





## 1 Introduction

## 2 Logistic Regression

### ③ Evaluation Metrics

#### 4 Extra reading

## 5 References



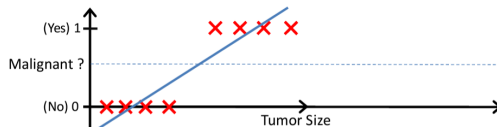
# Linear Regression for Classification

- A natural approach is to use linear regression:

$$h_{\theta}(x) = \theta_0 + \theta_1 x$$

and define a threshold at 0.5 for prediction:

$$\hat{y} = \begin{cases} 1, & h_{\theta}(x) \geq 0.5 \\ 0, & h_{\theta}(x) < 0.5 \end{cases}$$













## Introduction (cont.)

- $\sigma(w^T x)$ : probability that  $y = 1$  given  $x$  (parameterized by  $\mathbf{w}$ )

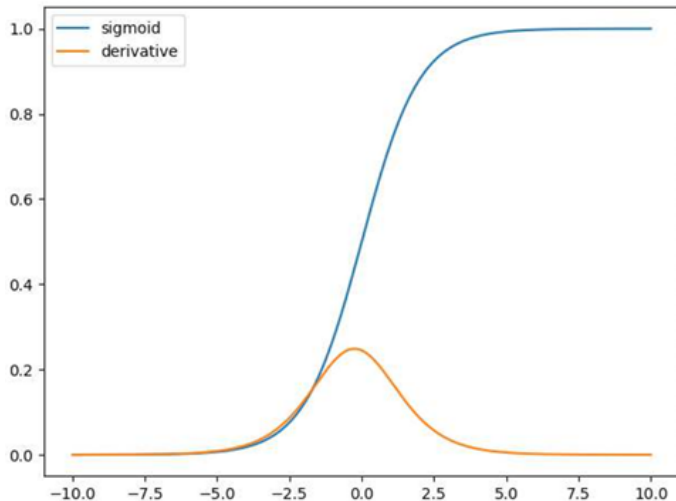
$$P(y = 1|x, \mathbf{w}) = \sigma(\mathbf{w}^T x)$$

$$P(y = 0|x, \mathbf{w}) = 1 - \sigma(\mathbf{w}^T x)$$

- We need to look for a function which gives us an output in the range  $[0, 1]$ . (like a probability).
- Let's denote this function with  $\sigma(\cdot)$  and call it the **activation function**.



# Sigmoid function & its derivative



## Introduction (cont.)

- The sigmoid function takes a number as input but we have:

$$x = [x_0 = 1, x_1, \dots, x_d]$$

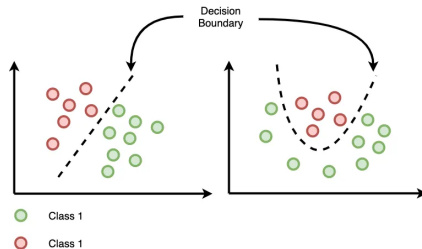
$$w = [w_0, w_1, \dots, w_d]$$

- So we can use the **dot product** of  $x$  and  $w$ .
- We have  $0 \leq \sigma(\mathbf{w}^T x) \leq 1$ . which is the estimated probability of  $y = 1$  on input  $x$ .
- An Example : A basketball game (Win, Lose)
  - $\sigma(\mathbf{w}^T x) = 0.7$
  - In other terms 70 percent chance of winning the game.



# Decision surface

- Decision surface or decision boundary is the region of a problem space in which the output label of a classifier is ambiguous. (could be linear or non-linear)
- In binary classification it is where the probability of a sample belonging to each  $y = 0$  and  $y = 1$  is equal.



- Decision boundary hyperplane always has **one less dimension** than the feature space.

## Decision surface (cont.)

- An example of linear decision boundaries:

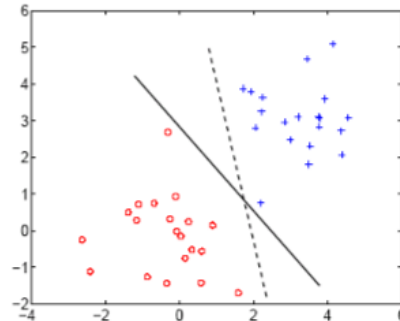
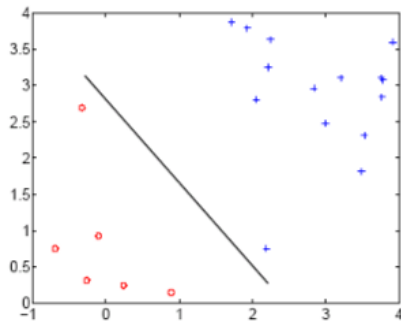


Figure adapted from Eric Xing, Machine Learning, CMU



## Decision surface (cont.)

- Back to our logistic regression problem.
- Decision surface  $\sigma(\mathbf{w}^T x) = \mathbf{constant}$ .

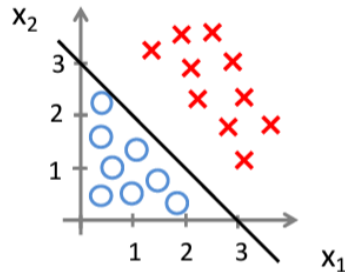
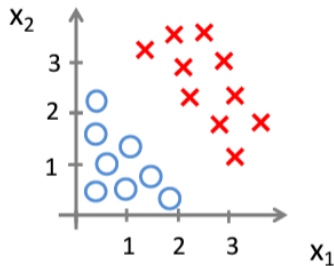
$$\sigma(\mathbf{w}^T x) = \frac{1}{1 + e^{-(\mathbf{w}^T x)}} = 0.5$$

- Decision surfaces are **linear functions** of  $x$ 
  - if  $\sigma(\mathbf{w}^T x) \geq 0.5$  then  $\hat{y} = 1$ , else  $\hat{y} = 0$
  - Equivalently, if  $\mathbf{w}^T x + w_0 \geq 0.5$  then decide  $\hat{y} = 1$ , else  $\hat{y} = 0$

**$\hat{y}$  is the predicted label**

## Decision boundary example

$$\sigma(\mathbf{w}^T x) = \sigma(w_0 + w_1 x_1 + w_2 x_2)$$

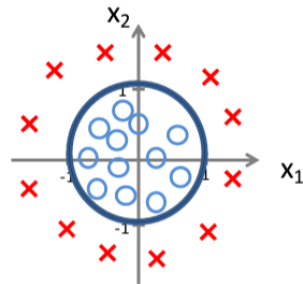
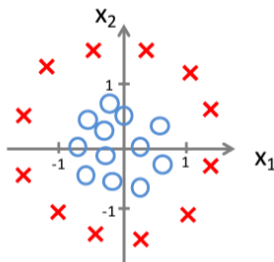


Predict  $y = 1$  if  $-3 + x_1 + x_2 \geq 0$

# Non-linear decision boundary example

$$\sigma(\mathbf{w}^T x) = \sigma(w_0 + w_1 x_1 + w_2 x_2 + w_3 x_1^2 + w_4 x_2^2)$$

We can learn more complex decision boundaries when having higher order terms



$$\text{Predict } y = 1 \text{ if } -1 + x_1^2 + x_2^2 \geq 0$$

## 1 Introduction

## ② Logistic Regression

## Fundamentals

## Decision surface

MLE and MAP

## Gradient descent

## Multi-class logistic regression

### ③ Evaluation Metrics

#### 4 Extra reading

## 5 References

# Maximum Likelihood Estimation (MLE)

- For  $n$  independent samples, the likelihood is:

$$L(\mathbf{w}) = \prod_{i=1}^n P(y^{(i)} | x^{(i)}, \mathbf{w})$$

- For binary classification ( $y \in \{0, 1\}$ ):

$$P(y^{(i)} | x^{(i)}, \mathbf{w}) = \sigma(\mathbf{w}^T x^{(i)})^{y^{(i)}} (1 - \sigma(\mathbf{w}^T x^{(i)}))^{1-y^{(i)}}$$

- Log-likelihood:

$$\ell(\mathbf{w}) = \sum_{i=1}^n [y^{(i)} \log \sigma(\mathbf{w}^T x^{(i)}) + (1 - y^{(i)}) \log(1 - \sigma(\mathbf{w}^T x^{(i)}))]$$

## From Likelihood to Cost Function

- Maximizing the likelihood  $\Leftrightarrow$  minimizing the negative log-likelihood (NLL):

$$J_{\text{MLE}}(\mathbf{w}) = -\ell(\mathbf{w})$$

- Can be written as an integral over the data distribution:

$$J_{\text{MLE}}(\mathbf{w}) = - \int p(x, y) \log P(y|x, \mathbf{w}) \, dx \, dy$$

- Empirical estimate (training data):

$$J_{\text{MLE}}(\mathbf{w}) = -\frac{1}{n} \sum_{i=1}^n \log P(y^{(i)}|x^{(i)}, \mathbf{w})$$

- This is exactly the **cross-entropy loss** used in classification.

# Maximum A Posteriori Estimation (MAP)

- MAP incorporates prior knowledge about parameters:

$$\hat{\mathbf{w}}_{\text{MAP}} = \arg \max_{\mathbf{w}} P(\mathbf{w}|D)$$

- Using Bayes' rule:

$$P(\mathbf{w}|D) \propto P(D|\mathbf{w})P(\mathbf{w})$$

- Equivalently:

$$\hat{\mathbf{w}}_{\text{MAP}} = \arg \max_{\mathbf{w}} \left[ \log P(D|\mathbf{w}) + \log P(\mathbf{w}) \right]$$

- Cost function:

$$J_{\text{MAP}}(\mathbf{w}) = - \sum_{i=1}^n \log P(y^{(i)}|x^{(i)}, \mathbf{w}) - \log P(\mathbf{w})$$

# MAP with Gaussian Prior (L2 Regularization)

- Gaussian prior:  $\mathbf{w} \sim \mathcal{N}(0, \tau^2 I)$

$$P(\mathbf{w}) \propto \exp\left(-\frac{\|\mathbf{w}\|^2}{2\tau^2}\right)$$

- MAP cost function:

$$J_{\text{MAP}}(\mathbf{w}) = J_{\text{MLE}}(\mathbf{w}) + \frac{1}{2\tau^2} \|\mathbf{w}\|^2$$

- Let  $\lambda = 1/\tau^2$ :

$$J_{\text{MAP}}(\mathbf{w}) = J_{\text{MLE}}(\mathbf{w}) + \frac{\lambda}{2} \|\mathbf{w}\|^2$$

- Equivalent to **L2-regularized logistic regression**.



# MAP with Laplace Prior (L1 Regularization)

- Laplace prior:  $\mathbf{w} \sim \text{Laplace}(0, b)$

$$P(\mathbf{w}) \propto \exp\left(-\frac{\|\mathbf{w}\|_1}{b}\right)$$

- MAP cost:

$$J_{\text{MAP}}(\mathbf{w}) = J_{\text{MLE}}(\mathbf{w}) + \lambda \|\mathbf{w}\|_1$$

where  $\lambda = 1/b$

- L1 penalty encourages **sparsity** (feature selection).

## Regularization Effects (L1 vs L2)

- **L2 (Ridge):** smooths weights, keeps all features but shrinks them.
- **L1 (Lasso):** drives some weights to zero, performs feature selection.
- Both prevent overfitting by penalizing model complexity.

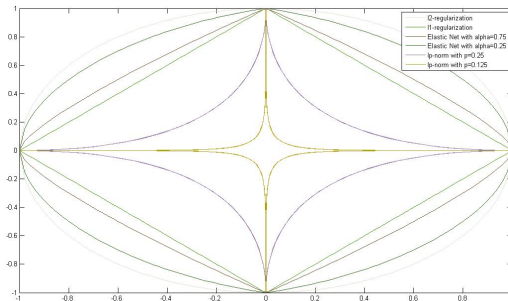


Image adapted from CS 4780, Cornell University

# Regularizers: Types and Properties

Type	Form	Properties / Advantages	Disadvantages / Effect
L2	$r(\mathbf{w}) = \ \mathbf{w}\ _2^2$	Strictly convex, differentiable	Dense solutions (uses all features)
L1	$r(\mathbf{w}) = \ \mathbf{w}\ _1$	Convex, encourages sparsity	Not differentiable at 0
Lp	$r(\mathbf{w}) = \ \mathbf{w}\ _p = (\sum_i  v_i ^p)^{1/p}, 0 < p \leq 1$	Very sparse solutions, initialization dependent	Non-convex, not differentiable

*Note:* Choice of regularizer affects sparsity and smoothness of learned weights.

# MLE vs MAP Comparison

	<b>MLE</b>	<b>MAP</b>
Objective	Maximize likelihood $P(D \mathbf{w})$	Maximize posterior $P(D \mathbf{w})P(\mathbf{w})$
Prior	None (uniform)	Explicit prior $P(\mathbf{w})$
Cost function	$-\log P(D \mathbf{w})$	$-\log P(D \mathbf{w}) - \log P(\mathbf{w})$
Regularization	None	Arises from prior
Overfitting control	No	Yes



# Gradient descent

- Remember from previous slides:

$$J(w) = \sum_{i=1}^n -y^{(i)} \log(\sigma(\mathbf{w}^T x^{(i)})) - (1 - y^{(i)}) \log(1 - \sigma(\mathbf{w}^T x^{(i)}))$$

- Update rule for **gradient descent**:

$$w^{t+1} = w^t - \eta \nabla_w J(w^t)$$

- With  $J(w)$  definition for logistic regression we get:

$$\nabla_w J(w) = \sum_{i=1}^n (\sigma(\mathbf{w}^T x^{(i)}) - y^{(i)}) x^{(i)}$$

# Gradient descent

- Compare the gradient of **logistic regression** with the gradient of **SSE** in **linear regression** :

$$\nabla_w J(w) = \sum_{i=1}^n (\sigma(\mathbf{w}^T x^{(i)}) - y^{(i)}) x^{(i)}$$

$$\nabla_w J(w) = \sum_{i=1}^n (\mathbf{w}^T x^{(i)} - y^{(i)}) x^{(i)}$$

## Loss function

- Loss function is a single overall measure of loss incurred for taking our decisions (over entire dataset).
- We have:

$$Loss(y, \sigma(\mathbf{w}^T x)) = -y \times \log(\sigma(\mathbf{w}^T x)) - (1 - y) \times \log(1 - \sigma(\mathbf{w}^T x))$$

- Since in binary classification either  $y = 1$  or  $y = 0$  we have:

$$Loss(y, \sigma(\mathbf{w}^T x)) = \begin{cases} -\log(\sigma(\mathbf{w}^T x)) & \text{if } y = 1 \\ -\log(1 - \sigma(\mathbf{w}^T x)) & \text{if } y = 0 \end{cases}$$

- How is it related to zero-one loss? ( $\hat{y}$  is the predicted label and  $y$  is the true label)

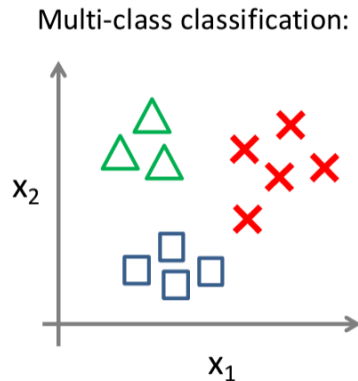
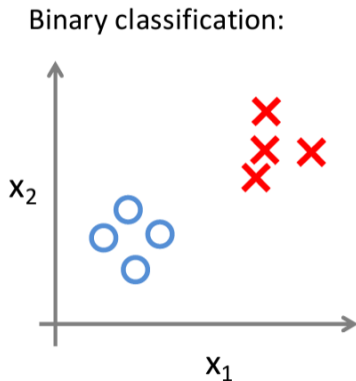
$$Loss(y, \hat{y}) = \begin{cases} 1 & \text{if } y \neq \hat{y} \\ 0 & \text{if } y = \hat{y} \end{cases}$$





# Multi-class logistic regression

- Now consider a problem where we have  $K$  classes and every sample only belongs to one class (for simplicity).



## Multi-class logistic regression (cont.)

- For each class  $k$ ,  $\sigma_k(x; \mathbf{W})$  predicts the probability of  $y = k$ .
  - i.e.,  $P(y = k|x, \mathbf{W})$
- For each data point  $x_0$ ,  $\sum_{k=1}^K P(y = k|x_0, \mathbf{W})$  must be 1
  - $W$  denotes a matrix of  $w_i$ 's in which each  $w_i$  is a weight vector dedicated for class label  $i$ .
- On a new input  $x$ , to make a prediction, we pick the class that maximizes  $\sigma_k(x; \mathbf{W})$ :

$$\alpha(x) = \operatorname{argmax}_{k=1,\dots,K} \sigma_k(x; \mathbf{W})$$

**if  $\sigma_k(x; \mathbf{W}) > \sigma_j(x; \mathbf{W}) \ \forall j \neq k$  then decide  $C_k$**

## Multi-class logistic regression (cont.)

- $K > 2$  and  $y \in \{1, 2, \dots, K\}$

$$\sigma_k(x, \mathbf{W}) = P(y = k|x) = \frac{\exp(w_k^T x)}{\sum_{j=1}^K \exp(w_j^T x)}$$

- Normalized exponential (Aka **Softmax**)
- if  $w_k^T x \gg w_j^T x$  for all  $j \neq k$  then  $P(C_k|x) \approx 1$  and  $P(C_j|x) \approx 0$
- Note : remember from Bayes theorem:

$$P(C_k|x) = \frac{P(x|C_k)P(C_k)}{\sum_{j=1}^K P(x|C_j)P(C_j)}$$

## Multi-class logistic regression (cont.)

- Softmax function **smoothly** highlights the maximum probability and is differentiable.
- Compare it with `max(.)` function which is strict and non-differentiable
- Softmax can also handle negative values because we are using exponential function
- And it gives us probability for each class since:

$$\sum_{k=1}^K \frac{\exp(w_k^T x)}{\sum_{j=1}^K \exp(w_j^T x)} = 1$$

- An example of applying softmax (note that  $z_i = w^T x_i$ ):



## Multi-class logistic regression (cont.)

- Again we set  $J(W)$  as negative of log likelihood.
- We need  $\hat{W} = \arg \min_W J(W)$

$$\begin{aligned} J(W) &= -\log \prod_{i=1}^n P(y^{(i)} | x^{(i)}, \mathbf{W}) \\ &= -\log \prod_{i=1}^n \prod_{k=1}^K \sigma_k(x^{(i)}; \mathbf{W})^{y_k^{(i)}} \\ &= -\sum_{i=1}^n \sum_{k=1}^K y_k^{(i)} \log(\sigma_k(x^{(i)}; \mathbf{W})) \end{aligned}$$

- If **i-th** sample belongs to class  $k$  then  $y_k^{(i)}$  is 1 else 0.
- Again no closed-form solution for  $\hat{W}$











## Accuracy and Error Rate

**Accuracy:** proportion of correctly classified instances

$$\text{Accuracy} = \frac{TP + TN}{TP + TN + FP + FN}$$

**Error Rate:** complement of accuracy

$$\text{Error Rate} = \frac{FP + FN}{TP + TN + FP + FN} = 1 - \text{Accuracy}$$

- Accuracy is reliable for balanced class distributions.
- Limited interpretability for skewed datasets.

## Precision and Recall

### Precision (Positive Predictive Value)

$$P = \frac{TP}{TP + FP}$$

- Fraction of predicted positives that are truly positive.
- Sensitive to false positives.

### Recall (Sensitivity / True Positive Rate)

$$R = \frac{TP}{TP + FN}$$

- Fraction of actual positives correctly identified.
- Sensitive to false negatives.

## F1-Score and Trade-Off Analysis

**F1-Score:** harmonic mean of precision and recall

$$F1 = 2 \cdot \frac{P \cdot R}{P + R}$$

- Provides a single measure balancing precision and recall.
- Useful in imbalanced datasets or when both types of error are significant.
- Trade-off: Increasing recall typically decreases precision and vice versa.

## Comparative Overview

Metric	Formula	Focus	Typical Use-Case
Accuracy	$\frac{TP+TN}{TP+FP+FN+TN}$	Overall correctness	Balanced classes
Precision	$\frac{TP}{TP+FP}$	False positive control	Spam detection, IR
Recall	$\frac{TP}{TP+FN}$	False negative control	Medical diagnosis, anomaly detection
F1-Score	$2 \frac{PR}{P+R}$	Balance	Imbalanced or skewed datasets

- Metric choice should reflect operational objectives and misclassification costs.
- Always interpret metrics in conjunction with the confusion matrix.

## 1 Introduction

## ② Logistic Regression

### ③ Evaluation Metrics

#### 4 Extra reading

## Probabilistic view in classification

## Probabilistic classifiers

## 5 References



- 1 Introduction
- 2 Logistic Regression
- 3 Evaluation Metrics
- 4 **Extra reading**
  - Probabilistic view
  - Probabilistic classification
- 5 References

## Probabilistic view in classification

## Probabilistic view in classification problem

- In a classification problem:
  - Each **feature** is a **random variable** (e.g. a person's height)
  - The **class label** is also considered a **random variable** (e.g. a person could be overweight or not)
- We observe the feature values for a random sample and intend to find its class label
  - Evidence: Feature vector  $x$
  - Objective: Class label

## Definitions

- Posterior probability : The probability of a class label  $C_k$  given a sample  $x$

$$P(C_k|x)$$

- Likelihood or class conditional probability : PDF of feature vector  $x$  for samples of class  $C_k$

$$P(x|C_k)$$

- Prior probability : Probability of the label be  $C_k$

$$P(C_k)$$

- $P(x)$ : PDF of feature vector  $x$ 
  - From total probability theorem:

$$P(x) = \sum_{k=1}^K P(x|C_k)P(C_k)$$

## 1 Introduction

## 2 Logistic Regression

### ③ Evaluation Metrics

#### 4 Extra reading

## Probabilistic view in classification

## Probabilistic classifiers

## 5 References

## Probabilistic classifiers

- Probabilistic approaches can be divided in two main categories:
  - Generative
    - Estimate PDF  $P(x, C_k)$  for each class  $C_k$  and then use it to find  $P(C_k|x)$ . Alternatively estimate both PDF  $P(x|C_k)$  and  $P(C_k)$  to find  $P(C_k|x)$ .
  - Discriminative
    - Directly estimate  $P(C_k|x)$  for class  $C_k$

## Probabilistic classifiers (cont.)

- Let's assume we have input data  $x$  and want to classify the data into labels  $y$ .
- A generative model learns the **joint** probability distribution  $P(x, y)$ .
- A discriminative model learns the **conditional** probability distribution  $P(y|x)$

## Discriminative vs. Generative : example

- Suppose we have the following dataset in form of  $(x, y)$ :

 $(1, 0), (1, 0), (2, 0), (2, 1)$ 

- $P(x, y)$  is :

	$y = 0$	$y = 1$
$x = 1$	$\frac{1}{2}$	0
$x = 2$	$\frac{1}{4}$	$\frac{1}{4}$

- $P(y|x)$  is :

	$y = 0$	$y = 1$
$x = 1$	1	0
$x = 2$	$\frac{1}{2}$	$\frac{1}{2}$





## Generative approach

## 1 Inference

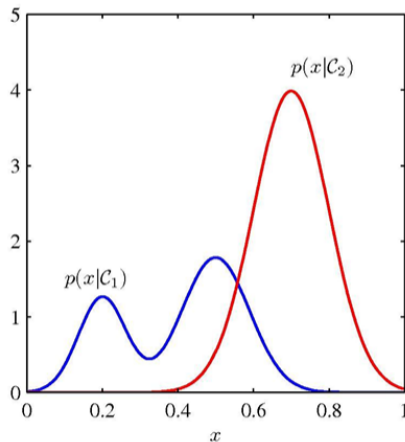
- Determine class conditional densities  $P(x|C_k)$  and priors  $P(C_k)$
- Use Bayes theorem to find  $P(C_k|x)$

## 2 Decision

- Make optimal assignment for new input (after learning the model in the inference stage)
- if  $P(C_i|x) > P(C_j|x) \forall j \neq i$ , then decide  $C_i$ .

## Generative approach (cont.)

- Generative approach for a binary classification problem:



Figures adapted from Machine Learning and Pattern Recognition, Bishop

## Discriminative approach

## 1 Inference

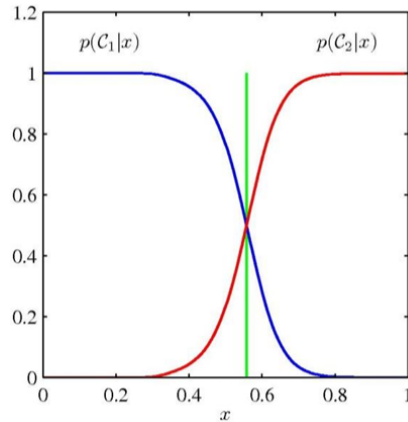
- Determine the posterior class probabilities  $P(C_k|x)$  directly.

## ② Decision

- Make optimal assignment for new input (after learning the model in the inference stage)
- if  $P(C_i|x) > P(C_j|x) \forall j \neq i$ , then decide  $C_i$ .

## Discriminative approach (cont.)

- Discriminative approach for a binary classification problem:



Figures adapted from Machine Learning and Pattern Recognition, Bishop

- 1 Introduction
- 2 Logistic Regression
- 3 Evaluation Metrics
- 4 Extra reading
- 5 References

## Contributions

- **These slides are authored by:**
  - Danial Gharib

