\{1, \dots, m\}$

$$i_n = \begin{cases} 1 & \text{w.p. } \frac{1}{m} \\ \vdots & \\ m & \text{---} \end{cases}$$

$$\Theta_{n+1} = \Theta_n - \alpha_n \nabla f_{i_n}(\Theta_n) \xrightarrow{\substack{\text{One sample gradient} \\ \text{update}}}$$

"Specializing SGD for linear model":

Linear regression:

$$f(\theta) = \frac{1}{2m} \sum_{i=1}^m (x_i^\top \theta - f_i)^2 \quad \text{--- (1)}$$

Assume  $A$  has full col. rank

$$\text{Let } A = \frac{1}{m} \sum_{i=1}^m x_i x_i^\top, \quad b = \frac{1}{m} \sum_{i=1}^m y_i x_i$$

Then, minimizer  $\hat{\theta}$  of (1) is  $\hat{\theta} = A^{-1} b$

$A \rightarrow d \times d$  matrix

Incremental inverse computation:  $\mathcal{O}(d^2m)$

$$\text{Woodbury's Identity: } (A + BCD)^{-1} = A^{-1} - A^{-1}B(C + DA^{-1}B)^{-1}DA^{-1}$$

$$\text{Sherman-Morrison Lemma: } (A + uv^T)^{-1} = A^{-1} - \frac{A^{-1}uv^TA^{-1}}{1 + v^TA^{-1}u}$$

SHD alternative:

$$\theta_{k+1} = \theta_k - \alpha_k \nabla f_{i_m}(\theta_k)$$

pick a sample  $(x_{i_m}, y_{i_m})$  wif. at random from  
 $\{ (x_i, y_i), i=1 \dots n \}$

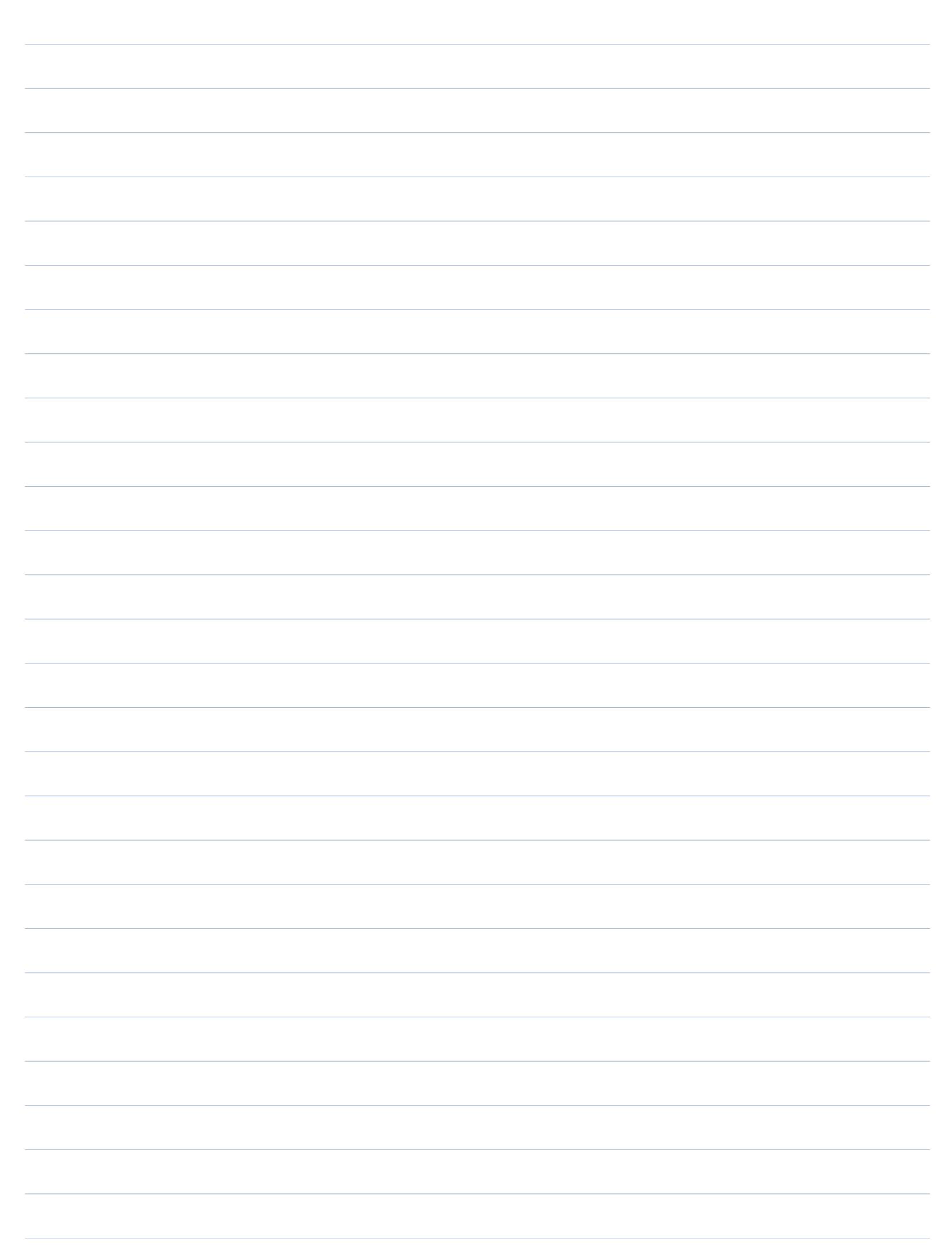
$$\theta_{k+1} = \theta_k + \alpha_k (y_{i_m} - \theta_k^T x_{i_m}) x_{i_m}$$

$$\Theta_k \rightarrow \hat{\Theta} \text{ as } k \rightarrow \infty.$$

(\*) is a LSA algorithm.

$$\frac{1}{m} \sum_{i=1}^m (y_i - \theta^\top x_i) x_i = 0 \quad (\Rightarrow) \quad A\hat{\theta} = b$$

Take expectation & equate to



## Some notation before example 2:

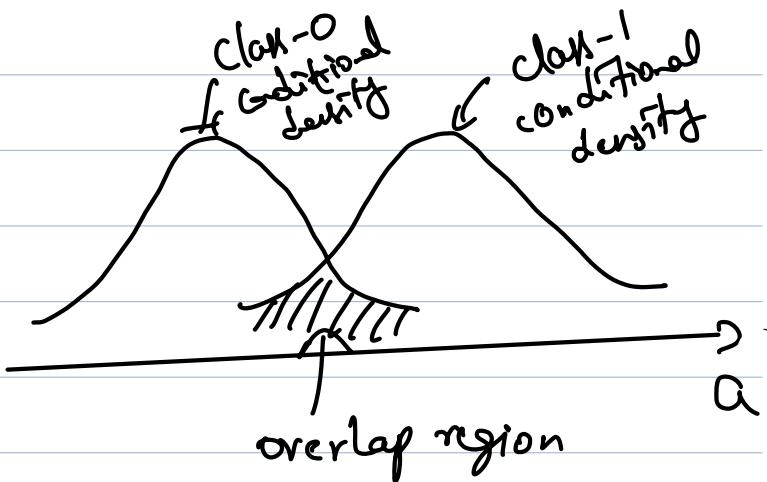
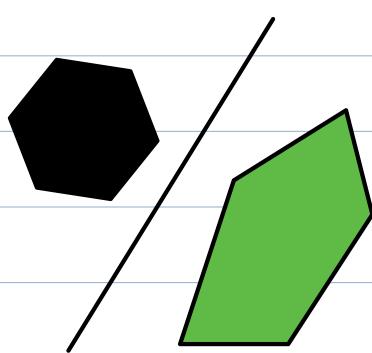
Formally,

$f_i$ : conditional density of features from class - i

Let  $\mathbf{x} = (x_1, \dots, x_d) \in \mathbb{R}^d$  represent a feature vector

$f_i(\mathbf{x})$ : joint density of  $(x_1, \dots, x_d)$  given that  $\mathbf{x}$  is in class - i

Note!: A feature vector " $\mathbf{x}$ " can belong to different classes with different probabilities



## Example 2: Logistic Regression

Linear models for classification

Ref: Sec 4.2 of Bishop's book

Probabilistic Generative model:

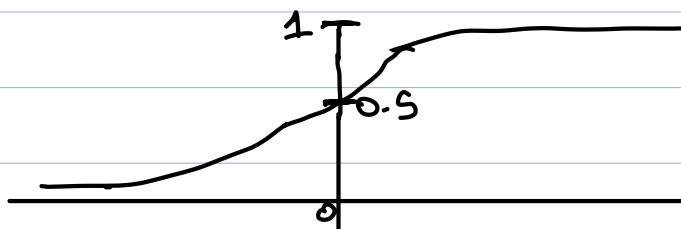
~~Bayes' Classifier~~ Model the class-conditional densities & prior  
Then, compute the posterior

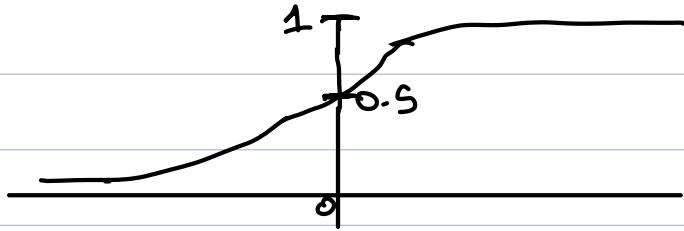
Two-class classification problem:-

$$q_1(x) = \frac{f_0(x)p_0}{f_0(x)p_0 + f_1(x)p_1}$$

$$= \frac{1}{1 + \frac{f_1(x)p_1}{f_0(x)p_0}}$$

$$q_1(x) = \frac{1}{1 + \exp(-w)}, \quad w = \log \frac{f_0(x)p_0}{f_1(x)p_1}$$
$$= \sigma(w)$$





- Note: ①  $\frac{d\sigma(w)}{da} = \sigma(w)(1-\sigma(w))$  ← check this
- ②  $\sigma(-w) = 1 - \sigma(w)$
- ③ Logit function:  $w = \log\left(\frac{\sigma}{1-\sigma}\right)$

Classification task is easier if "a" has a simple functional form, such as "linear".

Example:- Class conditional densities are Gaussian with the same covariance matrix

$$f_i(x) = \frac{1}{(2\pi)^{d/2}} \frac{1}{\sqrt{|\Sigma|}} \exp\left(-\frac{1}{2} (x - \mu_i)^T \Sigma^{-1} (x - \mu_i)\right), \quad i=0, 1$$

Calculating the sigmoid parameter "a" in this case:

$$a = \log \frac{f_0(x)}{f_1(x)} + \log \frac{p_0}{p_1}$$

$$\begin{aligned} \log \frac{f_0(x)}{f_1(x)} &= -\frac{1}{2} (x - \mu_0)^T \Sigma^{-1} (x - \mu_0) \\ &\quad + \frac{1}{2} (x - \mu_1)^T \Sigma^{-1} (x - \mu_1) \\ &= x^T \Sigma^{-1} (\mu_0 - \mu_1) - \frac{1}{2} \mu_0^T \Sigma^{-1} \mu_0 + \frac{1}{2} \mu_1^T \Sigma^{-1} \mu_1 \end{aligned}$$

$$= \omega^T x + \omega_0, \text{ where}$$

$$\begin{aligned} \omega &= \Sigma^{-1} (\mu_0 - \mu_1), \text{ and} \\ \omega_0 &= -\frac{1}{2} \mu_0^T \Sigma^{-1} \mu_0 + \frac{1}{2} \mu_1^T \Sigma^{-1} \mu_1 \end{aligned}$$

$$\text{So, } q_{V_i}(x) = \sigma(\omega^T x + \omega_0)$$


---

Max-likelihood estimation:

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

$\uparrow$   
class labels  $\in \{0, 1\}$

$y_i = +1$   
if class-0

To estimate:  $p_0, \mu_0, \mu_1, \Sigma$

For a data point from class-0,

$$f(x_i, \text{class-0}) = p_0 f_0(x_i)$$

$$\text{Similarly, } f(x_i, \text{class-1}) = (1-p_0) f_1(x_i)$$

Likelihood function:

$$L(p_0, \mu_0, \mu_1, \Sigma) = \prod_{i=1}^n \left( p_0 f_0(x_i) \right)^{y_i} \left( (1-p_0) f_1(x_i) \right)^{1-y_i}$$

$f_0 \sim N(\mu_0, \Sigma)$ 
 $f_1 \sim N(\mu_1, \Sigma)$

Estimate  $p_0$ :

$$\text{log-likelihood: } l(p_0, \mu_0, \mu_1, \Sigma) = \log L(p_0, \mu_0, \mu_1, \Sigma)$$

In  $l(p_0, \mu_0, \mu_1, \Sigma)$ , the terms that depend on  $p_0$  are as follows:-

$$\nabla_{p_0} \left( \sum_{i=1}^n y_i \log p_0 + (1-y_i) \log (1-p_0) \right) = 0$$

$$p_0 = \frac{1}{n} \sum_{i=1}^n y_i = \frac{n_0}{n_0+n_1}, \text{ where}$$

$n_0 = \# \text{ of points in class-0}$

$n_1 = \# \text{ of points in class-1}$

Estimate  $\mu_0$ :

$$\nabla_{\mu_0} [l(p_0, \mu_0, \mu_1, \Sigma)] = 0$$

$$\nabla_{\mu_0} \left( -\frac{1}{2} \sum_{i=1}^n y_i (x_i - \mu_0)^T \Sigma^{-1} (x_i - \mu_0) \right) = 0$$

$$\mu_0 = \frac{1}{n_0} \sum_{i=1}^n y_i x_i$$

$y_i = 1 \text{ if } x_i \in \text{class-0.}$

Similarly  $\mu_1 = \frac{1}{n_1} \sum_{i=1}^n (1-y_i) x_i$

$$\text{Estimating the covariance matrix } \Sigma:$$

The terms involving  $\Sigma$  in the log-likelihood function are

$$\begin{aligned} & \left( -\frac{1}{2} \sum_{i=1}^n y_i \log |\Sigma| - \frac{1}{2} \sum_{i=1}^n y_i (x_i - \mu_0)^T \Sigma^{-1} (x_i - \mu_0) \right. \\ & \quad \left. - \frac{1}{2} \sum_{i=1}^n (1-y_i) \log |\Sigma| - \frac{1}{2} \sum_{i=1}^n (1-y_i) (x_i - \mu_1)^T \Sigma^{-1} (x_i - \mu_1) \right) \\ & = -\frac{n}{2} \log |\Sigma| - \frac{n}{2} \text{Tr}(\Sigma^{-1} S), \text{ where} \end{aligned}$$

$$S = \frac{n_0}{n} S_0 + \frac{n_1}{n} S_1, \quad \text{---} \circled{1}$$

$$S_0 = \frac{1}{n_0} \sum_{i \in \text{class-0}} (x_i - \mu_0) (x_i - \mu_0)^T$$

$$S_1 = \frac{1}{n_1} \sum_{i \in \text{class-1}} (x_i - \mu_1) (x_i - \mu_1)^T$$

The ML-estimate of  $\Sigma$  is " $S$ "

Note: In the univariate case, it is easy to infer the  $S$  in  $\circled{1}$  is the ML estimate.

$$+ n = -\frac{1}{2} \sum y_i (x_i - \mu_0)^2 \times \left(-\frac{1}{S_0}\right)$$

IS

$$-\frac{1}{2} \sum_{i=1}^n (y_i - \mu_1)^2 + (-\frac{1}{2}) \sum_{i=n+1}^m (x_i - \mu_2)^2$$

$$S = \frac{1}{n} \left[ \sum_{i=1}^n (y_i - \mu_1)^2 + \sum_{i=n+1}^m (x_i - \mu_2)^2 \right]$$

$$S = \left( \frac{n_0}{n} \sum_{i \in \text{class 0}} (x_i - \mu_0)^2 + \frac{n_1}{n} \sum_{i \in \text{class 1}} (x_i - \mu_1)^2 \right)$$
$$= \frac{n_0}{n} S_0 + \frac{n_1}{n} S_1, \text{ where } S_0 = \frac{1}{n_0} \sum_{i \in \text{class 0}} (x_i - \mu_0)^2$$

$$S_1 = \frac{1}{n_1} \sum_{i \in \text{class 1}} (x_i - \mu_1)^2$$

## Probabilistic discriminative models

We have seen that

- (i) For a 2-class classification problem, the posterior probability of each class can be written as a logistic sigmoid acting on a linear function of the feature vector.

ML approach can be adopted to estimate density parameters.

- (ii) Alternatively, we the functional form of the generalized linear model & determine the parameters (weights) directly.

$$q_0(x) = \sigma(\omega^T x), \quad q_1(x) = 1 - q_0(x)$$

ML-approach:-  $x \in \mathbb{R}^d$

Take Gaussian CDF,

# parameters to estimate are

(i)  $\mu_0, \mu_1 \Rightarrow 2d$  parameters

(ii)  $\Sigma$  (symmetric)  $\Rightarrow \frac{d(d+1)}{2}$  parameters

(iii) prior prob  $p_0$

Overall ML approach requires  $O(d^2)$  parameters to estimate

If we use the generalized linear model & estimate parameters directly, then the complexity of the task reduces to optimizing  $O(d)$  # of parameters.

## Logistic regression:

Dataset  $\{(x_i, y_i), i=1\dots n\}$   $y_i \in \{0, 1\}$

$$\text{likelihood } L(\omega) = \prod_{i=1}^n \left( \sigma(\omega^T x_i) \right)^{y_i} \left( 1 - \sigma(\omega^T x_i) \right)^{1-y_i}$$

(cost-entropy)

$$\text{Loss } l(\omega) = -\log \mathcal{L}(\omega)$$

$$l(\omega) = - \left[ \sum_{i=1}^n y_i \log(\sigma(\omega^\top x_i)) + (1-y_i) \log(1-\sigma(\omega^\top x_i)) \right]$$

$$\nabla l(\omega) = - \sum_{i=1}^n \frac{y_i}{\sigma(\omega^\top x_i)} \sigma(\omega^\top x_i)(1-\sigma(\omega^\top x_i)) x_i \\ - \sum_{i=1}^n \frac{(1-y_i)}{(1-\sigma(\omega^\top x_i))} \sigma(\omega^\top x_i)(1-\sigma(\omega^\top x_i)) x_i$$

$$= - \sum_{i=1}^n \left[ y_i (1-\sigma(\omega^\top x_i)) x_i - (1-y_i) \sigma(\omega^\top x_i) x_i \right]$$

$$\nabla l(\omega) = \sum_{i=1}^n (\sigma(\omega^\top x_i) - y_i) x_i$$

The gradient descent update

$$\omega_{k+1} = \omega_k - \alpha_k [\sigma(\omega_k^\top x(k)) - y(k)] x(k)$$

where  $(x(k), y(k))$  could be chosen such that

we cycle through the dataset or  
we could pick  $(x(k), y(k))$  uniformly at random  
from the dataset

Claim:  $l$  is strictly convex

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

## Example 3: Matrix factorization

"Sparse" data

Want: Low-rank matrix estimate

Notation:

$A_j \in n \times p$  matrix,  $j=1, \dots, m$ , "Given"

$X \in n \times p$  matrix "Unknown"

Aim:  $\min_X \frac{1}{2m} \sum_{j=1}^m (\langle A_j, X \rangle - y_j)^2 \rightarrow (\star)$

↓  
trace of  $A_j^T X$

Regularized variant:

$$\min_X \frac{1}{2m} \sum_{j=1}^m (\langle A_j, X \rangle - y_j)^2 + \lambda \|X\|_*$$

↙  
nuclear norm = sum of  
singular values

This is a convex-optimization problem

Another variation: Suppose  $X = L R^T$ ,  $L \in n \times r$  matrix  
 $R \in p \times r$  matrix

(\*) becomes

"r"  $\in$  rank  $r \ll \min(n, p)$

$$\min_{L, R} \frac{1}{2m} \sum_{j=1}^m (\langle A_j, L R^T \rangle - y_j)^2 \rightarrow (\star\star)$$

In (\*\*), # elements =  $(n+p)r \ll np$  (in  $x$ )

(\*\*)  $\leftarrow$  non-convex optimization problem.

Non-negative matrix factorization requires  $L \in R$  entries to be  $\geq 0$ .

### Example 4: optimization in ML training

$$\min_x \sum_{i=1}^m f_i(x) + \lambda R(x)$$

$\uparrow$   
regularizer

$f_i, R \rightarrow$  convex in typical training problems (e.g. linear regression, SVM)

$f_i \rightarrow$  non-convex in deep learning

e.g. Image dataset

$$\{a_j, j=1 \dots m\}$$

$a_j \rightarrow$  n-pixel image

Binary classification problem



Dataset  $\{(a_j, b_j), j=1 \dots m\}$

$$\min_x f(x) = \sum_{j=1}^m l(f_x(a_j), b_j)$$

$$\max(0, 1 - b_j(x^T a_j))$$

Hinge-loss

sum

deep-learning

"non-convex"  
loss function