# Machine Learning I: Supervised Methods

B. Keith Jenkins

## Announcements

- No Slido poll questions today

- Homework 7 is due Friday

- Homework 8 will be posted early next week

- Project report template and instructions will be posted soon

## Reading

- Bishop 2.5 (density estimation)

## Today's lecture

- Statistical classification and Bayes Decision Theory

  - Minimum-error classifiers and $P_e$ (C=2)

  - Minimum-error classifiers and $P_e$ (C>2)

  - Summary

  - Minimum-risk criterion

- Mahalanobis distance

- Classifiers: Gaussian density case

  ‣ Linear Bayes (LDA)

  ‣ Quadratic Bayes (QDA)

- Density estimation techniques for machine learning (time permitting)

  - Preliminaries

(1)
$$P_e = P(S_1) \int_{\Gamma_2} p(\underline{x}|S_1)\, d\underline{x} + P(S_2) \int_{\Gamma_1} p(\underline{x}|S_2)\, d\underline{x}$$
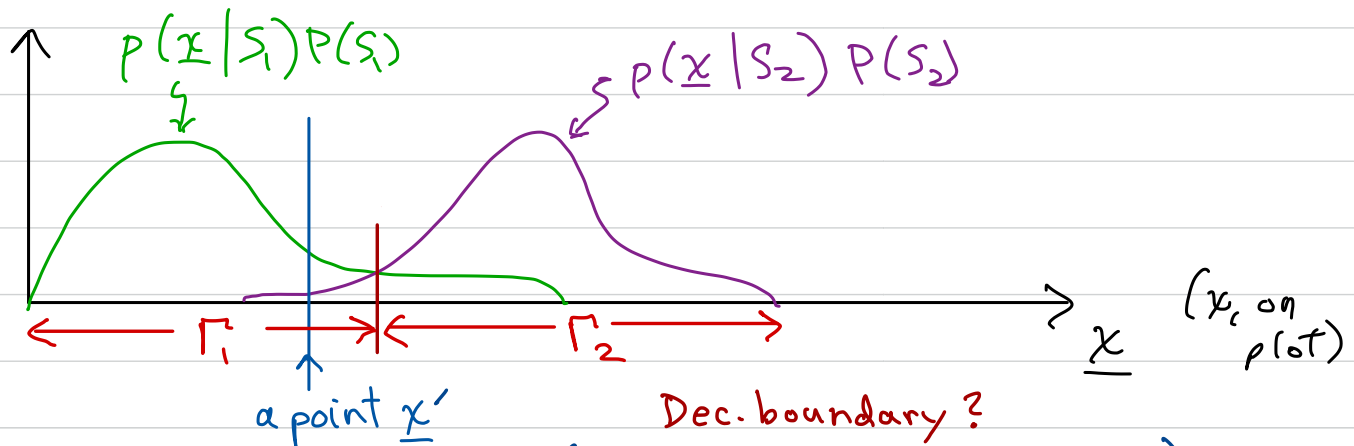
→ $P_e$ is our criterion fcn.

Goal: find $\Gamma_1$ and $\Gamma_2$ that minimize $P_e$.

Requirements: $\Gamma_1$ and $\Gamma_2$ must be non-overlapping.

$\Gamma_1 \cup \Gamma_2$ covers all of feature space.
(except possibly boundaries that have area (or volume) = 0)



$p(\underline{x}|S_1) P(S_1)$

$p(\underline{x}|S_2) P(S_2)$

$\Gamma_1$   $\Gamma_2$

$\underline{x}$   ($x_i$ on plot)

a point $\underline{x}'$           Dec. boundary ?

⇒ assign to $\Gamma_1$ (for min. contribution to $P_e$).

In $P_e$, each pt. $\underline{x}$ must be included in 1 of the 2 terms; pick the smaller term.

⇒ minimizes $P_e$.

∴ Assign $\underline{x}$ to $\Gamma_1$ if $p(\underline{x}|S_2)P(S_2) < p(\underline{x}|S_1)P(S_1)$

Assign $\underline{x}$ to $\Gamma_2$ if $p(\underline{x}|S_2)P(S_2) > p(\underline{x}|S_1)P(S_1)$
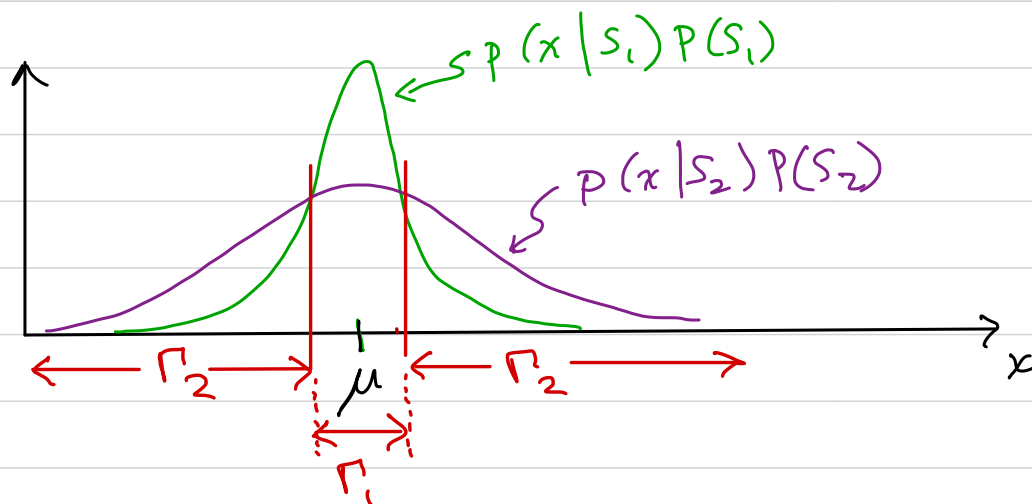
or

(2)
$$p(\underline{x}|S_1)P(S_1) \underset{\Gamma_2}{\overset{\Gamma_1}{\gtrless}} p(\underline{x}|S_2)P(S_2)$$

↳ Bayes decision rule for min. error  ($c=2$)

or $\ln [\cdot] \underset{\Gamma_2}{\overset{\Gamma_1}{\gtrless}} \ln [\cdot]$

Example: suppose $\mu_1 = \mu_2 = \mu$, $\sigma_1 < \sigma_2$ :    (Normal densities $p(\underline{x} | S_i)$)



$\leftarrow p(x|S_1) P(S_1)$

$p(x|S_2) P(S_2)$
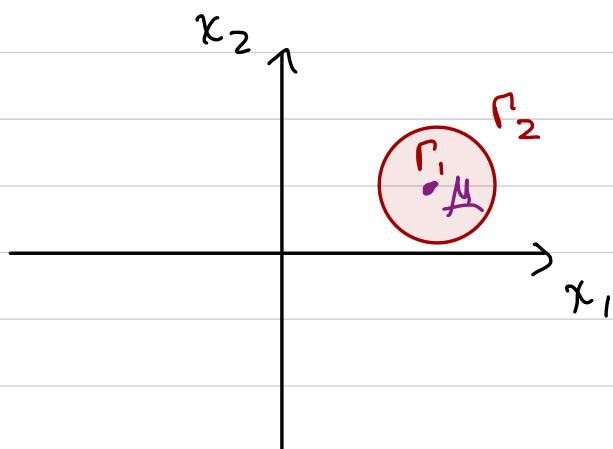
$\Gamma_2 \quad \mu \quad \Gamma_2$

$\Gamma_1$

> Decision boundaries and regions? ↗

>| Is this classifier linear? → No.

2D version: $p(\underline{x}|S_i) = N(\underline{x} | \underline{\mu}, \sigma_i^2 \underline{\underline{I}})$
$P(S_1) = P(S_2)$

$\sigma_1^2 < \sigma_2^2$



$x_2$

$\Gamma_2$

$\Gamma_1$ $\quad \cdot \underline{\mu}$

$x_1$

Use Bayes theorem:

$$P(S_k|\underline{x})P(\underline{x}) = P(\underline{x}|S_k)P(S_k)$$

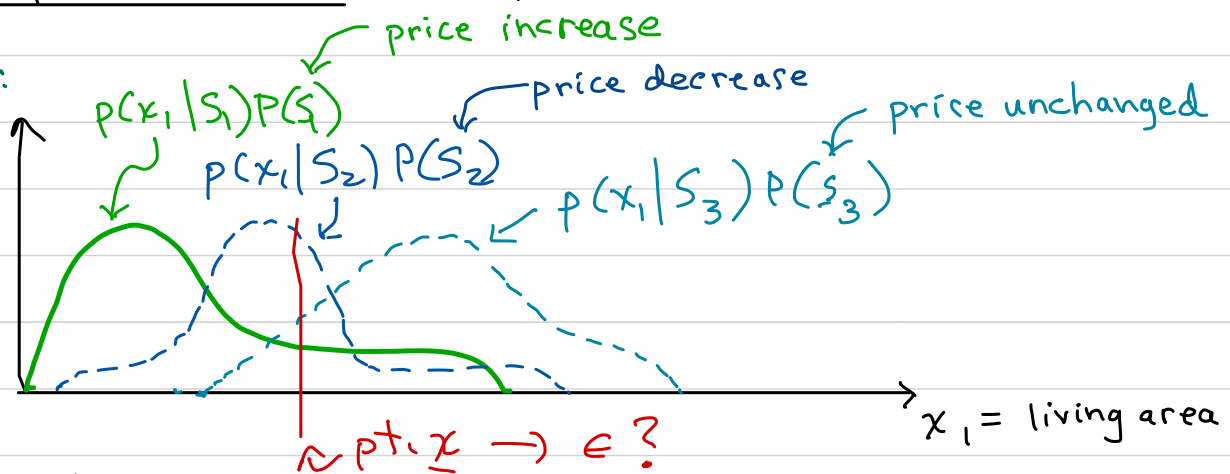$$P(\underline{x}) = \sum_{i=1}^{C} P(\underline{x}|S_i)P(S_i)$$

$$\Rightarrow \quad P(S_1|\underline{x})P(\underline{x}) \underset{\Gamma_2}{\overset{\Gamma_1}{\gtrless}} P(S_2|\underline{x})P(\underline{x})$$

(3)
$$\boxed{P(S_1|\underline{x}) \underset{\Gamma_2}{\overset{\Gamma_1}{\gtrless}} P(S_2|\underline{x})}$$

$\widehat{C}$ Bayes min-error dec-rule, in terms of posterior probabilities.

Bayes min error $(C > 2)$

Ex: $C=3$:



- price increase
- price decrease
- price unchanged

$P(x_1|S_1)P(S_1)$

$P(x_1|S_2)P(S_2)$

$P(x_1|S_3)P(S_3)$

$x_1 = $ living area

$\sim$ pt. $\underline{x} \rightarrow \epsilon$ ?

$$P_e = 1 - P_{correct}$$

$$P_{correct} = \sum_{i=1}^{C} P(S_i) \int_{\Gamma_i} P(\underline{x}|S_i)\, d\underline{x}$$

$$P_e = 1 - P_{correct} = 1 - \left\{ P(S_1)\int_{\Gamma_1} P(\underline{x}|S_1)\,d\underline{x} + P(S_2)\int_{\Gamma_2} P(\underline{x}|S_2)\,d\underline{x} \right.$$
$$\left. + P(S_3)\int_{\Gamma_3} P(\underline{x}|S_3)\,dx_1 \right\}$$

$\Rightarrow$

$g_k(x)$        $g_j(x)$

$$p(\underline{x}|S_k)P(S_k) > p(\underline{x}|S_j)P(S_j) \qquad \forall \, j \neq k$$

$$\Rightarrow \underline{x} \in \Gamma_k$$

Bayes min. error decision rule, $C > 2$.

Is this OvR, OvO, MvM, or something different?

This is a maximal value method.

Can define:

$$g_i(\underline{x}) = p(\underline{x}|S_i)P(S_i)$$

or

could choose $\tilde{g}_i(\underline{x}) = \ln\left[p(\underline{x}|S_i)P(S_i)\right]$

$$= \ln p(\underline{x}|S_i) + \ln P(S_i)$$

# SUMMARY OF BAYES DECISION THEORY SO FAR

1. Bayes minimum error classifier

Decision rule:

$$p(\underline{x}|S_i)P(S_i) > p(\underline{x}|S_j)P(S_j) \quad \forall \, j \neq i \implies \underline{x} \in \Gamma_i$$

2-class case:

$$p(\underline{x}|S_1)P(S_1) \underset{\Gamma_2}{\overset{\Gamma_1}{\gtrless}} p(\underline{x}|S_2)P(S_2)$$

2. Probability of error

$$P_E = 1 - P_{correct} = 1 - \sum_{i=1}^{c} \int_{\Gamma_i} p(\underline{x}|S_i)P(S_i)\, d\underline{x}$$

2-class case:

$$P_e = \int_{\Gamma_2} p(\underline{x}|S_1)P(S_1)\, d\underline{x} + \int_{\Gamma_1} p(\underline{x}|S_2)P(S_2)\, d\underline{x}$$

3. Note: for discrete-valued features $x_k$,

$$\int dx_k \quad becomes \quad \sum_{x_k}$$

and same decision rules apply.

# MINIMUM RISK CRITERION [Bishop 1.5.2]

For cases in which minimizing $P_e$ is not optimal; e.g. misclassifying $S_1$ data pt. as $S_2$ is significantly more costly than vice versa.

Ex: medical test that screens for cancer. $S_1 =$ positive ($\Rightarrow$ cancer)

False negative is worse than false positive. $S_2 =$ negative ($\Rightarrow$ no cancer)
↳ (miss the cancer)

→ Allow different costs for different kinds of error.

Let $L_{ji} =$ loss of assigning $\underline{x}$ to $\Gamma_i$ when it actually belongs to $S_j$.

Total expected loss is then: $E\{L\} = \sum\limits_{j=1}^{c} \sum\limits_{i=1}^{c} \left[ \int\limits_{\Gamma_i} L_{ji} \, p(s_j | \underline{x}) d\underline{x} \right]$

average over all classes

Expected loss of assigning $\underline{x}$ to $\Gamma_i$.

$\Rightarrow$ Instead of $p(S_i | \underline{x}) > p(S_k | \underline{x}) \; \forall k \neq i \Rightarrow \underline{x} \in \Gamma_i$

we have:

Decision rule

$$\sum\limits_{j=1}^{c} L_{ji} \, p(S_j | \underline{x}) < \sum\limits_{j=1}^{c} L_{jk} \, p(S_j | \underline{x}) \quad \forall k \neq i \Rightarrow \underline{x} \in \Gamma_i$$

$R(\alpha_i | \underline{x}) \triangleq$ conditional risk of taking action $\alpha_i$ ($\underline{x} \in \Gamma_i$) given $\underline{x} \in S_j$.

$\underline{\underline{L}} = \begin{bmatrix} L_{11} & L_{12} & \\ L_{21} & L_{22} & \\ & & \ddots \\ & & & L_{cc} \end{bmatrix}$, typically $L_{ii} = 0 \; \forall i$.

# MAHALANOBIS DISTANCE [Bishop 2.3.0]

$$d_M^2(\underline{x}, \underline{m}) = (\underline{x}-\underline{m})^T \underline{\underline{\Sigma}}^{-1} (\underline{x}-\underline{m})$$

$$\underline{\underline{\Sigma}} = \begin{bmatrix} \sigma_{11} & \sigma_{12} & \cdots & \\ \sigma_{21} & \ddots & & \\ \vdots & & & \\ & & & \sigma_{DD} \end{bmatrix}$$

CASE 1 :

$$\underline{\underline{\Sigma}} = \begin{bmatrix} \sigma_1^2 & & & \\ & \sigma_2^2 & & O \\ & & \ddots & \\ O & & & \sigma_D^2 \end{bmatrix}$$

$$\underline{\underline{\Sigma}}^{-1} = \begin{bmatrix} \frac{1}{\sigma_1^2} & & & \\ & \frac{1}{\sigma_2^2} & & O \\ & & \ddots & \\ O & & & \frac{1}{\sigma_D^2} \end{bmatrix}$$

$$d_M^2(\underline{x}, \underline{m}) = \sum_{i=1}^{D} \frac{1}{\sigma_i^2} (x_i - m_i)^2$$

$$d_M^2(\underline{x}, \underline{m}) = \text{const.} \Rightarrow \text{?}$$

2-space: $d_M^2 = \dfrac{(x_1 - m_1)^2}{\sigma_1^2} + \dfrac{(x_2 - m_2)^2}{\sigma_2^2} = $ const.



$\uparrow x_2 = $ dist. cc.

$\sigma_2$

$\sim d_m = 1$

$\underline{m}$

$\sigma_1$

$x_1 = $ living area

$\longrightarrow$ Ellipse (2D)

$\longrightarrow$ Hyperellipsoid (D dim.)

## CASE 2:

$$\underline{\underline{\Sigma}} = \text{general}$$

$$= \begin{bmatrix} \sigma_{11} & \sigma_{12} & - - - \\ \sigma_{21} & \ddots & \\ \vdots & & \\ \vdots & & \sigma_{DD} \end{bmatrix}$$

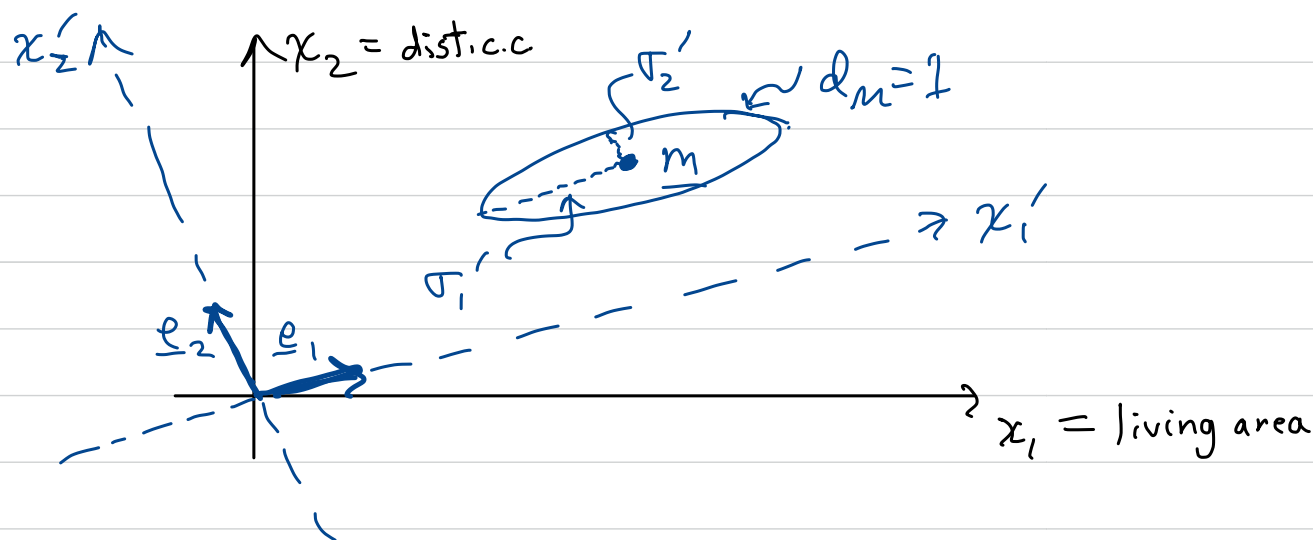$$d_M^2(\underline{x}, \underline{m}) = (\underline{x} - \underline{m})^T \underline{\underline{\Sigma}}^{-1} (\underline{x} - \underline{m})$$

Apply orthonormal transformation (rotate basis):

$$\underline{x}' = \underline{\underline{E}}^T \underline{x}$$

$$\underline{\underline{\Sigma}}' = \underline{\underline{E}}^T \underline{\underline{\Sigma}} \underline{\underline{E}} = \underline{\underline{\Lambda}} = \text{diagonal} \longrightarrow \text{CASE 1.}$$

$$\Rightarrow d_M^2 = \text{const.} \Rightarrow \text{hyperellipsoids (axes rotated)}$$

$x_2' \nearrow$    $\uparrow x_2 = dist.c.c$    $\sigma_2'$  $\sim d_M = 1$

$m$

$\sigma_1'$    $\dashrightarrow x_1'$

$\underline{e}_2$  $\underline{e}_1$

$x_1 = living\ area$

---

BAYES MIN. ERROR CLASSIFIERS — GAUSSIAN DENSITY CASE

[Bishop 2.3.3 — optional reading]

$$p(\underline{x}\,|\,S_k) = N(\underline{x}\,|\,\underline{m}_k,\,\underline{\underline{\Sigma}}_k)$$

$$= \frac{1}{(2\pi)^{D/2}\,|\underline{\underline{\Sigma}}_k|^{1/2}}\,\exp\left\{-\frac{1}{2}\left[(\underline{x}-\underline{m}_k)^T\,\underline{\underline{\Sigma}}_k^{-1}\,(\underline{x}-\underline{m}_k)\right]\right\}$$

$\underset{i}{Maximize}\ \ p(\underline{x}\,|\,S_i)\,P(S_i)$

$$g_i(\underline{x}) = \ln\left\{p(\underline{x}\,|\,S_i)P(S_i)\right\} = \ln p(\underline{x}\,|\,S_i) + \ln P(S_i)$$

(1) $\boxed{\;g_i(\underline{x}) = -\frac{1}{2}\ln|\underline{\underline{\Sigma}}_i| - \frac{1}{2}(\underline{x}-\underline{m}_i)^T\,\underline{\underline{\Sigma}}_i^{-1}\,(\underline{x}-\underline{m}_i) + \ln P(S_i)\;}$
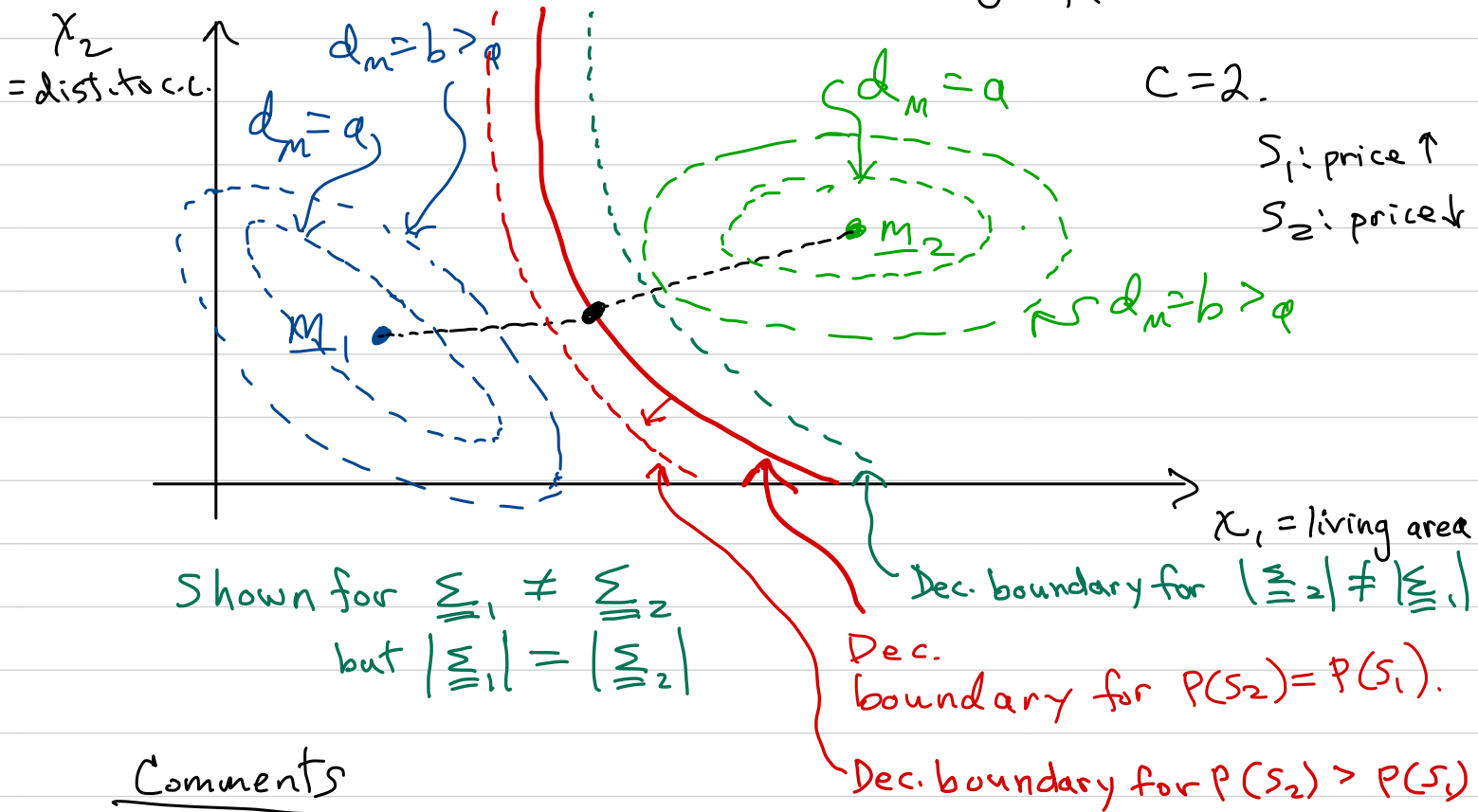
$Let:\ \ |\underline{\underline{\Sigma}}_1| = |\underline{\underline{\Sigma}}_2| = \cdots$

$and:\ \ P(S_1) = P(S_2) = \cdots$

Then:

$$g_i(\underline{x}) = -(\underline{x}-\underline{m}_i)^T \underline{\underline{\Sigma}}_i^{-1} (\underline{x}-\underline{m}_i) = -d_M^2(\underline{x},\underline{m}_i)$$

$\implies$ Nearest-means classifier using $d_M$ instead of $d_E$.



$X_2$ = dist. to c.c.

$d_m = b > a$

$d_m = a$

$d_m = a$

$d_m = b > a$

$\underline{m}_2$

$\underline{m}_1$

$C = 2$.

$S_1$: price $\uparrow$
$S_2$: price $\downarrow$

$X_1$ = living area

Shown for $\underline{\underline{\Sigma}}_1 \neq \underline{\underline{\Sigma}}_2$
but $|\underline{\underline{\Sigma}}_1| = |\underline{\underline{\Sigma}}_2|$

Dec. boundary for $|\underline{\underline{\Sigma}}_2| \neq |\underline{\underline{\Sigma}}_1|$

Dec. boundary for $P(S_2) = P(S_1)$.

Dec. boundary for $P(S_2) > P(S_1)$

Comments

1. Include: $P(S_i) \neq P(S_j)$

$\implies$ Boundary shifts
e.g.: $P(S_2) > P(S_1)$ (see plot)

2. Include: $|\underline{\underline{\Sigma}}_i| \neq |\underline{\underline{\Sigma}}_j|$

$\implies$ incorporates differences in ellipsoid
volumes from class to class.

CASE A  $\underline{\underline{\Sigma}}_i = \underline{\underline{\Sigma}}$

$$g_i(\underline{x}) = -\frac{1}{2}(\underline{x} - \underline{m}_i)^t \underline{\underline{\Sigma}}^{-1}(\underline{x} - \underline{m}_i) + \ln P(S_i)$$

$$= -\frac{1}{2}\left[ \underline{x}^t \underline{\underline{\Sigma}}^{-1}\underline{x} - 2\underline{m}_i^t \underline{\underline{\Sigma}}^{-1}\underline{x} + \underline{m}_i^t \underline{\underline{\Sigma}}^{-1}\underline{m}_i \right] + \ln P(S_i)$$
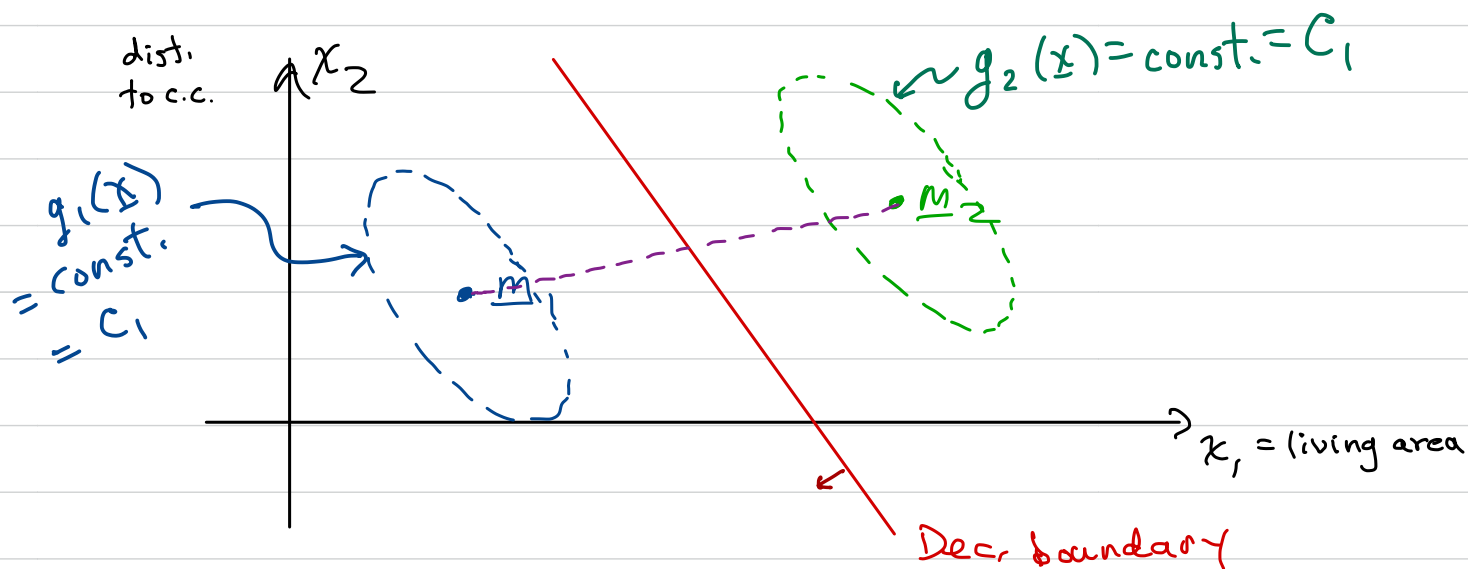
Let $g_i(\underline{x}) = \underline{m}_i^T \underline{\underline{\Sigma}}^{-1}\underline{x} - \frac{1}{2}\underline{m}_i^t \underline{\underline{\Sigma}}^{-1}\underline{m}_i + \ln P(S_i)$

$$= \underline{w}^T \underline{x} + w_0$$

↑ is this classifier linear? → Yes!

$\Rightarrow$ Classifier is linear.

Often called **Linear Bayes**.

dist. to c.c.

$g_1(\underline{x})$ = const. = $C_1$

$\sim g_2(\underline{x}) = $ const. $= C_1$



$x_1$ = living area

Decr. boundary

Also called: **LDA: linear discriminant analysis**.
(if parameters are estimated from the data)

<u>Case B:</u>  $\underline{\underline{\Sigma}}_i = $ arbitrary

$$g_i(\underline{x}) = -\frac{1}{2} \ln\left|\underline{\underline{\Sigma}}_i\right| - \frac{1}{2} d_M^2 (\underline{x}, \underline{m}_i) + \ln P(S_i)$$

$$\underbrace{\phantom{-\frac{1}{2} d_M^2 (\underline{x}, \underline{m}_i)}}$$

$$(\underline{x}-\underline{m}_i)^T \underline{\underline{\Sigma}}_i^{-1} (\underline{x}-\underline{m}_i)$$
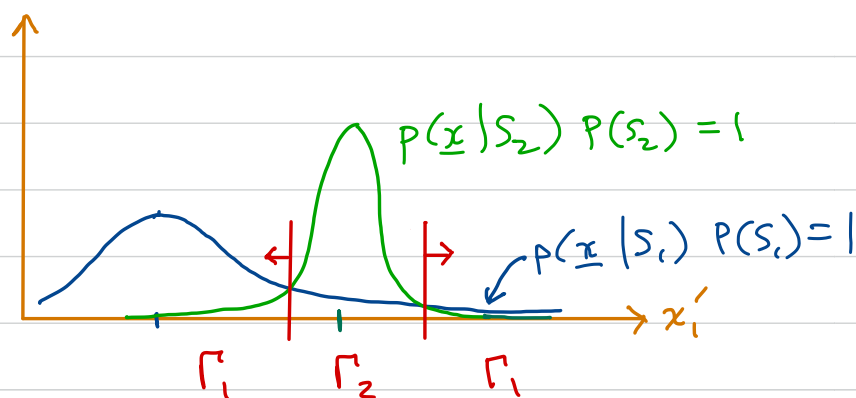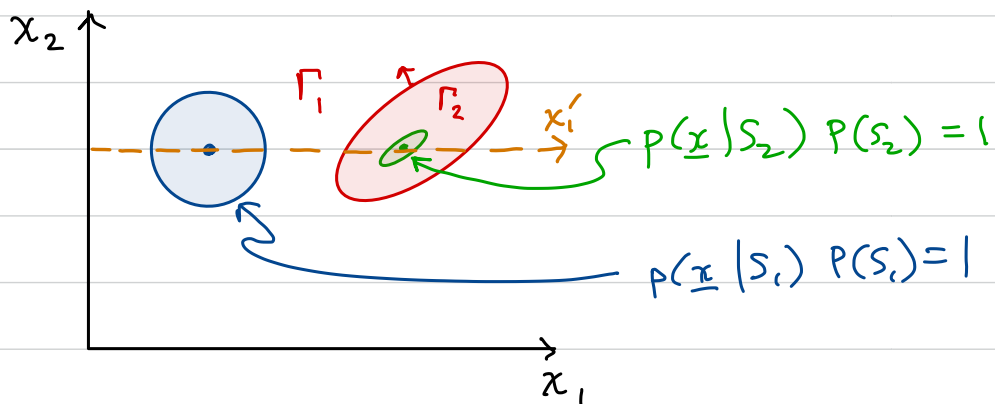
$g_i(\underline{x})$ is quadratic fcn. of $\underline{x}$.

$\Rightarrow$ Dec. boundaries are quadratic

$\rightarrow$ <u>Quadratic Bayes</u>

Also called <u>QDA: quadratic discriminant analysis</u>
(if parameters are estimated from the data)

$\leftarrow g_1(\underline{x}) = g_2(\underline{x})$

$S_1$: price increase

$S_2$: price decrease



$x_2 = $ dist. to c.c.

$d_M = 2$
$d_M = 1$

$\underline{m}_1$

$\underline{m}_2$

$d_M = 1$
$d_M = 2$

$x_1 = $ living area

$\leftarrow$ dec. boundary (quadratic)
ex:  hyperparabola
hyperhyperbola
hyperellipsoid
hypersphere

Ex 1:

$x_2$

$\Gamma_1$  $\Gamma_2$  $x_1'$

$P(\underline{x} | S_2) P(S_2) = 1$

$P(\underline{x} | S_1) P(S_1) = 1$

$x_1$

$P(\underline{x} | S_2) P(S_2) = 1$

$P(\underline{x} | S_1) P(S_1) = 1$

$x_1'$

$\Gamma_1$  $\Gamma_2$  $\Gamma_1$

Ex 2:

$x_2$

$\Gamma_2$  $\Gamma_1$

$\Gamma_1$

$x_1$