

# Machine Learning I: Supervised Methods

B. Keith Jenkins

## Announcements

- Homework 8 is due today
- Project Reports and code are due Friday 4/26
  - See piazza for late submission policy
- Final exam is Wed., May 1

## Today's lecture

- Summary and orientation for estimating  $p(x | S_k)$  for classification
- Parameter estimation for classification
  - Problem set-up and notation
- Ad hoc estimates
- View 1: parameters are deterministic
  - Maximum likelihood estimate (MLE)
- View 2: parameters are random
  - MAP estimate
  - Prior on parameters
  - Use in a classifier (2 approaches)
    - MAP as point estimate
    - Full Bayesian (integrating)
- d.o.f. and constraints in parameter estimation

## Summary so far (probabilistic classification) (2-class case)

L22: Bayes dec. rule for 2-class:

$$p(\underline{x} | S_1) P(S_1) \stackrel{r_1}{\geq} p(\underline{x} | S_2) P(S_2)$$

$\downarrow$

$$\frac{N_1}{N}$$

$\downarrow$

$$\frac{N_2}{N}$$

Estimate priors  $\hat{P}(S_i)$

Ex: frequency of occurrence  $\hat{P}(S_i) =$

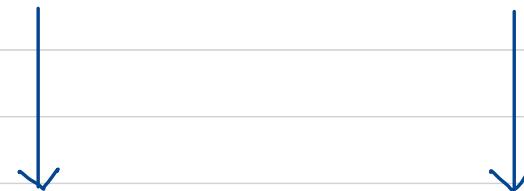
L23: Estimate  $p(\underline{x} | S_i)$

1. Density estimation

Ex: • kNN

• KDE

$$\hat{p}(\underline{x} | S_1) \quad \hat{p}(\underline{x} | S_2)$$



$$f_1(\underline{x}, \hat{\underline{\theta}}_{MLE}^{(1)})$$

$$f_2(\underline{x}, \hat{\underline{\theta}}_{MLE}^{(2)})$$

$$f_1(\underline{x}, \hat{\underline{\theta}}_{MAP}^{(1)})$$

$$f_2(\underline{x}, \hat{\underline{\theta}}_{MAP}^{(2)})$$

$$p(\underline{x} | S_1, \underline{z}_1)$$

$$p(\underline{x} | S_2, \underline{z}_2)$$

↑  
data labelled  $S_1$

↑  
data labelled  $S_2$

$f_k(\underline{x}, \underline{\theta}^{(k)}) =$  known or assumed functional form of  $p(\underline{x} | S_k)$ .

## PARAMETER ESTIMATION

### Motivation:

1. Reduce d.o.f. by making more assumptions
2. Include any knowledge of  $p(\underline{x} | S_k)$

(1) Assume:  $p(\underline{x} | S_k) = f_k(\underline{x}, \underline{\theta}_k) \leftarrow$  functional form of  $f_k$   
is known or assumed.  
 $\underline{\theta}_k$  are parameters that are unknown.

$$\text{Ex: } p(\underline{x} | S_k) = N(\underline{x} | \underline{m}_k, \underline{\Sigma}_k) \triangleq f_k(\underline{x}, \underline{\theta}_k)$$

with  $\underline{m}_k$  unknown, and/or  $\underline{\Sigma}_k$  unknown.  
→ Unknown parameters  $\underline{\theta}_k$  are learned from the training data.

### Problem set-up and notation

Let  $\hat{\underline{\theta}}_k$  be an estimate of  $\underline{\theta}_k$ . (parameters of  $p(\underline{x} | S_k)$ )

Let  $\underline{x}_1^{(k)}, \underline{x}_2^{(k)}, \dots, \underline{x}_n^{(k)}$  be the data pts., drawn i.i.d. from  $p(\underline{x} | S_k)$ .

Let  $\underline{\mathbf{x}} = [\underline{x}_1 \ \underline{x}_2 \ \underline{x}_3 \cdots \underline{x}_n]^T$  denote all training data

$\underline{\mathbf{x}}_k = [\underline{x}_1^{(k)} \ \underline{x}_2^{(k)} \ \cdots \ \underline{x}_{n_k}^{(k)}]^T$  denote training data  
labeled  $S_k$ .

## Class independence assumption

Assume: To estimate  $\underline{\theta}_k$ , we only need data  $\underline{z}_k$  from  $S_k$ .

When considering estimates for only 1 class  $S_k$ , will simplify notation:

$$p(\underline{x} | \underline{\theta}_k, S_k) \rightarrow p(\underline{x} | \underline{\theta})$$

$$\underline{x}_m^{(k)} \rightarrow \underline{x}_m, \underline{z}^{(k)} \rightarrow \underline{z}$$

Is  $\underline{\theta}$  deterministic (but unknown), or random?

→ Both views are valid.

View 1:  $\underline{\theta}$  is deterministic

$$\hat{\underline{\theta}} = G(x_1, x_2, \dots, x_n) = G(\underline{x}) = \hat{\underline{\theta}}(\underline{x})$$

Two properties of an estimate  $\hat{\underline{\theta}}$

Unbiased

$$\text{If } E_{\underline{x}} \{ \hat{\underline{\theta}} \} \stackrel{\Delta}{=} \int \hat{\underline{\theta}}(\underline{z}) p(\underline{z}) d\underline{z} = \underline{\theta}$$

Then  $\hat{\underline{\theta}}$  is an unbiased estimate of  $\underline{\theta}$ .

Otherwise,  $\hat{\underline{\theta}}$  is biased.

Consistent

If  $\lim_{n \rightarrow \infty} P \{ \| \hat{\underline{\theta}} - \underline{\theta} \| < \epsilon \} = 1$  for every  $\epsilon > 0$ ,

then  $\hat{\underline{\theta}}$  is consistent.

("convergence in probability")

## Ad hoc estimates

Sample mean

$$\hat{m} \triangleq \frac{1}{n} \sum_{i=1}^n x_i$$

- unbiased?

$$\begin{aligned} E\{\hat{m}\} &= E\left\{\frac{1}{n} \sum x_i\right\} = \frac{1}{n} \sum E\{x_i\} \\ &= \frac{1}{n} \sum_{i=1}^n m = \frac{1}{n}(nm) = m \end{aligned}$$

$\Rightarrow \hat{m}$  is unbiased.

- consistent?

$$\text{Can show } \lim_{n \rightarrow \infty} P\left\{ \|\hat{m} - m\| < \epsilon \right\} = 1$$

$\Rightarrow \hat{m}$  is consistent.

Sample covariance

$$\hat{\Sigma} = a_n \sum_{i=1}^n (x_i - \hat{m})(x_i - \hat{m})^T$$

$$a_n = \frac{1}{n} \quad (\text{biased})$$

$$a_n = \frac{1}{n-1} \quad (\text{unbiased})$$

## Maximum Likelihood Estimate (MLE)

[Bishop pp.26 - middle of p.28]

Let  $p(\underline{z} | \underline{\theta})$  be the likelihood of parameter values  $\underline{\theta}$  based on training data  $\underline{z}$ .

The MLE of  $\underline{\theta}$  is:

$$(2) \quad \hat{\underline{\theta}}_{MLE} \triangleq \underset{\underline{\theta}}{\operatorname{argmax}} \{ p(\underline{z} | \underline{\theta}) \} = \underset{\underline{\theta}}{\operatorname{argmax}} \{ \ln p(\underline{z} | \underline{\theta}) \}$$

$\Rightarrow \hat{\underline{\theta}}_{MLE}$  yields the  $\underline{\theta}$  that maximizes the probability of obtaining the data that was actually observed.

o Criterion fcn.  $\tilde{J}_{MLE}(\underline{\theta}) = -p(\underline{z} | \underline{\theta})$

$$(3) \quad \text{or } J_{MLE}(\underline{\theta}) = -\ln p(\underline{z} | \underline{\theta}) \quad (\text{"negative log likelihood"})$$

$$\begin{aligned} &= -\ln \prod_{i=1}^n p(z_i | \underline{\theta}) \\ &= -\sum_{i=1}^n \ln p(z_i | \underline{\theta}). \end{aligned}$$

Is  $J_{MLE}$  convex?

→ Depends on  $f_k(x, \underline{\theta}_k)$ .

Ex:

Let  $p(\underline{x}|\theta) = N(\underline{x}|\mu, \Sigma)$ ,  $\mu = \text{unknown}$ ,  $\Sigma = \text{known}$ .

$$\begin{aligned}\hat{\theta}_{MLE} &= \underset{\theta}{\operatorname{argmax}} \{ \ln p(\underline{x}|\theta) \} = \underset{\mu}{\operatorname{argmax}} \{ \ln p(\underline{x}_1, \dots, \underline{x}_n | \mu) \} \\ \hat{\mu}_{MLE} &= \underset{\mu}{\operatorname{argmax}} \left\{ \sum_{i=1}^n \ln p(\underline{x}_i | \mu) \right\} \\ &= \underset{\mu}{\operatorname{argmax}} \left\{ \sum_{i=1}^n \underbrace{-\frac{1}{2} \ln(2\pi)}_{\text{constant}} - \underbrace{\frac{1}{2} \ln |\Sigma|}_{\text{constant}} - \underbrace{\frac{1}{2} (\underline{x}_i - \mu)^T \Sigma^{-1} (\underline{x}_i - \mu)}_{\text{quadratic term}} \right\} \\ &\quad + \underbrace{\frac{1}{2} \underline{x}_i^T \Sigma^{-1} \mu - \frac{1}{2} \mu^T \Sigma^{-1} \mu}_{\text{linear term}}\end{aligned}$$

drop(constants of  $\mu$ ): ↑ and  $-\frac{1}{2} \underline{x}_i^T \Sigma^{-1} \underline{x}_i$ .

$$\hat{\mu}_{ML} = \underset{\mu}{\operatorname{argmin}} J, \quad J = -\sum_{i=1}^n \left[ \mu^T \Sigma^{-1} \underline{x}_i - \frac{1}{2} \mu^T \Sigma^{-1} \mu \right]$$

• Is  $J$  convex? yes ✓ ✓

$$\min J: \quad \nabla_{\mu} J(\mu) = 0$$

$$\nabla_{\mu} \{ \cdot \} = \sum_{i=1}^n \left[ \Sigma^{-1} \underline{x}_i - \Sigma^{-1} \mu \right]$$

$$\sum_{i=1}^n \left[ \Sigma^{-1} \underline{x}_i - \Sigma^{-1} \hat{\mu}_{MLE} \right] = 0$$

$$\Sigma \left[ \Sigma^{-1} \left( \sum_i \underline{x}_i \right) \right] = \Sigma \left[ n \Sigma^{-1} \hat{\mu}_{MLE} \right]$$

$$\underbrace{\frac{1}{n} \sum_i \underline{x}_i}_{\hat{\mu}_{MLE}} = \hat{\mu}_{MLE}$$

## MLE in a Bayes (min. error) classifier

Bayes min.error dec. rule:

$$(4) \quad p(\underline{x} | S_k) P(S_k) > p(\underline{x} | S_l) P(S_l) \quad \forall l \neq k \Rightarrow \underline{x} \in S_k$$

$$(5) \quad \text{Given: } p(\underline{x} | S_k) = f_k(\underline{x}, \underline{\theta}_k)$$

$\uparrow$  known       $\uparrow$  unknown

Use:  $\underline{\theta}_k \approx \hat{\underline{\theta}}_{MLE}^{(k)}$ , so that  $\hat{p}(\underline{x} | S_k) = f_k(\underline{x}, \hat{\underline{\theta}}_{MLE}^{(k)})$ .

Bayes min.error decision rule becomes:

$$\hat{p}(\underline{x} | S_k) P(S_k) > \hat{p}(\underline{x} | S_l) P(S_l) \quad \forall k \neq l \Rightarrow \underline{x} \in S_k$$

$$(6) \quad \text{with } \hat{p}(\underline{x} | S_j) = f_j(\underline{x}, \hat{\underline{\theta}}_{MLE}^{(j)}) \quad \forall j.$$

Can also use  $P(S_j) \approx \hat{P}(S_j)$  if  $P(S_j)$  are unknown

## Comments

1.  $\hat{p}(\underline{x} | S_k)$  is obtained by plugging in the estimate (MLE) of  $\underline{\theta}$  into the assumed or known  $p(\underline{x} | S_k) = f_k(\underline{x}, \underline{\theta}_k)$
2. Linear Bayes and quadratic Bayes usually mean Bayes (min.error) classifier with assumed normal densities for  $p(\underline{x} | S_j)$ , with  $\underline{m}_j$  and  $\underline{\Sigma}_j$  estimated from the training data using (usually) maximum likelihood estimates.  
It's linear if  $\underline{\Sigma}_i = \underline{\Sigma}_j \forall i, j$ , and generally quadratic otherwise.
3. Also, linear discriminant analysis (LDA) and quadratic discriminant analysis (QDA) often have similar meaning as linear Bayes and quadratic Bayes  
- except LDA and QDA are often also used as a method of feature-space transformation for dimensionality reduction (similar to FLD and MDA). (The "analysis" term)

## View 2: $\underline{\theta}$ is random

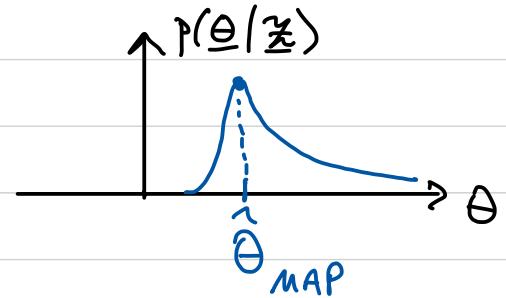
$\underline{\theta}$  has underlying pdf:  $p(\underline{\theta})$ .  
 → Bayesian estimation.

### Maximum a-Posteriori (MAP) estimate

$$(7) \quad \hat{\underline{\theta}}_{MAP} \triangleq \underset{\underline{\theta}}{\operatorname{argmax}} \left\{ p(\underline{\theta} | \underline{z}) \right\} = \underset{\underline{\theta}}{\operatorname{argmax}} \left\{ \ln p(\underline{\theta} | \underline{z}) \right\}$$

$p(\underline{\theta} | \underline{z})$  = posterior density of  $\underline{\theta}$ .

$\hat{\underline{\theta}}_{MAP}$  is the  $\underline{\theta}$  that is most probable, given the data  $\underline{z}$  (mode of  $p(\underline{\theta} | \underline{z})$ ).



Write  $p(\underline{\theta} | \underline{z})$  in terms of  $p(\underline{z} | \underline{\theta})$ :

$$(8) \quad \begin{cases} p(\underline{\theta} | \underline{z}) = \frac{p(\underline{z} | \underline{\theta}) p(\underline{\theta})}{p(\underline{z})} \\ \ln p(\underline{\theta} | \underline{z}) = \ln p(\underline{z} | \underline{\theta}) + \ln p(\underline{\theta}) - \underbrace{\ln p(\underline{z})}_{\text{const. of } \underline{\theta}} \end{cases}$$

$$(9) \quad \hat{\underline{\theta}}_{MAP} = \underset{\underline{\theta}}{\operatorname{argmax}} \left\{ \ln p(\underline{z} | \underline{\theta}) + \ln p(\underline{\theta}) \right\}$$

known  
or assumed  
(MLE)

prior  
on  $\underline{\theta}$

## Prior on $\Theta$

→ Allows us to insert our prior knowledge of the problem.

Ex: Classify tumors:  $S_1 = \text{malignant}$ ;  $S_2 = \text{benign}$ .

Assume you normalize or standardize the data so that

$$0 \leq x_j \leq 1 \quad \forall j \quad \text{or} \quad \sigma_j = 1 \quad \forall j.$$

We have  $D$  features. Let  $\underline{\theta} = \underline{w}$ .

We will use a linear model so  $g(\underline{x}) = \underline{w}^T \underline{x}$  with  $g(\underline{x}) \stackrel{\Gamma_1}{\geq} 0$ .

Obtain priors from domain experts: a panel of radiologists.

Q: On a scale of 0 to 10, how significant is tumor size in deciding malignancy? (0 = insignificant (not relevant), 10 = very important in deciding malignancy)

A: 8

Q: Does a larger tumor size usually mean a higher or lower chance it is malignant?

A: higher

$$\Rightarrow w_1 = +8$$

Q: Same questions for tumor symmetry?

A: 5, and more symmetry means less likely to be malignant.

$$\Rightarrow w_2 = -5$$

Q: Overall, how certain are you of your numeric answers?

Can you give a numeric range - for example a significance of 5 could be a range of 4-6, or 3-7, or?

A: 3-7

$$\Rightarrow \text{Prior } p(\underline{w}) = N(\underline{w} | \underline{m}_w, \Sigma_w), \Sigma_w = \sigma_w^2 \mathbb{I} \text{ (assumptions)}$$

$$(\underline{m}_w)_1 = 8, (\underline{m}_w)_2 = -5, \Sigma_w = 2^2 \mathbb{I}$$

## Bayesian estimation in a Bayes (min.error) classifier

Bayes min.error dec. rule:

$$(4) \quad p(\underline{x} | S_k) P(S_k) > p(\underline{x} | S_\ell) P(S_\ell) \quad \forall \ell \neq k \Rightarrow \underline{x} \in S_k$$

$$(1) \quad \text{Given: } p(\underline{x} | S_k) = f_k(\underline{x}, \underline{\theta}_k) = p(\underline{x} | \underline{\theta}_k, S_k)$$

$\uparrow$  known       $\uparrow$  unknown

Two approaches

1. MAP as a point estimate for  $\underline{\theta}_k$

$$\underline{\theta}_k \approx \hat{\underline{\theta}}_{MAP}^{(k)}$$

(4) becomes:

$$\hat{p}(\underline{x} | S_k) P(S_k) > \hat{p}(\underline{x} | S_\ell) P(S_\ell) \quad \forall \ell \neq k \Rightarrow \underline{x} \in S_k$$

$$(10) \quad \text{with } \hat{p}(\underline{x} | S_j) = f_j(\underline{x}, \hat{\underline{\theta}}_{MAP}^{(j)}) \quad \forall j.$$

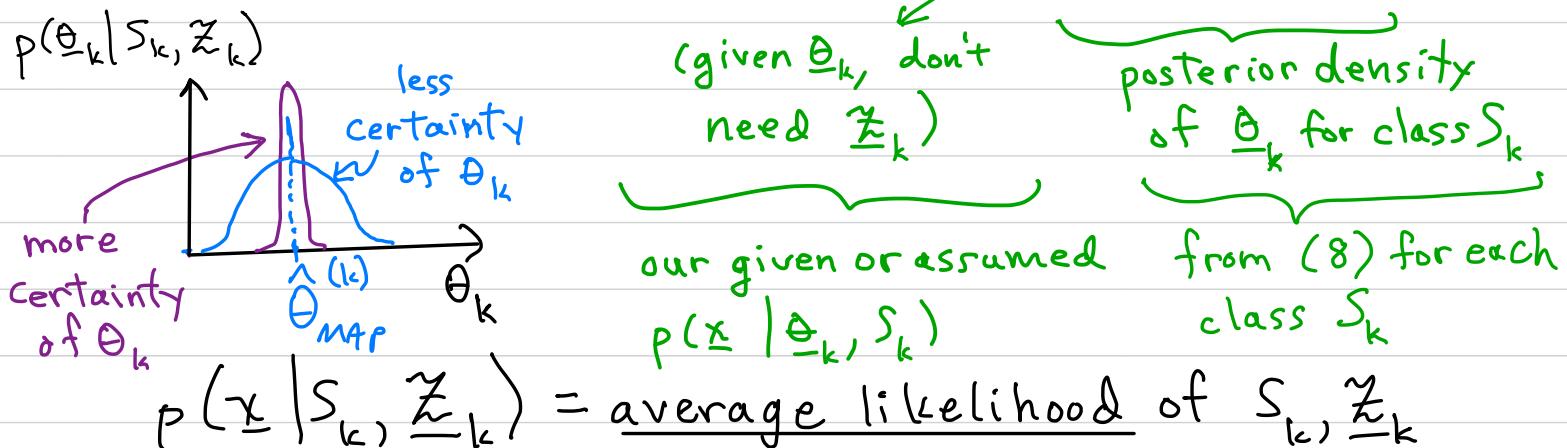
Can also use  $P(S_j) \approx \hat{P}(S_j)$  if  $P(S_j)$  are unknown

## 2. Integrate over parameter posterior

Use posterior density of  $\underline{\theta}$ ,  $p(\underline{\theta} | \underline{z})$

Integrate over  $\underline{\theta}$ , instead of setting  $\underline{\theta} \approx \hat{\underline{\theta}}_{MAP}$   
 $\rightarrow$  Full Bayesian estimate for  $p(\underline{x} | S_k)$

$$(11) \quad p(\underline{x} | S_k, \underline{z}_k) = \int p(\underline{x} | \underline{\theta}_k, S_k, \underline{z}_k) p(\underline{\theta}_k | S_k, \underline{z}_k) d\underline{\theta}_k$$



(4) becomes:

$$\hat{p}(\underline{x} | S_k) P(S_k) > \hat{p}(\underline{x} | S_l) P(S_l) \quad \forall l \neq k \Rightarrow \underline{x} \in S_k$$

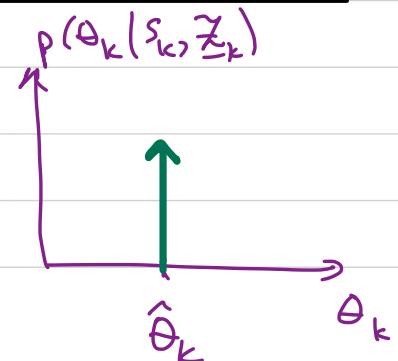
$$(12) \quad \text{with } \hat{p}(\underline{x} | S_j) = p(\underline{x} | S_j, \underline{z}_j) \quad \forall j.$$

Can also use  $P(S_j) \approx \hat{P}(S_j)$  if  $P(S_j)$  are unknown

Note: if  $p(\underline{\theta}_k | S_k, \underline{z}_k) = \delta(\underline{\theta}_k - \hat{\underline{\theta}}_k)$

$$\Rightarrow p(\underline{x} | S_k, \underline{z}_k) = p(\underline{x} | S_k, \hat{\underline{\theta}}_k).$$

(same as approach 1).



## Degrees of freedom and constraints in parameter estimation

Ex1: estimate  $\underline{m}_k$  only. ( $k=1, \dots, C$ ).  $\Rightarrow D$  d.o.f. for each class  
 $(\rho(\underline{x} | S_k), k=1, \dots, C)$ .  
# constraints =  $N_k$  for class  $S_k$   
 $\Rightarrow CD$  total.

$N$  total

[ Ex2: estimate  $\hat{\Sigma}_k$  using sample covariance matrix.

$$\hat{\Sigma}_k = \frac{1}{N_k} \sum_{i=1}^{N_k} (\underline{x}_i - \hat{\underline{m}}_k) (\underline{x}_i - \hat{\underline{m}}_k)^T,$$

$\hat{A}_i$

Each matrix  $\hat{A}_i$  is rank 1.

$\sum_i$  will have rank  $\leq D$

what if  $N < D$ ?

The max. rank of  $\hat{\Sigma}$  is  $N_k - 1$ .

$\Rightarrow$  rank of  $\hat{\Sigma}_k \leq \min\{N_k - 1, D\}$ .

$$\left( \# \text{d.o.f.} \right)_k \leq D (\text{rank of } \sum_k) \leq D \min \{ N_k - 1, D \}.$$

$$\text{Total: } D \sum_{k=1}^c \min \{ N_k - 1, D \} \leq cD^2.$$

#constraints = N. 

$$\text{Est. } m_k \text{ and } \sum_k, k=1, \dots, c \Rightarrow \text{d.o.f.} = cD^2 + cD = \underline{cD(D+1)}$$

(dens-est., histogr. method:  $M^D$ ).  
 $\rightarrow c$  classes  $\Rightarrow c \cdot M^D$