

Quiz



# Gaussian Discriminant Analysis

Austin Welch

- Generative algorithm
- models  $p(x|y)$  &  $p(y)$
- likelihood                  prior
- can get  $p(y|x)$  from Bayes' Rule
- Posterior
- assumes  $p(x|y)$  is Gaussian

RV  $Z \sim N(\bar{\mu}, \Sigma)$

$$P(z) = \frac{1}{(2\pi)^{d/2} |\Sigma|^{1/2}} \exp \left\{ -\frac{1}{2} (\bar{x} - \bar{\mu})^T \Sigma^{-1} (\bar{x} - \bar{\mu}) \right\}$$

$$\Sigma = E[(\bar{x} - \bar{\mu})(\bar{x} - \bar{\mu})^T]$$

e.g. for  $d=2$ , if  $\Sigma$  = diagonal, variables are uncorrelated

Diagonals shrink/stretch dimensions (sharper/flatter peak)

if  $\Sigma$  not diagonal, variables become correlated which has effect of flattening out the Gaussian in  $x=y$  or  $x=-y$  dimension (X)

( $\Sigma$  is symmetric and positive-semi-definite)

Say  $p(y) = \phi^y (1-\phi)^{1-y}$  (Bernoulli)

then, for  $d=2$ ,

$$P(x|y=0) = \frac{1}{2\pi |\Sigma|^{1/2}} \exp \left( -\frac{1}{2} (x - \mu_0)^T \Sigma^{-1} (x - \mu_0) \right)$$

$$P(x|y=1) = \frac{1}{2\pi |\Sigma|^{1/2}} \exp \left( -\frac{1}{2} (x - \mu_1)^T \Sigma^{-1} (x - \mu_1) \right)$$

And,

$$\begin{aligned} \ell(\phi, \mu_0, \mu_1, \Sigma) &= \log \prod_{i=1}^m P(x^{(i)}, y^{(i)}) \\ &= \log \prod_{i=1}^m P(x^{(i)}, y^{(i)}) p(y^{(i)}) \end{aligned}$$

QDA in words: map  $X$  to the class whose mean vector  $\mu_y$  is closest to it. Here, closeness is measured by a suitably weighted and offset-adjusted distance between  $x$  &  $\mu_y$ .

Linear boundary if  $\Sigma_y$ 's of adjacent classes are equal. (otherwise parabolic)

$$\begin{aligned} h_{MAP}(x) &= \arg \max_{y=1, \dots, m} P(y|x, \theta) = h_{QDA}(x) \\ &= \arg \min_x [-\ln P(x|y, \theta) - \ln P(y)] \\ &= \arg \min_y \left[ \frac{1}{2} (x - \mu_y)^T \Sigma_y^{-1} (x - \mu_y) + \frac{1}{2} \ln \det(\Sigma_y) - \ln P(y) \right] \end{aligned}$$

multivariate Gaussian

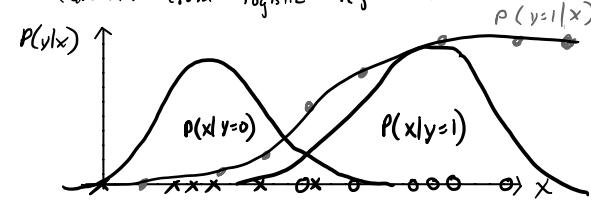
Class-dependent quadratic in  $x$

remember  $\mu_y$  is a vector!

Note:  $\sqrt{(x-\mu)^T \Sigma^{-1} (x-\mu)}$  is the Mahalanobis distance ( $m-d$  std's from mean)

Connection to logistic regression

when you assume  $P(x|y)$  is Gaussian in GDA, if you plot the probabilities of  $P(y|x)$ , you end up with the sigmoid function from logistic regression.



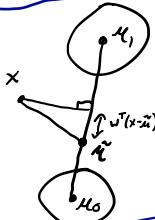
But, GDA will choose different position & steepness for the sigmoid than logistic regression.

$X|y \sim \text{Gaussian} \Rightarrow$  logistic posterior for  $P(y|x)$

$\Sigma = \Sigma_y$  (stronger assumption)

LDA:  $\Sigma_y = \Sigma$  for all  $y$   
Then quadratic rule becomes linear rule:  
 $h_{LDA}(x) = \arg \max_y [B_y^T x + \gamma_y]$

$$\begin{aligned} \text{2-class LDA} \\ h_{LDA}(x) &= 2(g(x) > r) \\ g(x) &= W^T x \\ r &= W^T \bar{\mu} \end{aligned}$$



$w^T x$ : inner product of 2 vectors is proportional to the correlation between their components

LDA in words: Map  $x$  to the class whose mean vector  $\mu_y$  is maximally correlated with it. Here, the correlation is measured by a suitably weighted and offset-adjusted correlation between the components of  $\mu_y$  and  $x$ .

Posterior PMF:  
 $P_{LDA}(y|x, \theta) = \frac{e^{\theta^T x + \eta_y}}{\sum_i e^{\theta^T x + \eta_i}} = \text{Softmax}(y|\eta)$

where  $\eta = (B_1^T x + \gamma_1, \dots, B_m^T x + \gamma_m)^T = (\eta_1, \dots, \eta_m)^T$

Unbiased Estimates

$$\begin{aligned} QDA: \hat{\Sigma}_y &\rightarrow \frac{n_y}{n_y - 1} \hat{\Sigma}_y \\ LDA: \hat{\Sigma} &\rightarrow \frac{n}{n - 2} \hat{\Sigma} \end{aligned}$$

$$\begin{aligned} \hat{N}_y &= \sum_{j=1}^n \mathbb{1}(y_j = y) \\ \hat{P}(y) &= \hat{N}_y / n, \quad \hat{\mu}_y = \frac{1}{\hat{N}_y} \sum_{j \in \{y\}} x_j \\ \text{QDA } \hat{\Sigma}_y &= \frac{1}{\hat{N}_y} \sum_{j \in \{y\}} (x_j - \hat{\mu}_y)(x_j - \hat{\mu}_y)^T \\ \text{LDA } \hat{\Sigma} &= \sum_{y=1}^m \hat{P}(y) \hat{\Sigma}_y \end{aligned}$$

SAME AS  
ML ESTIMATES!

RDA Key Idea  
A: non-invertible (positive semi-definite)  
B: symmetric, positive-definite (so B is invertible)  
Then for  $C = A + B$ ,  
 $C$  is symmetric pos. def. and therefore, invertible.

Special case with  $\lambda=1$ :  $\hat{\Sigma}_{reg} = \text{diag}(\hat{\Sigma}_{ML})$

Equivalent to assuming that all  $d$  components are independent.  $\Rightarrow$  "Gaussian Naive Bayes"

Reg strategy #1:  $\hat{\Sigma}_{reg} = \lambda I_d + (1-\lambda)(\hat{\Sigma}_{ML})$

Reg strat. #2:  $\hat{\Sigma}_{reg} = \lambda \text{diag}(\hat{\Sigma}_{ML}) + (1-\lambda)(\hat{\Sigma}_{ML})$

Need for Regularization:  
Overfitting: Complexity of model is greater than the amount of data, d, n  
can still be ill-conditioned when  $d < n$

**Bayesian**

- $\theta$ : parameters
- $x$ : evidence
- $p(x)$ : evidence
- $p(\theta)$ : priors
- $p(x|\theta)$ : likelihood
- $p(\theta|x)$ : posterior

Generative

$$P(x, \theta) = P(x|\theta)P(\theta)$$

discriminative

$$P(\theta|x)$$
 directly.

Posterior  $\propto$  like. likelihood  $\cdot$  prior (Proportional), NOT equal

$$p(\theta|x) \propto p(x|\theta)p(\theta)$$

Cov ( $x, y$ ) =  $E[(x - \mu_x)(y - \mu_y)]$   
 $= E[xy] - E[x]E[y]$

Cov  $\rightarrow$  uncorrelated  
 $\text{Var}(x) = E[x^2] - E[x]^2$   
 $\text{Cov}(x, x) = \text{Var}(x)$   
 conditional expectation  
 $E[x|Y=y] = \int x f_{XY}(x|y) dx$

Frequentist

$$\hat{\theta} = \underset{\theta}{\operatorname{argmax}} p(D|\theta)$$

$$= \underset{\theta}{\operatorname{argmax}} \left[ \prod_{j=1}^n p(x_j, y_j | \theta) \right]$$

$$= \underset{\theta}{\operatorname{argmin}} \left[ \frac{1}{n} \sum_{j=1}^n -\ln p(x_j, y_j | \theta) \right]$$

Bayesian

$$\hat{\theta} = \underset{\theta}{\operatorname{argmax}} P(\theta|D)$$

$$= \underset{\theta}{\operatorname{argmax}} \left[ P(D|\theta)\pi(\theta) \right]$$

$$= \underset{\theta}{\operatorname{argmin}} \left[ \frac{1}{n} \sum_{j=1}^n -\ln p(x_j, y_j | \theta) - \frac{1}{n} \ln \pi(\theta) \right]$$

Regression

- Supervised learning
- label = quantity
- $y \in Y = \mathbb{R}$
- $\ell(x, y, h) := (y - h(x))^2$
- Risk =  $E[(y - h(x))^2] = \text{MSE}$

$h_{\text{Bayes}}(x) = h_{\text{MLE}}(x) = E[Y|x=x]$

Ordinary Least Squares Regression

$$\underline{x} \in X = \mathbb{R}^d$$

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

$$\theta = \begin{pmatrix} w \\ b \end{pmatrix} \in \mathbb{R}^{d+1}$$

$$\hat{\theta}_{\text{OLS}} = (\hat{w}_{\text{OLS}}, \hat{b}_{\text{OLS}})^T = \underset{(w, b)}{\operatorname{argmin}} \sum_{j=1}^n (y_j - w^T x_j - b)^2$$

Solution

$$\hat{w}_{\text{OLS}} = (\tilde{\Sigma}_x)^{-1} \hat{\Sigma}_{xy}$$

$$\hat{b}_{\text{OLS}} = \hat{w}_y - \hat{\Sigma}_{yx} (\tilde{\Sigma}_x)^{-1} \hat{w}_x$$

Equivalent Expressions

$$\hat{w}_x = \frac{1}{n} [\underline{x}_1 \dots \underline{x}_n]^T$$

$$\hat{w}_y = \frac{1}{n} \tilde{\Sigma}_x^T$$

$$\hat{\Sigma}_x = \frac{1}{n} \tilde{\Sigma} \tilde{\Sigma}^T$$

$$\hat{\Sigma}_{xy} = \frac{1}{n} \tilde{\Sigma} \tilde{\Sigma}^T$$

$$\hat{\Sigma} = [y_1 - \hat{w}_y, \dots, y_n - \hat{w}_y]^T = \tilde{\Sigma} B$$

$$\hat{w}_{\text{OLS}} = (\hat{\Sigma} \hat{\Sigma}^T)^{-1} \tilde{\Sigma} \tilde{\Sigma}^T$$

$$\hat{b}_{\text{OLS}} = \frac{1}{n} P \cdot I_n - (\hat{w}_{\text{OLS}})^T \frac{1}{n} \tilde{\Sigma} I_n$$

Generative

$$P(x, \theta) = P(x|\theta)P(\theta)$$

discriminative

$$P(\theta|x)$$
 directly.

Posterior  $\propto$  like. likelihood  $\cdot$  prior (Proportional), NOT equal

$$p(\theta|x) \propto p(x|\theta)p(\theta)$$

Cov ( $x, y$ ) =  $E[(x - \mu_x)(y - \mu_y)]$   
 $= E[xy] - E[x]E[y]$

Independence

$$p(x, y) = p(x)p(y)$$

$$p(x|y) = p(x)$$

Conditional Covariance

$$\text{Cov}(Y|X=x) = \text{Cov}(Y, Y) - \text{Cov}(Y, X) \text{Cov}(X, X)^{-1} \text{Cov}(X, Y)$$

In this problem, then  $\text{Cov}(Y, Y+x) = \text{Cov}(Y, Y) + \text{Cov}(Y, x)$   
 and likelihood  $p(x|y) = N(y, \sigma_x^2)$

Risk of Decision Rule is exp. val. of loss

$$R(h) = E[\ell(x, y, h)]$$

Risk is expected loss

$$R_{\text{emp}}(h) = \frac{1}{n} \sum_{j=1}^n \ell(x_j, y_j, h) \xrightarrow{n \rightarrow \infty} E[\ell(x, y, h)] = R(h)$$

for ex.  $\ell = (y - h(x))^2$

E [ $\ell(h(x), y)$ ] (needs  $p(x, y)$ )

$$= \int \ell(h(x), y) p(x, y) dx dy$$

R Bayes: min possible risk over all decision rules when  $p(x, y)$  is known

h Bayes: Best decision rule when  $p(x, y)$  is known.

$R(h_{\text{Bayes}}) = R_{\text{Bayes}}$  (Bayes decision rule attains Bayes Risk)

Best in-family decision rule minimizes the in-family posterior risk.

$h_H(x) = \underset{h \in H}{\operatorname{argmin}} E[\ell(x, y, h) | X=x]$  Posterior risk given  $X$

$h_{\text{Bayes}}(x) = \underset{y}{\operatorname{argmin}} E[\ell(p, \hat{y}) | X=x]$

MMSE  $\ell(x, y, h) = \|y - h(x)\|^2$

$h_{\text{Bayes}}(x) = h_{\text{MAP}}(x) = \text{posterior mean} = E[Y|x=x]$

Bayes risk for MMSE is  $E[\|y - h(x)\|^2 | X=x]$

MMSE  $\ell(x, y, h) = \|y - h(x)\| \rightarrow$  (expected val of loss-fn)

$h_{\text{Bayes}} = h_{\text{MMSE}} = \text{Posterior median} = \underset{h_{\text{MAP}}(x)}{\operatorname{argmin}} \ell(x, y, h)$  More compactly,

$R_{\text{Bayes}} = E[\|y - \hat{y}\|^2 | X=x]$

$\hat{\theta}_{\text{OLS}} = \begin{pmatrix} \hat{w}_{\text{OLS}} \\ \hat{b}_{\text{OLS}} \end{pmatrix} = \frac{(\tilde{\Sigma} \tilde{\Sigma}^T)^{-1} \tilde{\Sigma} \tilde{\Sigma}^T}{(d+1)n}$

ex.  $d \times n = 1 \times 50$  ( $HW$ )

$x_j \rightarrow \begin{pmatrix} x_j \\ 1 \end{pmatrix}$  then  $\tilde{\Sigma} \rightarrow \tilde{\Sigma}_{\text{ext}} = \begin{pmatrix} x_1 & x_2 & \dots & x_n \\ 1 & 1 & \dots & 1 \end{pmatrix}$

$\tilde{\Sigma} = (y_1 \dots y_n)$

Whenever you see an optimization problem which looks like:

$$\underset{w}{\operatorname{argmin}} \sum_{i=1}^n (\hat{y}_i - w^T \tilde{\Sigma}_i)^2$$

Ridge Regression ( $\ell_2$ -regularized)

$$(\hat{w}_{\text{ridge}}, \hat{b}_{\text{ridge}}) := \underset{(w, b)}{\operatorname{argmin}} (\hat{w}_{\text{ridge}}, \hat{b}_{\text{ridge}})^T (\hat{w}_{\text{ridge}}, \hat{b}_{\text{ridge}})$$

$$\hat{w}_{\text{opt}} = (\tilde{\Sigma} \tilde{\Sigma}^T)^{-1} \tilde{\Sigma} \tilde{\Sigma}^T$$

$$\tilde{\Sigma} = (\tilde{x}_1 \dots \tilde{x}_n)$$

$$\tilde{\Sigma} = (\tilde{y}_1 \dots \tilde{y}_n)$$

N.L.L.  $y_i = w^T x_i + b + \text{noise}$   
 $w \sim N(0, \frac{1}{2I})$

This is equivalent to Bayesian Discriminant Learning with  $p(y | w, b, x) = N(w^T x + b, \frac{1}{2})$   
 $\pi(w) \sim N(0, \frac{1}{2I})$

Solution:

for any  $w$ , best  $b$  is given by  $(\hat{w}_y - w^T \hat{w}_x)$

the cost function (in terms of only  $w$ ) becomes

2)  $\hat{w}_{\text{ridge}} = (\tilde{\Sigma} \tilde{\Sigma}^T + \lambda I_n)^{-1} \tilde{\Sigma} \tilde{\Sigma}^T$

Regularization:  
 Penalize complexity  
 - Solve ill-posed prob.  
 - Reduce overfitting

# Naive Bayes

- Generative algorithm
  - motivating example: spam classifier
  - $d$  features in each feature vector
  - NB assumption: all features are conditionally independent given the class label.
- $$\forall y, P(\mathbf{x}|y, \theta) = P(x_1|y, \theta_1) \dots P(x_d|y, \theta_d)$$
- $$= \prod_{i=1}^d P(x_i|y, \theta_i), \theta = (\theta_0, \dots, \theta_d)^T$$
- and  $\theta_0 = (P(y=1), \dots, P(y=M))^T$

## Why Naive Bayes?

• Low Complexity

$(x_i) \rightarrow$  Binary feature (0/1)  
 $\vdots$   
 $x_d \rightarrow$  Binary feature (0/1)

$O(md)$ , relatively immune to overfitting

$$P(x_1, \dots, x_d | \theta)$$

$$\stackrel{x_i}{\times} \stackrel{d=2}{\times} \stackrel{(4)}{\times} \stackrel{(2)}{\times}$$

$$\stackrel{0=3}{\times} \stackrel{(8)}{\times} \stackrel{(1)}{\times}$$

$$\stackrel{0=3}{\times} \dots \stackrel{(1)}{\times}$$

$$\stackrel{0=3}{\times} \stackrel{(8)}{\times}$$

$$O(md)$$

$$\text{vs. joint of } O(m(2^{d-1}))$$

$\therefore 2^d$  distinct feature vectors

$\stackrel{d=2}{\times} \stackrel{2}{\times} \stackrel{2}{\times}$   
 $\# \text{ of params needed is growing exponentially fast}$

## Gaussian Naive Bayes

$$\forall i, y, P(x_i|y, \theta_i) \sim N(\mu_{iy}, \sigma_{iy}^2)$$

Categorical NB: All features are discrete

$$\text{N.B. assumption: } P(X_j|y, \theta_j) = \prod_{i=1}^d P(x_{ij}|y, \theta_{ij})$$

## Categorical NB

Bayes Classifier for  $\delta=1$  loss (MAP rule = MPE rule)

$$h_{\text{MPE}}(\mathbf{x}) = \arg \max_{y=1, \dots, M} \prod_{i=1}^d P(x_{i,y} | \theta_{iy})$$

$$h_{\text{MAP}}(\mathbf{x}) = \arg \max_{y=1, \dots, M} \left[ \sum_{w=1}^W \left( \frac{n_{w,y,i}}{n_y} \right) \ln \beta_{w,y,i} \right]$$

$n_y$ : # classify train exs

$n_{w,y,i}$ : # classify train exs where  $i$ th feature =  $w$

$$\forall w, y, i, \hat{\beta}_{w,y,i} = \frac{n_{w,y,i}}{n_y}$$

Notation

$y$  = class label 1, ...,  $M$

$j$  = sample index 1, ...,  $n_y$

$i$  = feature index 1, ...,  $d$

$w$  = feature value 1, ...,  $W$

## Document Classification

(2 classes; msa)

Vocabulary size =  $W = 3$  e.g. {a, b, c}

class 1 examples

$x_1$	$x_2$
a	b
b	c
c	a
a	b
c	c
c	a
doc 1	
length = 7	

length = 5

likelihoods  $P(\mathbf{x}|y)$

$$\hat{\beta}_{a,1} = \frac{2+1}{7+5} = \frac{3}{12}$$

$$\hat{\beta}_{b,1} = \frac{2+2}{12} = \frac{4}{12}$$

$$\hat{\beta}_{c,1} = \frac{3+2}{12} = \frac{5}{12}$$

$$\hat{p}(y=1) = \frac{2}{2+3} = \frac{2}{5}$$

Priors

$$\hat{p}(y=2) = \frac{3}{2+3} = \frac{3}{5}$$

Training set

c	c	c
c	a	b
a	c	c
b	c	c
c	a	a
L=5		

class 2 examples

c	c	c
c	c	b
a	c	c
b	c	c
c	a	a
L=3		

L=4

Test document

b	b	a	c	a	b	c
b	a	c	c	c	b	a
a	c	c	a	b	a	c
c	c	a	b	c	c	a
c	a	b	c	a	c	b
j=12						
p(class 1)				p(class 2)		

Assumption:  
 $P(b|1) = P(b|2)$   
 $P(a|1) = P(a|2)$   
 $P(c|1) = P(c|2)$

Don't mult. by 0 if d not in train set

OR can do a Bayesian Regularization

(Laplace smoothing - default counts)

Then compare with

$$\hat{p}(y=1) \cdot \hat{p}(x|y=1) = \hat{p}(y=2) \cdot \hat{p}(x|y=2)$$

$$= \frac{2}{5} \cdot \hat{\beta}_{a,1} \cdot \hat{\beta}_{b,1} \cdot \hat{\beta}_{c,1} = \frac{3}{5} \cdot (\hat{\beta}_{a,2})^3 \cdot (\hat{\beta}_{b,2})^4 \cdot (\hat{\beta}_{c,2})^3$$

## Dirichlet Prior

length = 7

length = 5

likelihoods  $P(\mathbf{x}|y)$

$$\hat{\beta}_{a,1} = \frac{2+1}{7+5} = \frac{3}{12}$$

$$\hat{\beta}_{b,1} = \frac{2+2}{12} = \frac{4}{12}$$

$$\hat{\beta}_{c,1} = \frac{3+2}{12} = \frac{5}{12}$$

$$\hat{p}(y=1) = \frac{2}{2+3} = \frac{2}{5}$$

Priors

$$\hat{p}(y=2) = \frac{3}{2+3} = \frac{3}{5}$$

Example: Binary classification,  $y = \text{lifestyle}$  Peaceful / Stressed Binary (0,1)

### Feature Vector

Values

$x = \begin{pmatrix} \text{Income} \\ \text{Charitable giving} \\ \text{Avg life sleep per day} \end{pmatrix} \rightarrow \text{Low/middle/high ... Ternary, say Categorical over } \{0, 1, 2\}$   
 $\rightarrow \text{No/Yes ... Binary, say Bernoulli over } \{0, 1\}$   
 $\rightarrow 2-10 \dots \text{Real values, say Gaussian (mean)}$

$d=3$

### Training data

	$y=0$ (peaceful)	$y=1$ (stressed)
(1)	$\begin{pmatrix} 2 \\ 2 \\ 2 \end{pmatrix}$	$\begin{pmatrix} 0 \\ 2 \\ 2 \end{pmatrix}$
(2)	$\begin{pmatrix} 1 \\ 1 \\ 1 \end{pmatrix}$	$\begin{pmatrix} 0 \\ 0 \\ 1 \end{pmatrix}$
(3)	$\begin{pmatrix} 9 \\ 8 \\ 6.5 \end{pmatrix}$	$\begin{pmatrix} 4 \\ 5 \\ 7 \end{pmatrix}$

## Naive Bayes

$$\theta_{0,y} = p(y) \quad \hat{p}_{ML} = \frac{\text{sample mean}}{14} = \hat{p}_{ML}(1)$$

First feature:

$$\theta_{1,y} = \text{pmf over } \{0, 1, 2\}$$

for class '0'  $\theta_{0,1}, \theta_{0,2}, \theta_{0,3}$

for class '1'  $\theta_{1,1}, \theta_{1,2}, \theta_{1,3}$

$\hat{\theta}_{0,1} = \frac{1}{3}$

$\hat{\theta}_{0,2} = \frac{2}{3}$

$\hat{\theta}_{0,3} = \frac{1}{3}$

$\hat{\theta}_{1,1} = \frac{2}{3}$

$\hat{\theta}_{1,2} = \frac{2}{3}$

$\hat{\theta}_{1,3} = \frac{3}{3}$

$\hat{p}_{ML} = \frac{1}{3}$

$\hat{p}_{ML} = \frac{2}{3}$

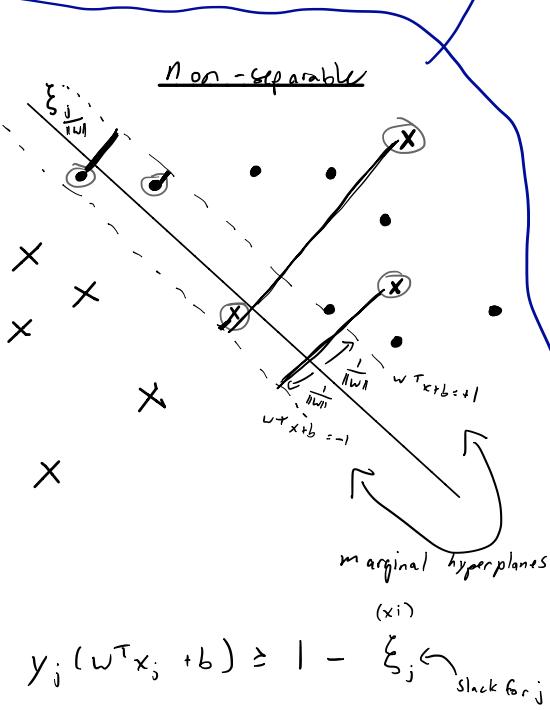
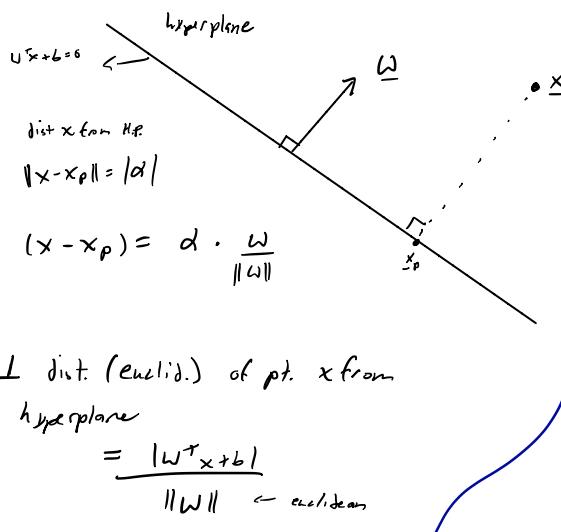
$\hat{p}_{ML} = \frac{2}{3}$

$\hat{p}_{ML} = \frac{1}{3}$

# SVM: Find max margin hyperplane (convex optimization problem)

Equation of hyperplane in  $\mathbb{R}^d$ :  $w^T x + b = 0$

$(w, b)$  and  $(\gamma w, \gamma b)$  define the same hyperplane if  $\gamma \neq 0$



objective function of SVM:

$$(w_{SVM}, b_{SVM}) = \underset{(w, b)}{\operatorname{argmin}} \left[ \frac{1}{2} \|w\|^2 + \text{penalty}(\xi_1, \dots, \xi_n) \right]$$

$$\forall j \quad y_j (w^T x_j + b) \geq 1 - \xi_j, \quad \xi_j \geq 0$$

$\{x_j : \alpha_j > 0\} \rightarrow$  Support Vectors

$\{x_j : \xi_j > 0\} \rightarrow$  Outliers

All outliers are support vectors

→ from complementary slackness condition #2:

$$\xi_j^* > 0 \Rightarrow \alpha_j^* = \xi_j^* > 0$$

Support Vectors →  $\alpha_j^* \neq 0$

## Two-class Linearly Separable SVM

- Label space  $Y = \{-1, +1\}$
- Feature vectors  $X = \mathbb{R}^d$
- Training set  $D = \{(x_1, y_1), \dots, (x_n, y_n)\}$
- Family of classifiers: "linear" classifiers

$$H = \{h(x) = \operatorname{sign}(w^T x + b) \quad w \in \mathbb{R}^d, b \in \mathbb{R}\}$$

Parameters

$$\Theta = \begin{pmatrix} w \\ b \end{pmatrix} \quad (d+1) \text{ scalar parameters}$$

$$w^T x_p + b = 0$$

$$\Rightarrow w^T (x - x_p) = d \frac{w^T w}{\|w\|} = d \frac{\|w\|^2}{\|w\|} = d \|w\|$$

$$\Rightarrow d = \frac{w^T w - w^T x_p}{\|w\|} = \frac{w^T x + b}{\|w\|} \Rightarrow |d| = \frac{|w^T x + b|}{\|w\|}$$

- Geometric margin of a hyperplane  $(w, b)$  wrt  $D$

$$\rho(w, b) := \min_{1 \leq j \leq n} \frac{|w^T x_j + b|}{\|w\|}$$

(distances of feature vector which is closest to the hyperplane)

- Linearly Separable training set  $D$ :

$$\exists \text{ exists } (w, b) \text{ st. } y_j (w^T x_j + b) > 0 \quad \forall j = 1, \dots, n$$

$$\Rightarrow \boxed{\rho(w, b) > 0}$$

Canonical Parameterization of hyperplane  $(w, b)$  wrt lin. sep.  $D$ :

If  $D$  is lin. sep. by  $(w, b)$  and  $\rho > 0$ , then the canonical parameterization  $(w_{\text{canon}}, b_{\text{canon}})$  is such that

$$\min_{1 \leq j \leq n} |w_{\text{canon}}^T x_j + b_{\text{canon}}| = 1$$

$$(w_{SVM}, b_{SVM}) = \underset{(w, b)}{\operatorname{argmax}} \frac{1}{\|w\|}$$

- $y_j (w^T x_j + b) > 0 \quad \forall j$
- $\min_{1 \leq j \leq n} |w^T x_j + b| = 1$

$$\rho_{\text{opt}} = \rho(w_{SVM}, b_{SVM}) = \max_{(w, b)} \frac{1}{\|w\|}$$