# Bayes Decision Theory - II

Nuno Vasconcelos

(Ken Kreutz-Delgado)

UCSD

# Nearest Neighbor Classifier
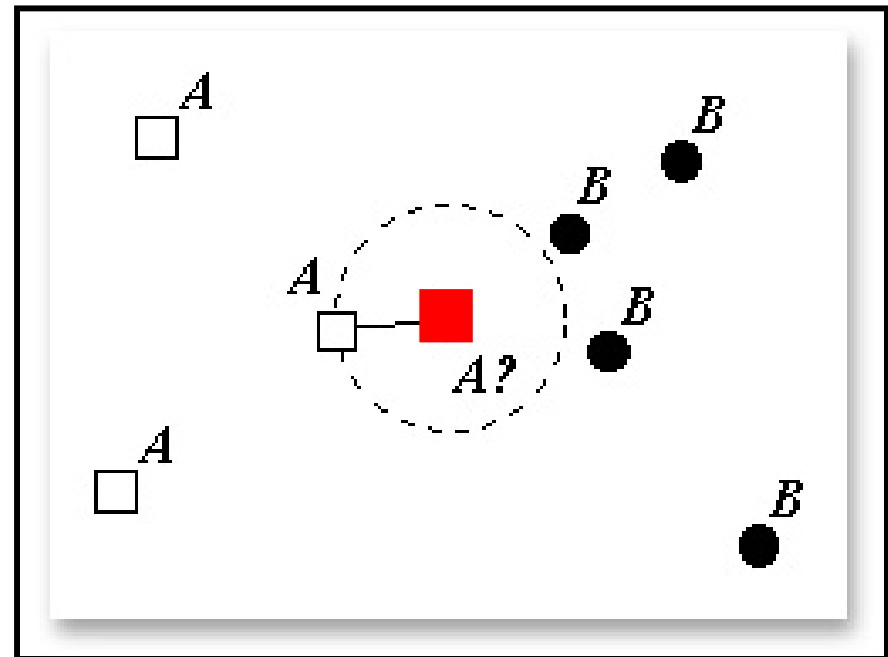
- We are considering *supervised* classification
- Nearest Neighbor (NN) Classifier
  - A *training set* $\mathcal{D} = \{(x_1, y_1), \ldots, (x_n, y_n)\}$
  - $x_i$ is a *vector of observations*, $y_i$ is the *corresponding class label*
  - a vector $x$ to classify
- The *"NN Decision Rule"* is
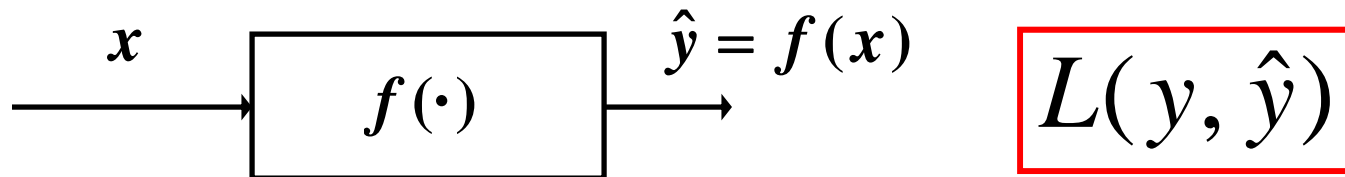
$$\text{Set} \quad y = y_{i*}$$
$$\text{where}$$
$$i* = \arg \min_{i \in \{1, \ldots, n\}} d(x, x_i)$$

  - argmin means: "the *i* that minimizes the distance"

# Optimal Classifiers

- We have seen that performance depends on metric
- Some metrics are "better" than others
- The meaning of "better" is connected to how well adapted the metric is to the properties of the data
- But can we be more rigorous? what do we mean by optimal?
- To talk about optimality we define cost or loss

$$x \longrightarrow \boxed{f(\cdot)} \longrightarrow \hat{y} = f(x) \qquad \boxed{L(y, \hat{y})}$$

  – Loss is the function that we want to minimize
  – Loss depends on true y and prediction $\hat{y}$
  – Loss tells us how good our predictor is

# Loss Functions

- Loss is a function of *classification errors*
  - What errors can we have?
  - Two types: *false positives* and *false negatives*
    - consider a face detection problem (decide "face" or "non-face")
    - if you see this and say

      "face"                    "non-face"

    - you have a
      false – positive          false-negative
      (false alarm)             (miss, failure to detect)
  - Obviously, we have corresponding sub-classes for non-errors
    - true-positives and true-negatives
  - positive/negative part reflects what we say or decide,
  - true/false part reflects the true class label ("true state of the world")

# (Conditional) Risk

- To weigh different errors differently
  - We introduce a loss function
  - Denote the cost of classifying *X* from class *i* as *j* by

$$L[i \rightarrow j]$$

  - One way to measure how good the classifier is to use the (data-conditional) expected value of the loss, aka the (conditional) Risk,

$$R(x,i) = E\{L[Y \rightarrow i] | x\} = \sum_j L[j \rightarrow i] P_{Y|X}(j|x)$$

- this means
  - risk of classifying x as *i* is equal to
  - sum, over all classes, of the loss of classifying as *i* when truth is *j*
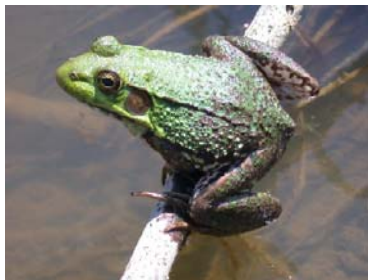  - times probability that true class is *j* (given x)

5

# Loss Functions

- example: two snakes and eating poisonous dart frogs
  - Regular snake will die
  - Frogs are a good snack for the predator dart-snake
  - This leads to the losses

| Regular snake | dart frog | regular frog |
|---|---|---|
| regular | ∞ | 0 |
| dart | 0 | 10 |

| Predator snake | dart frog | regular frog |
|---|---|---|
| regular | 10 | 0 |
| dart | 0 | 10 |

  - What is optimal decision when snakes find a frog like these?

# Minimum Risk Classification

- We have seen that
  - if both snakes have

$$P_{Y|X}(j \mid x) = \begin{cases} 0 & j = \textbf{dart} \\ 1 & j = \textbf{regular} \end{cases}$$

  then both say "regular"

  - However, if

$$P_{Y|X}(j \mid x) = \begin{cases} 0.1 & j = \textbf{dart} \\ 0.9 & j = \textbf{regular} \end{cases}$$

  then the vulnerable snake says "dart"

  while the predator says "regular"

- Its infinite loss for saying regular when frog is dart, makes the vulnerable snake much more cautious!

# Bayes decision rule

- Note that the definition of risk:
  - Immediately defines the *optimal classifier* as the one that *minimizes the conditional risk* for a *given observation x*
  - The *Optimal Decision* is the *Bayes Decision Rule (BDR)* :

$$i^*(x) = \arg\min_i R(x,i)$$

$$= \arg\min_i \sum_j L[j \to i] P_{Y|X}(j \mid x).$$

  - The BDR yields the *optimal (minimal) risk* :

$$R^*(x) = R(x,i^*) = \min_i \sum_j L[j \to i] P_{Y|X}(j \mid x)$$

# The 0/1 Loss Function

- An important special case of interest:
  - zero loss for no error and equal loss for two error types
- This is equivalent to the "zero/one" loss :

$$L[i \rightarrow j] = \begin{cases} 0 & i = j \\ 1 & i \neq j \end{cases}$$

| snake prediction | dart frog | regular frog |
|---|---|---|
| regular | 1 | 0 |
| dart | 0 | 1 |

- Under this loss the optimal Bayes decision rule (BDR) is

$$d^*(x) = i^*(x) = \arg\min_i \sum_j L[j \rightarrow i] P_{Y|X}(j \mid x)$$

$$= \arg\min_i \sum_{j \neq i} P_{Y|X}(j \mid x)$$

9

# 0/1 Loss yields MAP Decision Rule

- Note that :

$$i^*(x) = \arg\min_i \sum_{j \neq i} P_{Y|X}(j \mid x)$$

$$= \arg\min_i \left[ 1 - P_{Y|X}(i \mid x) \right]$$

$$= \arg\max_i P_{Y|X}(i \mid x)$$

- Thus the Optimal Decision for the 0/1 loss is :
  - Pick the class that is most probable given the observation *x*
  - *i\*(x)* is known as the Maximum *a Posteriori* Probability (MAP) solution

- This is also known as the Bayes Decision Rule (BDR) for the 0/1 loss
  - We will often simplify our discussion by assuming this loss
  - But you should always be aware that other losses may be used
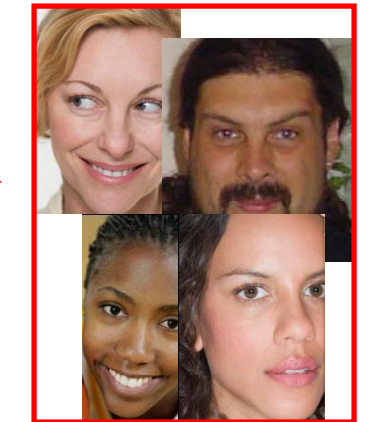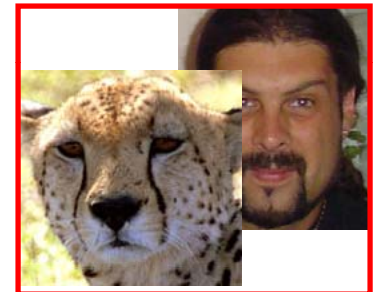
# BDR for the 0/1 Loss

- Consider the evaluation of the BDR for 0/1 loss

$$i^*(x) = \arg\max_i P_{Y|X}(i \mid x)$$

- This is also called the Maximum a Posteriori Probability (MAP) rule

- It is usually *not* trivial to evaluate the posterior probabilities $P_{Y|X}(i \mid x)$

- This is due to the fact that we are trying to infer the cause (class *i*) from the consequence (observation *x*) – i.e. we are trying to solve a nontrivial inverse problem

  - E.g. imagine that I want to evaluate

    $P_{Y|X}($ *person* | *"has two eyes"*$)$

  - This strongly depends on *what the other classes are*

# Posterior Probabilities and Detection

- If the two classes are "people" and "cars"
  - then $P_{Y|X}($ *person* | *"has two eyes"* $) = 1$



- But if the classes are "people" and "cats"
  - then $P_{Y|X}($ *person* | *"has two eyes"* $) = \frac{1}{2}$
    *if* there are equal numbers of cats and people
    to uniformly choose from [ this is additional info! ]



- How do we deal with this problem?
  - We note that it is much easier to infer consequence from cause
  - E.g., it is easy to infer that
    $$P_{X|Y}( \text{"has two eyes"} \mid person ) = 1$$
  - This does *not* depend on any other classes
  - We do *not* need any additional information
  - Given a class, *just count* the frequency of observation



12

# Bayes Rule

- How do we go from $P_{X|Y}(x \mid j)$ to $P_{Y|X}(j \mid x)$ ?
- We use *Bayes rule*:

$$P_{Y|X}(i \mid x) = \frac{P_{X|Y}(x \mid i) P_Y(i)}{P_X(x)}$$

- Consider the two-class problem, i.e. *Y=0* or *Y=1*
  - the BDR under 0/1 loss is

$$i^*(x) = \arg\max_i P_{Y|X}(i \mid x)$$

$$= \begin{cases} 0, & \text{if } P_{Y|X}(0 \mid x) \geq P_{Y|X}(1 \mid x) \\ 1, & \text{if } P_{Y|X}(0 \mid x) < P_{Y|X}(1 \mid x) \end{cases}$$

# BDR for 0/1 Loss Binary Classification

- Pick "0" when $P_{Y|X}(0\,|\,x) \geq P_{Y|X}(1\,|\,x)$ and "1" otherwise
- Using Bayes rule on both sides of this inequality yields

$$P_{Y|X}(0\,|\,x) \geq P_{Y|X}(1\,|\,x) \Leftrightarrow$$

$$\frac{P_{X|Y}(x\,|\,0)P_Y(0)}{P_X(x)} \geq \frac{P_{X|Y}(x\,|\,1)P_Y(1)}{P_X(x)}$$

 - Noting that $P_X(x)$ is a non-negative quantity this is the same as the rule pick "0" when

$$P_{X|Y}(x\,|\,0)P_Y(0) \geq P_{X|Y}(x\,|\,1)P_Y(1)$$

i.e.

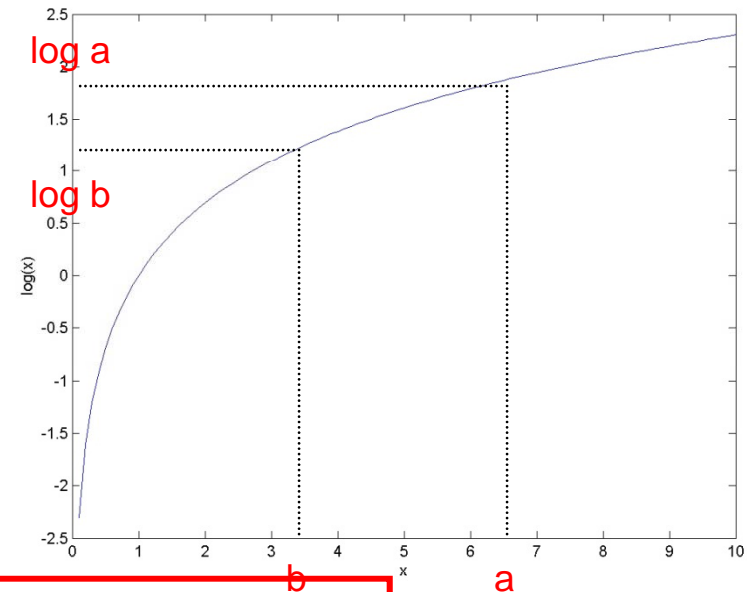$$i^*(x) = \arg\max_{i} P_{X|Y}(x\,|\,i)\,P_Y(i)$$

# The "Log Trick"

- Sometimes it's not convenient to work directly with pdf's
  - One helpful trick is to take logs
  - Note that the log is a monotonically increasing function

$$a > b \Leftrightarrow \log a > \log b$$

from which we have

$$i^*(x) = \arg\max_i P_{X|Y}(x \mid i) P_Y(i)$$

$$= \arg\max_i \log\left(P_{X|Y}(x \mid i) P_Y(i)\right)$$

$$= \arg\max_i \left(\log P_{X|Y}(x \mid i) + \log P_Y(i)\right)$$

$$= \arg\min_i \left(-\log P_{X|Y}(x \mid i) - \log P_Y(i)\right)$$

15

# "Standard" (0/1) BDR

- In summary
  - for the zero/one loss, the following three decision rules are optimal and equivalent

1) $$i^*(x) = \arg\max_i P_{Y|X}(i \mid x)$$

2) $$i^*(x) = \arg\max_i \left[ P_{X|Y}(x \mid i) P_Y(i) \right]$$

3) $$i^*(x) = \arg\max_i \left[ \log P_{X|Y}(x \mid i) + \log P_Y(i) \right]$$

The form 1) is usually hardest to use, 3) is frequently easier than 2)

# BDR - Example

- So far the BDR is an abstract rule
  - How does one implement the optimal decision in practice?
  - In addition to having a loss function, you need to *know, model, or estimate the probabilities*!
  - Example
    - Suppose that you run a gas station
    - On Mondays you have a promotion to sell more gas
    - Q: is the promotion working?    I.e., is $Y = 0$ (no) or $Y = 1$ (yes) ?
    - A good observation to answer this question is the interarrival time ($\tau$) between cars

high $\tau$: not working ($Y = 0$)          low $\tau$: working well ($Y = 1$)

# BDR - Example

- What are the class-conditional and prior probabilities?
  - the probability of arrival of a car follows a Poisson distribution
  - Poisson inter-arrival times are exponentially distributed
    - Hence

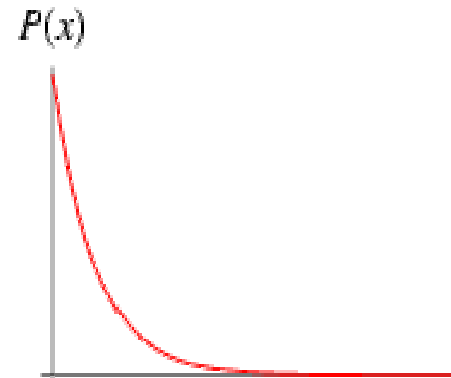$$P_{X|Y}(\tau \mid i) = \lambda_i e^{-\lambda_i \tau}$$

    where $\lambda_i$ is the arrival rate (cars/s).
    - The expected value of the interarrival time is

$$E_{X|Y}[x \mid y = i] = \frac{1}{\lambda_i}$$

    - Consecutive times are assumed to be independent :

$$P_{X_1,\ldots,X_n|Y}(\tau_1,\ldots,\tau_n \mid i) = \prod_{k=1}^{n} P_{X|Y}(\tau_k \mid i) = \prod_{k=1}^{n} \lambda_i e^{-\lambda_i \tau_k}$$

$P(x)$

# BDR - Example

- Let's assume that we
  - know $\lambda_i$ and the (prior) class probabilities $P_Y(i) = \pi_i$ , $i = 0,1$
  - Have measured a collection of times during the day, $\mathcal{D} = \{\tau_1,...,\tau_n\}$

- The probabilities are of exponential form
  - Therefore it is easier to use the log-based BDR

$$i^*(\mathrm{D}) = \arg\max_i \left[ \log P_{X|Y}(\mathrm{D} \mid i) + \log P_Y(i) \right]$$

$$= \arg\max_i \left[ \log\left( \prod_{k=1}^{n} \lambda_i e^{-\lambda_i \tau_k} \right) + \log \pi_i \right]$$

$$= \arg\max_i \left[ -\sum_{k=1}^{n} \lambda_i \tau_k + n \log \lambda_i + \log \pi_i \right]$$

$$= \arg\max_i \left[ -\sum_{k=1}^{n} \lambda_i \tau_k + n \log\left( \lambda_i \sqrt[n]{\pi_i} \right) \right]$$

# BDR - Example

- This means we pick "0" when

$$-\sum_{k=1}^{n} \lambda_0 \, \tau_k + n\log\left(\lambda_0 \sqrt[n]{\pi_0}\right) \geq -\sum_{k=1}^{n} \lambda_1 \, \tau_k + n\log\left(\lambda_1 \sqrt[n]{\pi_1}\right), \quad \text{or}$$

$$(\lambda_1 - \lambda_0)\sum_{k=1}^{n} \tau_k \geq n\log\left(\frac{\lambda_1 \sqrt[n]{\pi_1}}{\lambda_0 \sqrt[n]{\pi_0}}\right), \quad \text{or}$$

$$\frac{1}{n}\sum_{k=1}^{n} \tau_k \geq \frac{1}{(\lambda_1 - \lambda_0)}\log\left(\frac{\lambda_1 \sqrt[n]{\pi_1}}{\lambda_0 \sqrt[n]{\pi_0}}\right)$$

(reasonably taking $\lambda_1 > \lambda_0$)

and "1" otherwise

- Does this decision rule make sense?
  - Let's assume, for simplicity, that $\pi_1 = \pi_2 = 1/2$

# BDR - Example

- For $\pi_1 = \pi_2 = \frac{1}{2}$, we pick "promotion did not work" *(Y=0)* if

$$\frac{1}{n}\sum_{k=1}^{n}\tau_k \geq \frac{1}{(\lambda_1 - \lambda_0)}\log\left(\frac{\lambda_1}{\lambda_0}\right)$$

The left hand side is the (sample) average interarrival time for the day

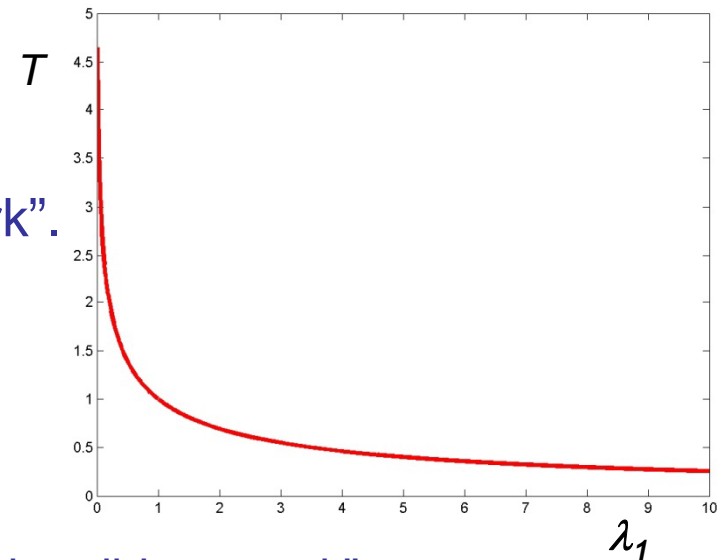  – This means that there is an optimal choice of a *"threshold"*

$$T = \frac{1}{(\lambda_1 - \lambda_0)}\log\left(\frac{\lambda_1}{\lambda_0}\right)$$



above which we say "promotion did not work".
This  makes sense!

  – What is the shape of this threshold?
    ▪ Assuming $\lambda_0 = 1$, it looks like this.
    ▪ Higher the $\lambda_1$, the more likely to say "promotion did not work".

21

# BDR - Example

- When $\pi_1 = \pi_2 = \frac{1}{2}$, we pick "did not work" *(Y=0)* when

$$\boxed{\frac{1}{n}\sum_{k=1}^{n}\tau_k \geq T} \qquad \boxed{T = \frac{1}{(\lambda_1 - \lambda_0)}\log\left(\frac{\lambda_1}{\lambda_0}\right)}$$



- – Assuming $\lambda_0 = 1$, *T* decreases with $\lambda_1$
- – I.e. for a given daily average,
  - Larger $\lambda_1$: easier to say "did not work"
- – This means that
  - As the expected rate of arrival for good days increases we are going to impose a tougher standard on the average measured interarrival times
  - The average has to be smaller for us to accept the day as a good one
- – Once again, *this makes sense!*
- – usually the case with the BDR (a good way to check your math)

22

# The Gaussian Classifier

- One important case is that of Multivariate Gaussian Classes
  - The pdf of class *i* is a Gaussian of mean $\mu_i$ and covariance $\Sigma_i$

$$P_{X|Y}(x \mid i) = \frac{1}{\sqrt{(2\pi)^d |\Sigma_i|}} \exp\left\{ -\frac{1}{2}(x - \mu_i)^T \Sigma_i^{-1}(x - \mu_i) \right\}$$

- The BDR is

$$i^*(x) = \arg\max_i \left[ -\frac{1}{2}(x - \mu_i)^T \Sigma_i^{-1}(x - \mu_i) \right.$$
$$\left. -\frac{1}{2}\log(2\pi)^d |\Sigma_i| + \log P_Y(i) \right]$$

# Implementation

- To design a Gaussian classifier (e.g. homework)
  - Start from a collection of datasets, where the $i$-th class dataset $\mathcal{D}^{(i)} = \{x_1^{(i)}, ..., x_n^{(i)}\}$ is a set of $n^{(i)}$ examples from class $i$
  - For each class *estimate* the Gaussian parameters :

$$\hat{\mu}_i = \frac{1}{n^{(i)}} \sum_j x_j^{(i)} \qquad \hat{\Sigma}_i = \frac{1}{n^{(i)}} \sum_j (x_j^{(i)} - \hat{\mu}_i)(x_j^{(i)} - \hat{\mu}_i)^T \qquad \hat{P}_Y(i) = \frac{n^{(i)}}{T}$$

where  T  is the total number of examples over all c classes

- the BDR is approximated as

$$i^*(x) = \arg\max_i \left[ -\frac{1}{2}(x - \hat{\mu}_i)^T \hat{\Sigma}_i^{-1}(x - \hat{\mu}_i) \right.$$
$$\left. -\frac{1}{2}\log(2\pi)^d \left| \hat{\Sigma}_i \right| + \log \hat{P}_Y(i) \right]$$

# Gaussian Classifier

- The Gaussian Classifier can be written as
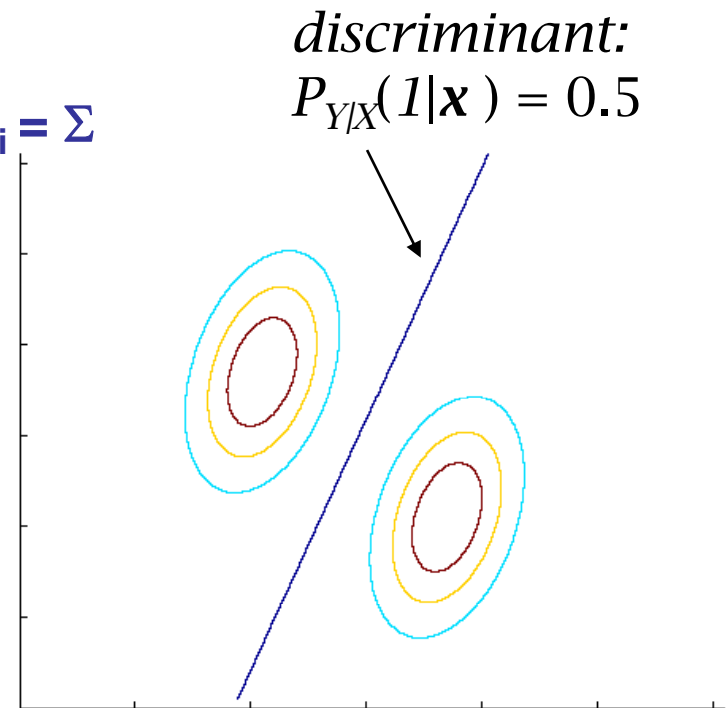
*discriminant:*
$$P_{Y|X}(1|\boldsymbol{x}) = 0.5$$

$$i^*(x) = \arg\min_i \left[ d^2{}_i(x, \mu_i) + \alpha_i \right]$$

with

$$d^2{}_i(x, y) = (x - y)^T \Sigma_i^{-1} (x - y)$$

$$\alpha_i = \log(2\pi)^d |\Sigma_i| - 2\log P_Y(i)$$

and can be seen as a nearest "class-neighbor" classifier with a "funny metric"

– Each class has its own "distance" measure:
  - Sum the Mahalanobis-squared for that class, then add the $\alpha$ constant.
  - We effectively have different "metrics" in the data (feature) space that are class i dependent.

25

# Gaussian Classifier

- A special case of interest is when
  - All classes have the same covariance $\Sigma_i = \Sigma$

$$i^*(x) = \arg\min_i \left[ d^2(x, \mu_i) + \alpha_i \right]$$

with

$$d^2(x, y) = (x - y)^T \Sigma^{-1}(x - y)$$

$$\alpha_i = -2\log P_Y(i)$$

*discriminant:*
$$P_{Y|X}(1|\boldsymbol{x}) = 0.5$$

- Note that:

  - $\alpha_i$ can be dropped when all classes have *equal prior probability*
  - This is reminiscent of the NN classifier with Mahalanobis distance
  - Instead of finding the *nearest data point neighbor* of *x*, it looks for the *nearest class "prototype*," (or "archetype," or "exemplar," or "template," or "representative", or "ideal", or "form") , defined as the class mean $\mu_i$

26

# Binary Classifier – Special Case

- Consider $\Sigma_i = \Sigma$ with two classes

  - One important property of this case is that the *decision boundary* is a hyperplane (Homework)

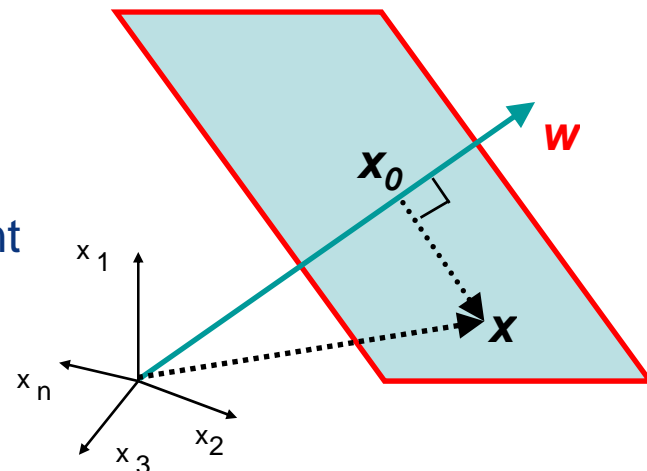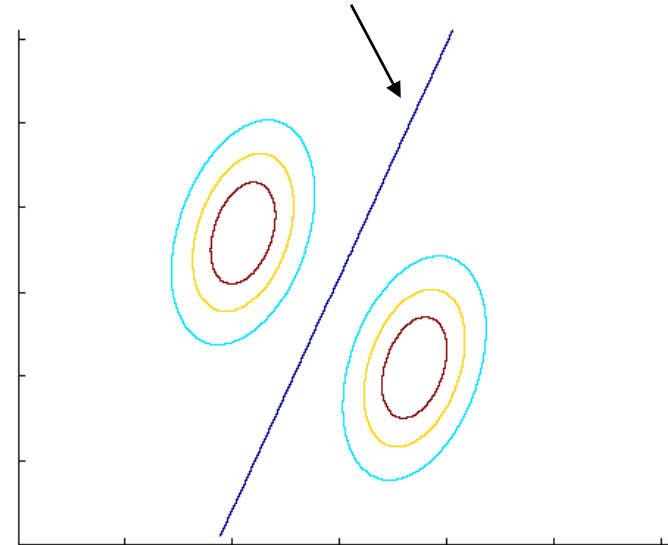  - This can be shown by computing the set of points *x* such that

$$d^2(x, \mu_0) + \alpha_0 = d^2(x, \mu_1) + \alpha_1$$

and showing that they satisfy

$$w^T(x - x_0) = 0$$

  - This is the equation of a *hyperplane* with normal *w*. $x_0$ can be *any* fixed point on the hyperplane, but it is *standard* to choose it to have minimum norm, in which case *w* and $x_0$ are then parallel

*Discriminant Surface:*
$$P_{Y|X}(1|x) = 0.5$$

27

# Gaussian Classifier

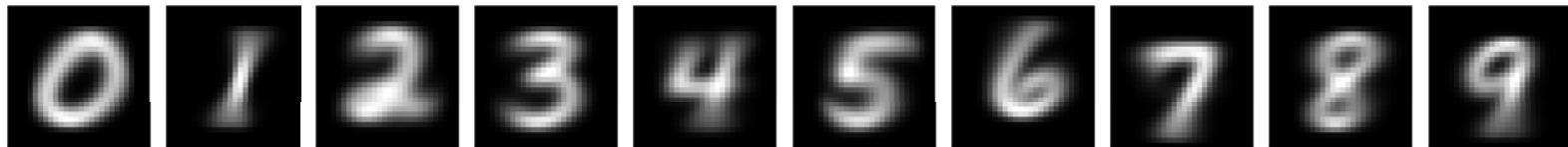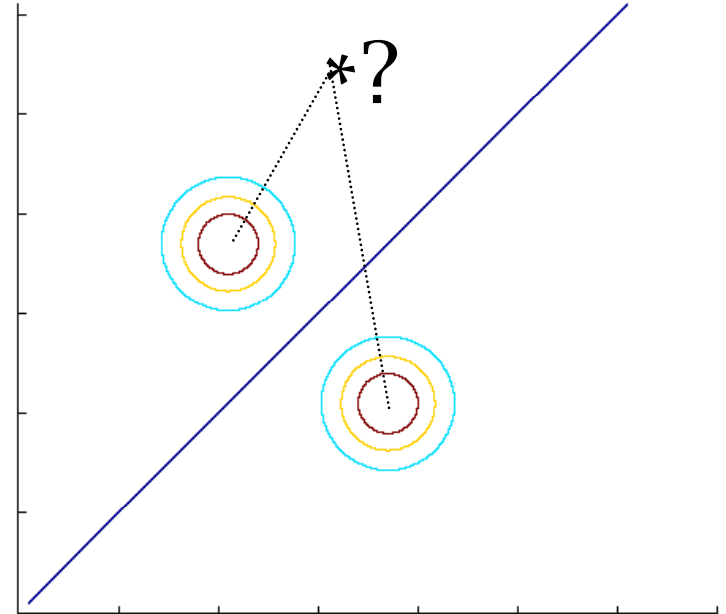- If *all* the class covariances are the identity, $\Sigma_i=I$, then

$$i^*(x) = \arg\min_i \left[ d^2(x, \mu_i) + \alpha_i \right]$$

with

$$d^2(x, y) = \| x - y \|^2$$

$$\alpha_i = -2\log P_Y(i)$$

- This is called template matching with class means as templates
  - E.g. for digit classification



Compare the complexity of this classifier to NN Classifier!

# The Sigmoid Function

- We have derived much of the above from the log-based BDR

$$i^*(x) = \arg\max_i \left[ \log P_{X|Y}(x|i) + \log P_Y(i) \right]$$

- When there are only two classes, $i = 0,1$, it is also interesting to consider the original definition

$$i^*(x) = \arg\max_i g_i(x)$$

where

$$g_i(x) = P_{Y|X}(i|x) = \frac{P_{X|Y}(x|i)P_Y(i)}{P_X(x)}$$

$$= \frac{P_{X|Y}(x|i)P_Y(i)}{P_{X|Y}(x|0)P_Y(0) + P_{X|Y}(x|1)P_Y(1)}$$

# The Sigmoid Function

- Note that this can be written as

$$i^*(x) = \arg\max_i g_i(x)$$

$$g_1(x) = 1 - g_0(x)$$

$$g_0(x) = \cfrac{1}{1 + \cfrac{P_{X|Y}(x|1)\, P_Y(1)}{P_{X|Y}(x|0)\, P_Y(0)}}$$

- For Gaussian classes, the posterior probabilities are

$$g_0(x) = \frac{1}{1 + \exp\left\{ d^2_{\,0}(x, \mu_0) - d^2_{\,1}(x, \mu_1) + \alpha_0 - \alpha_1 \right\}}$$

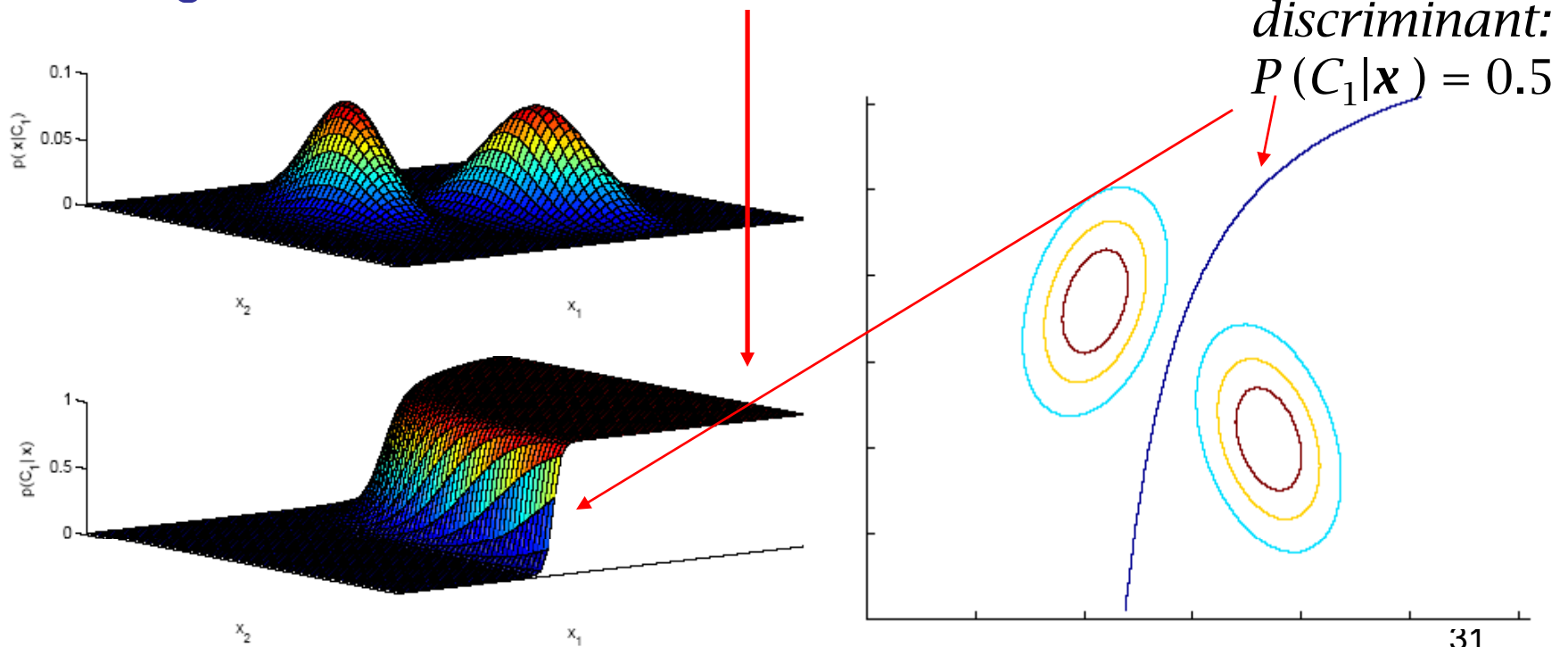where, as before,   $d^2_{\,i}(x, y) = (x - y)^T \Sigma_i^{-1} (x - y)$

$$\alpha_i = \log(2\pi)^d \left| \Sigma_i \right| - 2\log P_Y(i)$$

30

# The Sigmoid ("S-shaped") Function

- The posterior pdf for class *i = 0*,

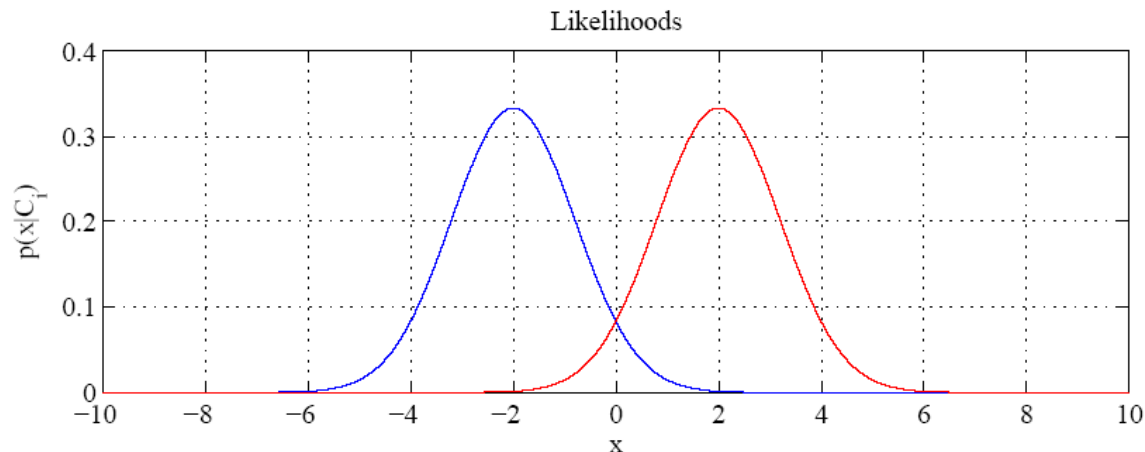$$g_0(x) = \frac{1}{1 + \exp\left\{d^2_{\,0}(x, \mu_0) - d^2_{\,1}(x, \mu_1) + \alpha_0 - \alpha_1\right\}}$$

is a sigmoid and looks like this

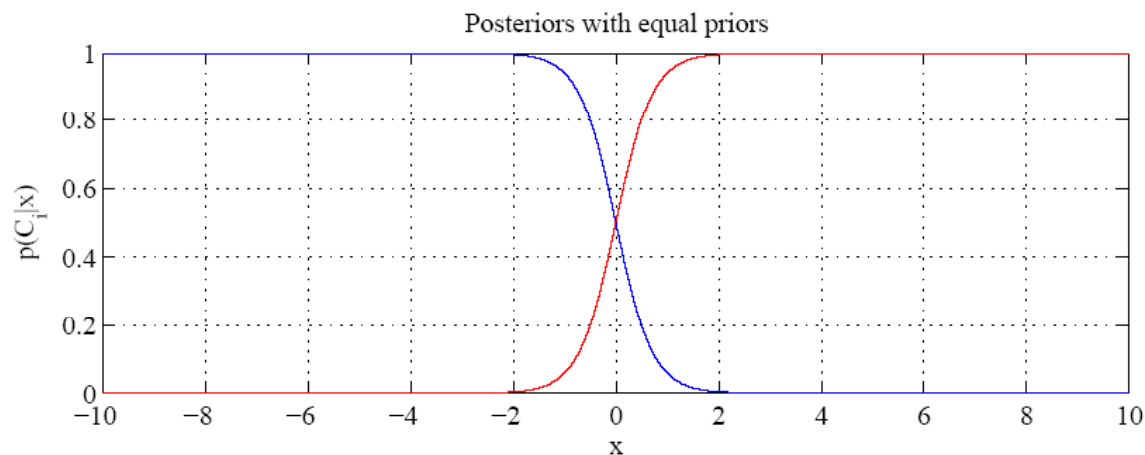*discriminant:*
$P(C_1|\boldsymbol{x}) = 0.5$

# The Sigmoid

- The sigmoid appears in neural networks, where it can be interpreted as a posterior pdf for a Gaussian binary classification problem when the covariances are the same
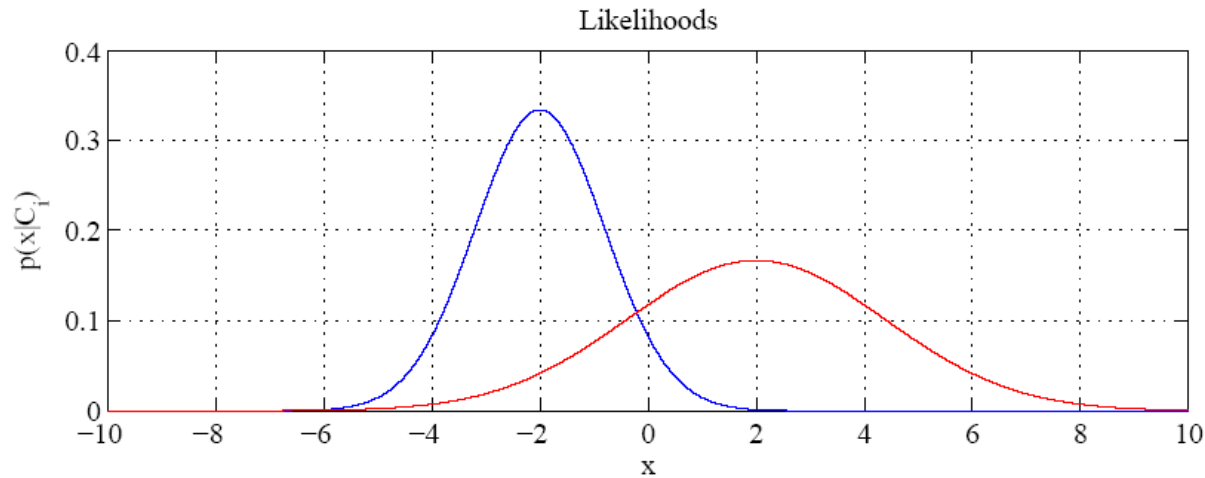


*Equal variances*
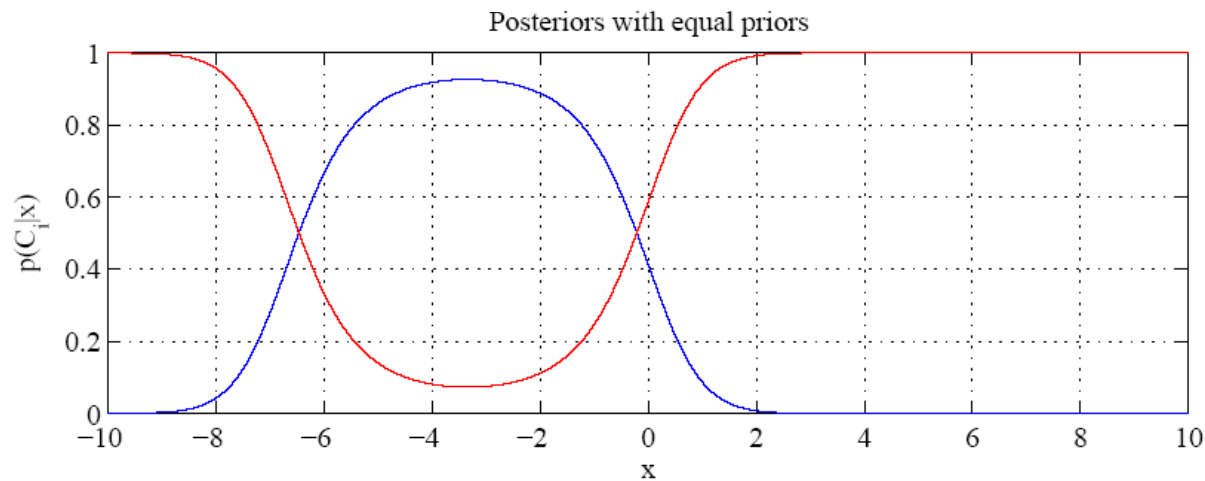
*Single boundary at halfway between means*

# The Sigmoid

- But not necessarily when the covariances are *different*



*Variances are different*

*Yields two boundaries*

Any questions?