

## CS 2750 Machine Learning Lecture 9

### Multi-way classification

Milos Hauskrecht  
[milos@cs.pitt.edu](mailto:milos@cs.pitt.edu)  
5329 Sennott Square

---

CS 2750 Machine Learning

### Administrative announcements

- **Homework 3 due today**
- **Homework 4 out**

---

CS 2750 Machine Learning

## Multi-way classification

- **Binary classification**  $Y = \{0,1\}$
- **Multi-way classification**
  - **K classes**  $Y = \{0,1,\dots, K-1\}$
  - **Goal:** learn to classify correctly K classes
  - Or **learn**  $f : X \rightarrow \{0,1,\dots, K-1\}$

- **Errors:**

- **Zero-one (misclassification) error for an example:**

$$Error_1(\mathbf{x}_i, y_i) = \begin{cases} 1 & f(\mathbf{x}_i, \mathbf{w}) \neq y_i \\ 0 & f(\mathbf{x}_i, \mathbf{w}) = y_i \end{cases}$$

- **Mean misclassification error (for a dataset):**

$$\frac{1}{n} \sum_{i=1}^n Error_1(\mathbf{x}_i, y_i)$$

---

CS 2750 Machine Learning

## Multi-way classification

### Approaches:

- **Discriminative approach**
  - Parametric discriminant functions
  - Learns discriminant functions **directly**
    - Much like the logistic regression model.
- **Generative model approach**
  - Generative model of the distribution  $p(\mathbf{x},y)$
  - Learns the parameters of the model through density estimation techniques
  - Discriminant functions defined on the top of the model
    - “Indirect” learning of a classifier

---

CS 2750 Machine Learning

## Generative model approach

**Indirect:**

1. Represent and learn the distribution  $p(\mathbf{x}, y)$
2. Define and use probabilistic discriminant functions

$$g_i(\mathbf{x}) = \log p(y = i | \mathbf{x})$$

**Model**  $p(\mathbf{x}, y) = p(\mathbf{x} | y)p(y)$

- $p(\mathbf{x} | y) =$  **Class-conditional distributions (densities)**

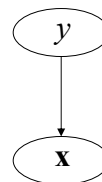
k class conditional distributions

$$p(\mathbf{x} | y = i) \quad \forall i \quad 0 \leq i \leq K - 1$$

- $p(y) =$  **Priors on classes**

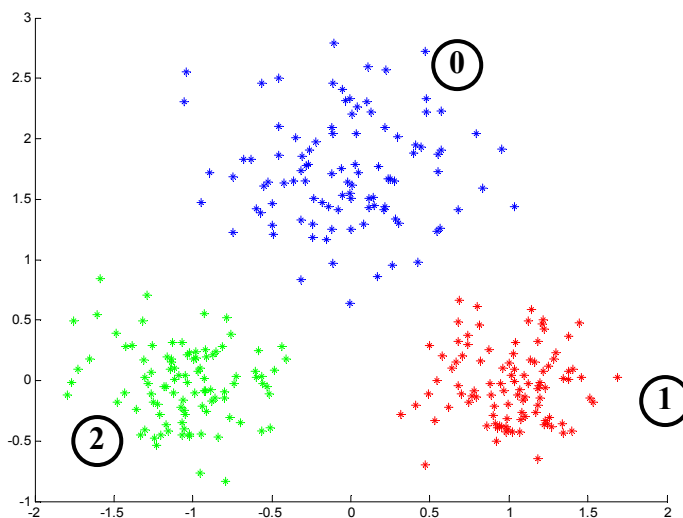
- - probability of class  $y$

$$\sum_{i=1}^{K-1} p(y = i) = 1$$



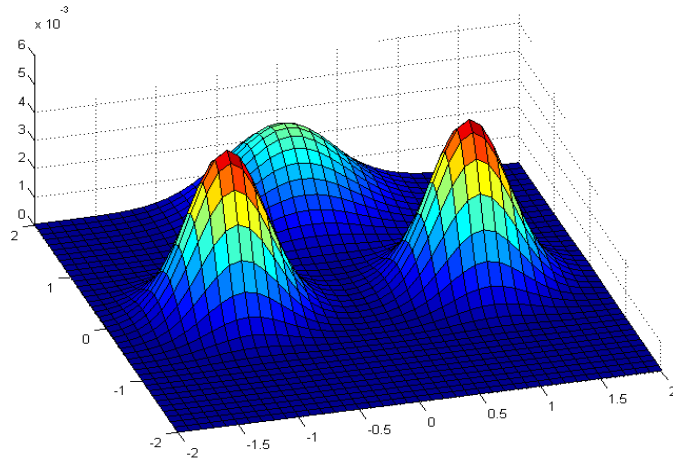
CS 2750 Machine Learning

## Multi-way classification. Example



CS 2750 Machine Learning

## Multi-way classification



CS 2750 Machine Learning

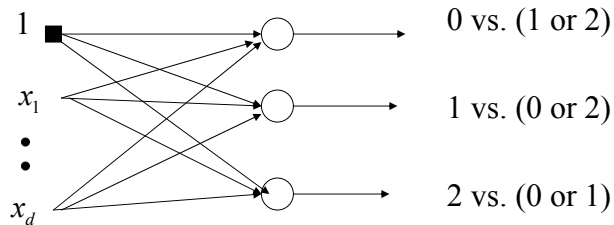
## Discriminative approach.

- **Parameteric model** of discriminant functions
- Learns the discriminant functions directly

How to learn to classify multiple classes, say 0,1,2?

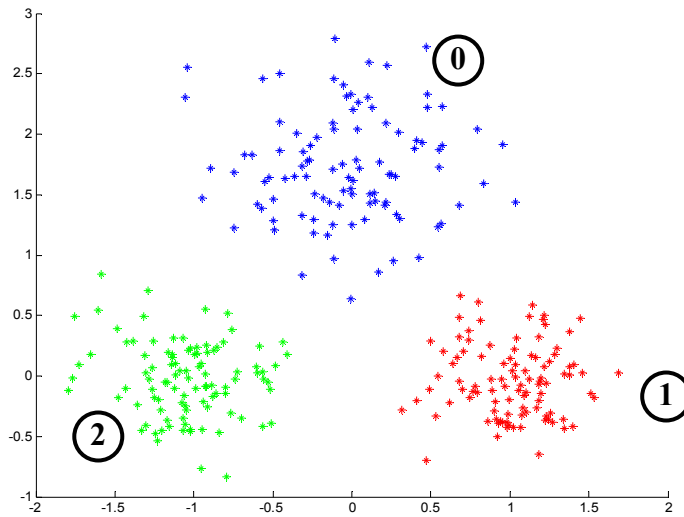
### Approach 1:

- Logistic regression on every class versus the rest



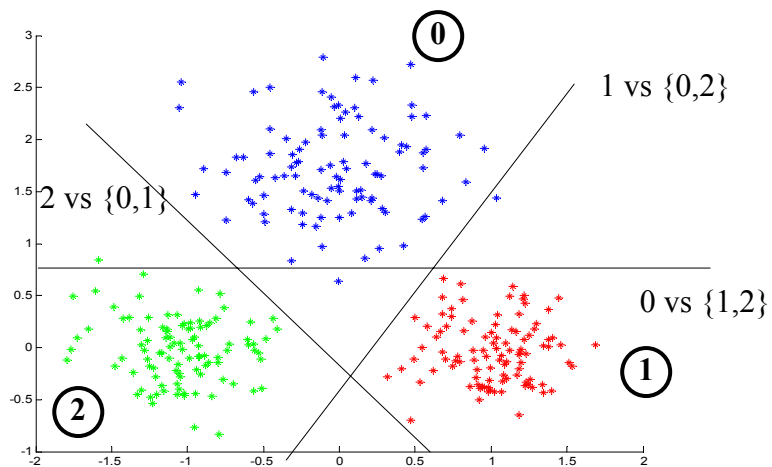
CS 2750 Machine Learning

## Multi-way classification. Example



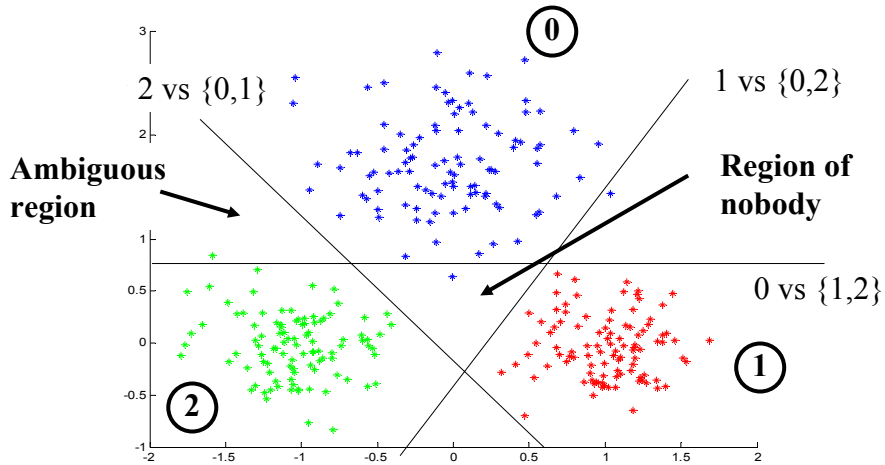
CS 2750 Machine Learning

## Multi-way classification. Approach 1.



CS 2750 Machine Learning

## Multi-way classification. Approach 1.



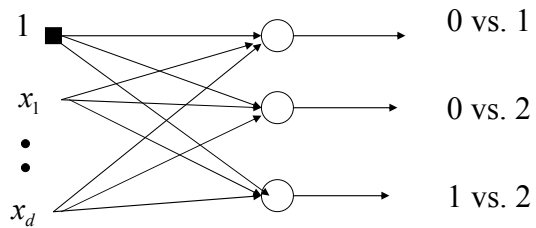
CS 2750 Machine Learning

## Discriminative approach.

How to learn to classify multiple classes, say 0,1,2 ?

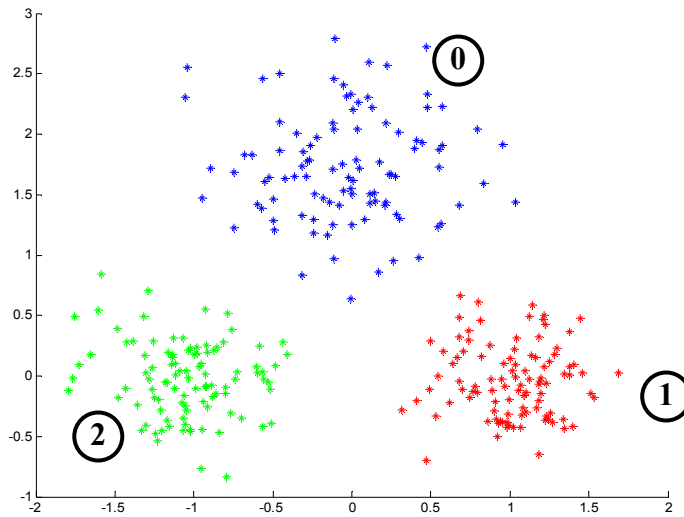
### Approach 2:

- Logistic regression on all pairs



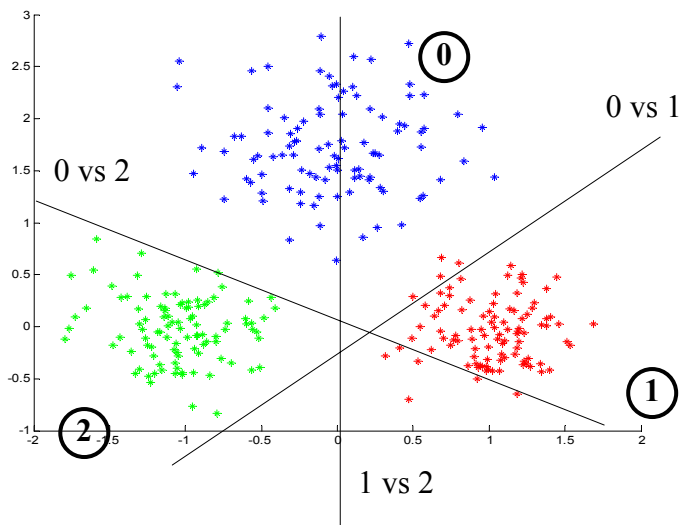
CS 2750 Machine Learning

## Multi-way classification. Example



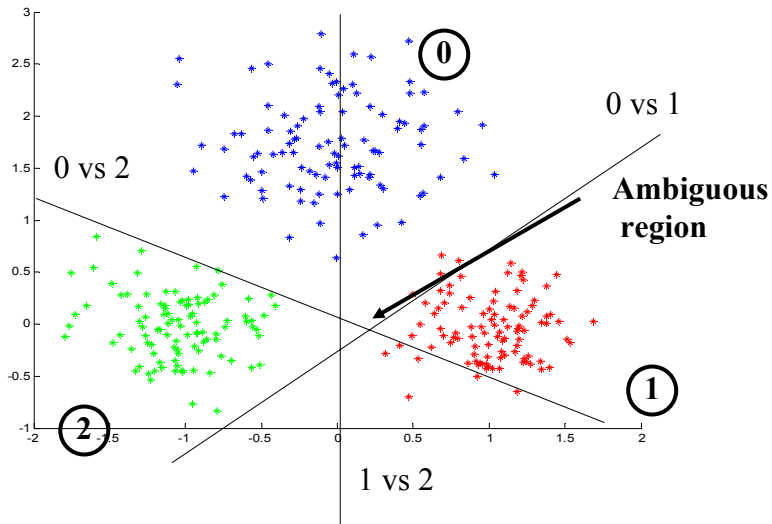
CS 2750 Machine Learning

## Multi-way classification. Approach 2



CS 2750 Machine Learning

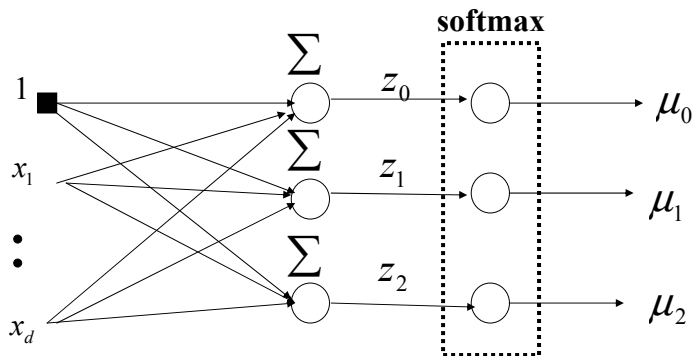
## Multi-way classification. Approach 2



CS 2750 Machine Learning

## Multi-way classification with softmax

- A solution to the problem of having an ambiguous region

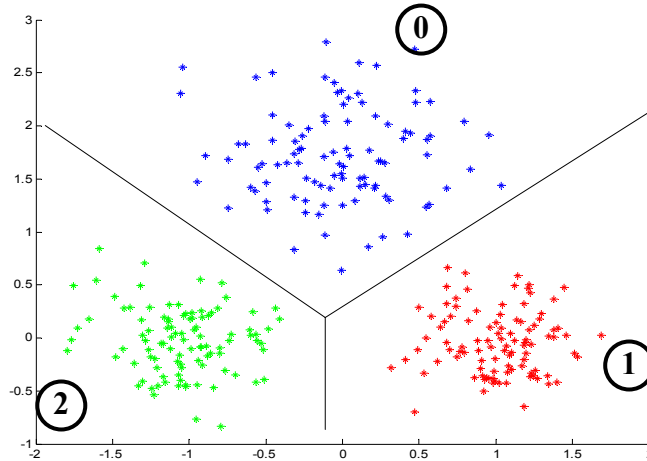


$$p(y = i | \mathbf{x}) = \mu_i = \frac{\exp(\mathbf{w}_i^T \mathbf{x})}{\sum_j \exp(\mathbf{w}_j^T \mathbf{x})} \quad \sum_i \mu_i = 1$$

CS 2750 Machine Learning



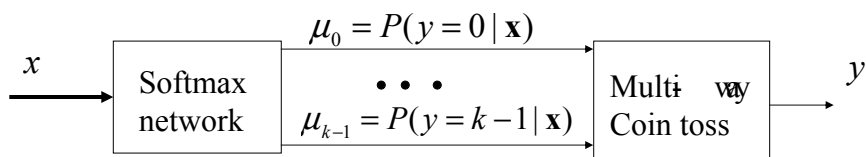
## Multi-way classification with softmax



CS 2750 Machine Learning

## Learning of the softmax model

- Learning of parameters  $\mathbf{w}$ : statistical view



Assume outputs  $y$  are transformed as follows

$$y \in \{0 \ 1 \ \dots \ k-1\} \quad \longrightarrow \quad y \in \left\{ \begin{array}{c} \begin{pmatrix} 1 \\ 0 \\ \dots \\ 0 \end{pmatrix} \quad \begin{pmatrix} 0 \\ 1 \\ \dots \\ 0 \end{pmatrix} \quad \dots \quad \begin{pmatrix} 0 \\ 0 \\ \dots \\ 1 \end{pmatrix} \end{array} \right\}$$

CS 2750 Machine Learning

## Learning of the softmax model

- Learning of the parameters  $\mathbf{w}$ : statistical view

- **Likelihood of outputs**

$$L(D, \mathbf{w}) = p(\mathbf{Y} | \mathbf{X}, \mathbf{w}) = \prod_{i=1, \dots, n} p(y_i | \mathbf{x}_i, \mathbf{w})$$

- We want parameters  $\mathbf{w}$  that maximize the likelihood

- **Log-likelihood trick**

– Optimize  $\log$  likelihood of outputs instead:

$$\begin{aligned} l(D, \mathbf{w}) &= \log \prod_{i=1, \dots, n} p(y_i | \mathbf{x}_i, \mathbf{w}) = \sum_{i=1, \dots, n} \log p(y_i | \mathbf{x}_i, \mathbf{w}) \\ &= \sum_{i=1, \dots, n} \sum_{q=0}^{k-1} \log \mu_i^{y_{i,q}} = \sum_{i=1, \dots, n} \sum_{q=0}^{k-1} y_{i,q} \log \mu_{i,q} \end{aligned}$$

- **Objective to optimize**

$$J(D_i, \mathbf{w}) = - \sum_{i=1}^n \sum_{q=0}^{k-1} y_{i,q} \log \mu_{i,q}$$

---

CS 2750 Machine Learning

## Learning of the softmax model

- **Error to optimize:**

$$J(D_i, \mathbf{w}) = - \sum_{i=1}^n \sum_{q=0}^{k-1} y_{i,q} \log \mu_{i,q}$$

- **Gradient**

$$\frac{\partial}{\partial w_{jk}} J(D_i, \mathbf{w}) = \sum_{i=1}^n -x_{i,j} (y_{i,j} - \mu_{i,j})$$

- The same very easy **gradient update** as used for the logistic regression

$$\mathbf{w}_j \leftarrow \mathbf{w}_j + \alpha \sum_{i=1}^n (y_{i,j} - \mu_{i,j}) \mathbf{x}_i$$

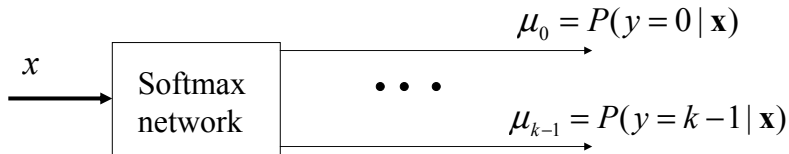
- But now we have to update weights of k networks

---

CS 2750 Machine Learning

## Multi-way classification

- When is the **softmax** the right model ?



- Assume:

$$p(\mathbf{x} | y = i) = h(\mathbf{x}, \boldsymbol{\varphi}) \exp \left\{ \frac{(\boldsymbol{\theta}_i^T \mathbf{x} - A(\boldsymbol{\theta}_i))}{a(\boldsymbol{\varphi})} \right\}$$

$\boldsymbol{\theta}_i$  - location parameter for class conditional  $i$

$\boldsymbol{\varphi}$  - scaling parameter (the same for all classes)

## Multi-way classification

- Class conditional:**

$$p(\mathbf{x} | y = i) = h(\mathbf{x}, \boldsymbol{\varphi}) \exp \left\{ \frac{(\boldsymbol{\theta}_i^T \mathbf{x} - A(\boldsymbol{\theta}_i))}{a(\boldsymbol{\varphi})} \right\}$$

- Class posterior:**

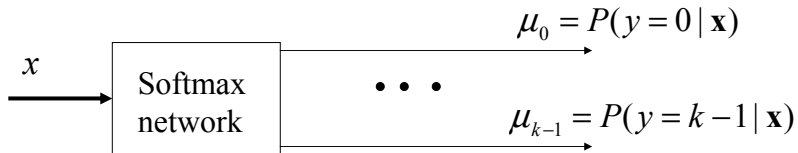
$$p(y = i | \mathbf{x}) = \frac{p(\mathbf{x} | y = i) p(y = i)}{\sum_j p(\mathbf{x} | y = j) p(y = j)}$$

$$= \frac{h(\mathbf{x}, \boldsymbol{\varphi}) \exp \left\{ \frac{(\boldsymbol{\theta}_i^T \mathbf{x} - A(\boldsymbol{\theta}_i))}{a(\boldsymbol{\varphi})} \right\} p(y = i)}{\sum_j h(\mathbf{x}, \boldsymbol{\varphi}) \exp \left\{ \frac{(\boldsymbol{\theta}_j^T \mathbf{x} - A(\boldsymbol{\theta}_j))}{a(\boldsymbol{\varphi})} \right\} p(y = j)} = \frac{\exp(\mathbf{w}_i^T \mathbf{x} + b_i)}{\sum_j \exp(\mathbf{w}_j^T \mathbf{x} + b_j)}$$

$$\mathbf{w}_i = \frac{\boldsymbol{\theta}_i}{a(\boldsymbol{\varphi})} \quad b_i = \frac{A(\boldsymbol{\theta}_i)}{a(\boldsymbol{\varphi})} + \ln p(y = i)$$

## Multi-way classification

- **Softmax model is an accurate model** when class conditional densities are represented with densities from the exponential family and the same scaling parameter



$$p(\mathbf{x} | y = i) = \exp \left\{ \frac{(\boldsymbol{\theta}_i^T \mathbf{x} - b(\boldsymbol{\theta}_i))}{a(\boldsymbol{\varphi})} + c(\mathbf{x}, \boldsymbol{\varphi}) \right\}$$

$\boldsymbol{\theta}_i$  - location parameter for class conditional  $i$

$\boldsymbol{\varphi}$  - scaling parameter (the same for all classes)

## CS 2750 Machine Learning Lecture 9b

### Bayesian decision theory.

Milos Hauskrecht  
[milos@cs.pitt.edu](mailto:milos@cs.pitt.edu)  
5329 Sennott Square

## Confusion matrix

- **Confusion matrix:**

- Counts of examples with:
- class label  $\omega_j$  that are classified with a label  $\alpha_i$

	$\alpha = 0$	$\alpha = 1$	$\alpha = 2$
$\omega = 0$	140	20	22
$\omega = 1$	17	54	8
$\omega = 2$	12	4	76

agreement

## Zero-one loss function

- **Misclassification error**

- Based on the zero-one loss function
  - Any misclassified example counts as 1
  - Correctly classified example counts as 0

	$\alpha = 0$	$\alpha = 1$	$\alpha = 2$
$\omega = 0$	140	20	22
$\omega = 1$	17	54	8
$\omega = 2$	12	4	76

Loss of 1 (pointing to the top-right triangle of the matrix)

Loss of 1 (pointing to the bottom-left triangle of the matrix)

Loss 0 (pointing to the diagonal elements of the matrix)

## General loss function

- **Error function based on a more general loss function**
  - Different misclassifications have different weight (loss)
  - $\alpha_i$  our choice
  - $\omega_j$  true label
  - $\lambda(\alpha_i | \omega_j)$  loss for classification

**Example:**

		$\alpha = 0$	$\alpha = 1$	$\alpha = 2$
$\lambda(\alpha_i   \omega_j)$	$\omega = 0$	0	1	5
	$\omega = 1$	3	0	2
	$\omega = 2$	3	1	0

CS 2750 Machine Learning

## Bayesian decision theory

- **More general loss function**
  - Different misclassifications have different weight (loss)

$$\lambda(\alpha_i | \omega_j)$$

- **Expected loss for choice** (action)  $\alpha_i$

$$R(\alpha_i | \mathbf{x}) = \sum_j \lambda(\alpha_i | \omega_j) P(\omega_j | \mathbf{x})$$

- Also called conditional risk

- **Decision rule:**  $\alpha(\mathbf{x})$ 
  - Chooses label (action) according to the input
- **Overall expected loss for the decision rule**

$$R(\alpha) = \int R(\alpha(\mathbf{x}), \mathbf{x}) P(\mathbf{x}) d\mathbf{x}$$

CS 2750 Machine Learning

## Bayesian decision theory

- The optimal decision rule

$$\alpha^*(\mathbf{x}) = \arg \min_{\alpha_i} \sum_j \lambda(\alpha_i | \omega_j) P(\omega_j | \mathbf{x})$$

How to modify classifiers to handle different loss functions?

- **Discriminative models:**

- Directly optimize the parameters according to the new loss function – estimate the loss based on examples

- **Generative models:**

- Learn probabilities as before
- Decisions about classes are biased to minimize the empirical loss (as seen above)

$$\alpha^*(\mathbf{x}) = \arg \min_{\alpha_i} \sum_j \lambda(\alpha_i | \omega_j) P(\omega_j | \mathbf{x})$$

---

CS 2750 Machine Learning

## Estimating the loss

- Counts in the **confusion matrix** can be used to estimate the **expected loss function**:

$N(\alpha_i   \omega_j)$	$\alpha = 0$	$\alpha = 1$	$\alpha = 2$
$\omega = 0$	140	20	22
$\omega = 1$	17	54	8
$\omega = 2$	12	4	76

agreement

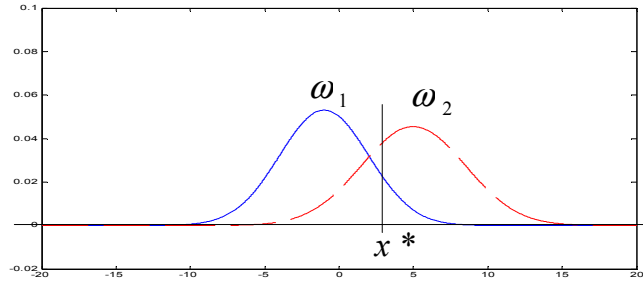
- **Loss**

$$\frac{1}{N} \sum_i \sum_j \lambda(\alpha_i | \omega_j) N(\alpha_i | \omega_j)$$

---

CS 2750 Machine Learning

## Binary decisions: accuracy.



- **Probabilities:**

- True positive (hit)  $p(x > x^* \mid \mathbf{x} \in \omega_2)$
- False positive (false alarm)  $p(x > x^* \mid \mathbf{x} \in \omega_1)$
- True negative (correct rejection)  $p(x < x^* \mid \mathbf{x} \in \omega_1)$
- False negative (a miss)  $p(x < x^* \mid \mathbf{x} \in \omega_2)$

CS 2750 Machine Learning

## Binary decision. Accuracy

- **Confusion matrix**

	$x \in \omega_2(+)$	$x \in \omega_1(-)$	
$x > x^*(+)$	140	10	$PPT=140/150$
$x < x^*(-)$	20	180	$NPV=180/200$
	$SENS=140/160$ $SPEC=180/190$		

- **Row and column quantities:**

- Sensitivity (SENS)
- Specificity (SPEC)
- Positive predictive value (PPV)
- Negative predictive value (NPV)

CS 2750 Machine Learning



## Classification accuracy

- Sensitivity**

$$\frac{p(x > x^* | \mathbf{x} \in \omega_2)}{p(x > x^* | \mathbf{x} \in \omega_2) + p(x < x^* | \mathbf{x} \in \omega_2)} = \frac{TP}{TP + FN}$$

- Specificity**

$$\frac{p(x < x^* | \mathbf{x} \in \omega_1)}{p(x < x^* | \mathbf{x} \in \omega_1) + p(x > x^* | \mathbf{x} \in \omega_1)} = \frac{TN}{TN + FP}$$

- Positive predictive value**

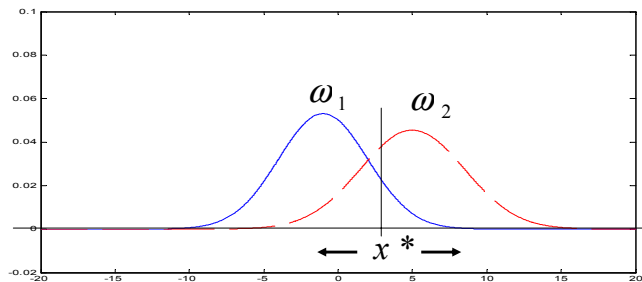
$$\frac{p(x > x^* | \mathbf{x} \in \omega_2)}{p(x > x^* | \mathbf{x} \in \omega_2) + p(x > x^* | \mathbf{x} \in \omega_1)} = \frac{TP}{TP + FP}$$

- Negative predictive value**

$$\frac{p(x < x^* | \mathbf{x} \in \omega_1)}{p(x < x^* | \mathbf{x} \in \omega_1) + p(x < x^* | \mathbf{x} \in \omega_2)} = \frac{TN}{TN + FN}$$

CS 2750 Machine Learning

## Decision threshold



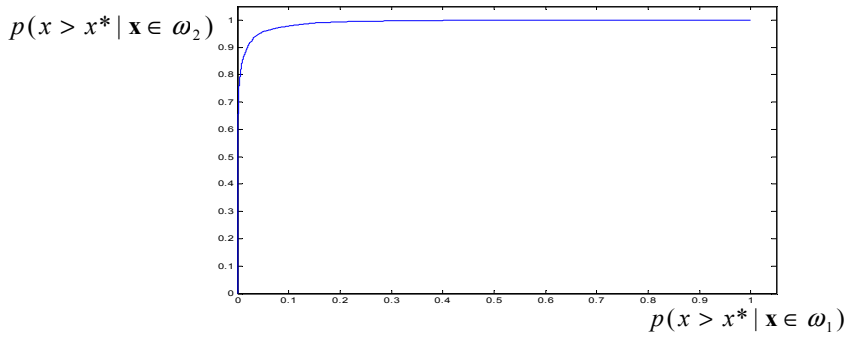
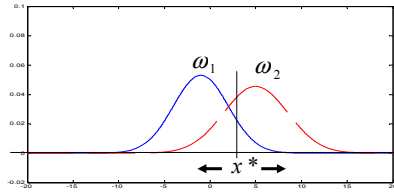
- Movement of  $x^*$  changes the probabilities:**

- True positive (hit)  $p(x > x^* | \mathbf{x} \in \omega_2)$
- False positive (false alarm)  $p(x > x^* | \mathbf{x} \in \omega_1)$
- True negative (correct rejection)  $p(x < x^* | \mathbf{x} \in \omega_1)$
- False negative (a miss)  $p(x < x^* | \mathbf{x} \in \omega_2)$

CS 2750 Machine Learning

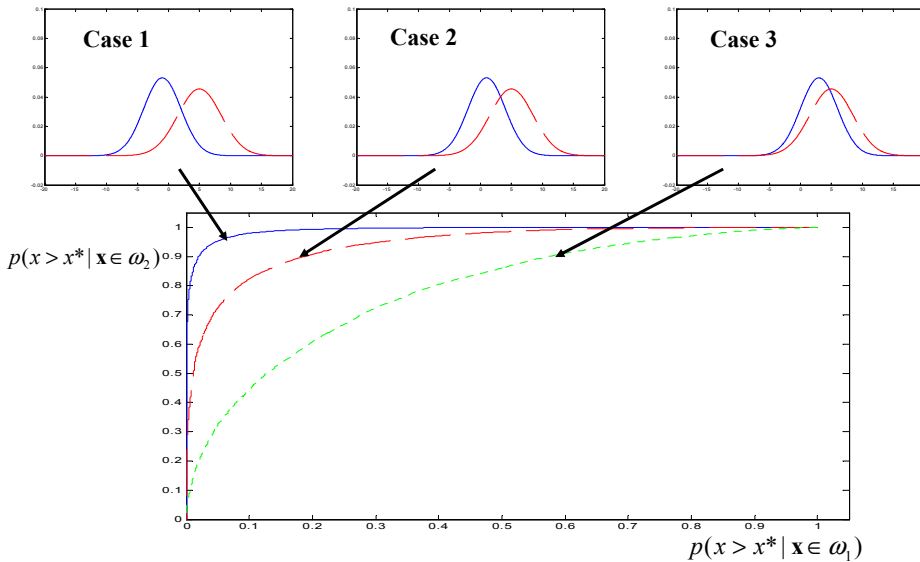
# Receiver Operating Characteristic (ROC)

- **ROC curve plots :**  
 $p(x > x^* | \mathbf{x} \in \omega_1)$   
 vs  $p(x > x^* | \mathbf{x} \in \omega_2)$   
**for different  $x^*$**



CS 2750 Machine Learning

## ROC curve



CS 2750 Machine Learning

## Receiver operating characteristic

- **ROC**
  - shows the discriminability between the two classes under different decision biases
- **Decision bias**
  - can be changed using different loss function