

Introduction to Machine Learning

Lecture 5 - Supervised Learning

1 Introduction

Consider the **training set**:

$$\mathcal{D} = \{\mathbf{x}_i, y_i\}_{i=1}^N$$

where:

- $\mathbf{x}_i \in \mathcal{X} \subseteq \mathbb{R}^D$ is the features vector (also known as a sample \ observation)
- $y_i \in \mathcal{Y} = \{C_1, C_2, \dots, C_K\}$ or just $y_i \in \{1, 2, \dots, K\}$ is the class (category) of \mathbf{x}_i .

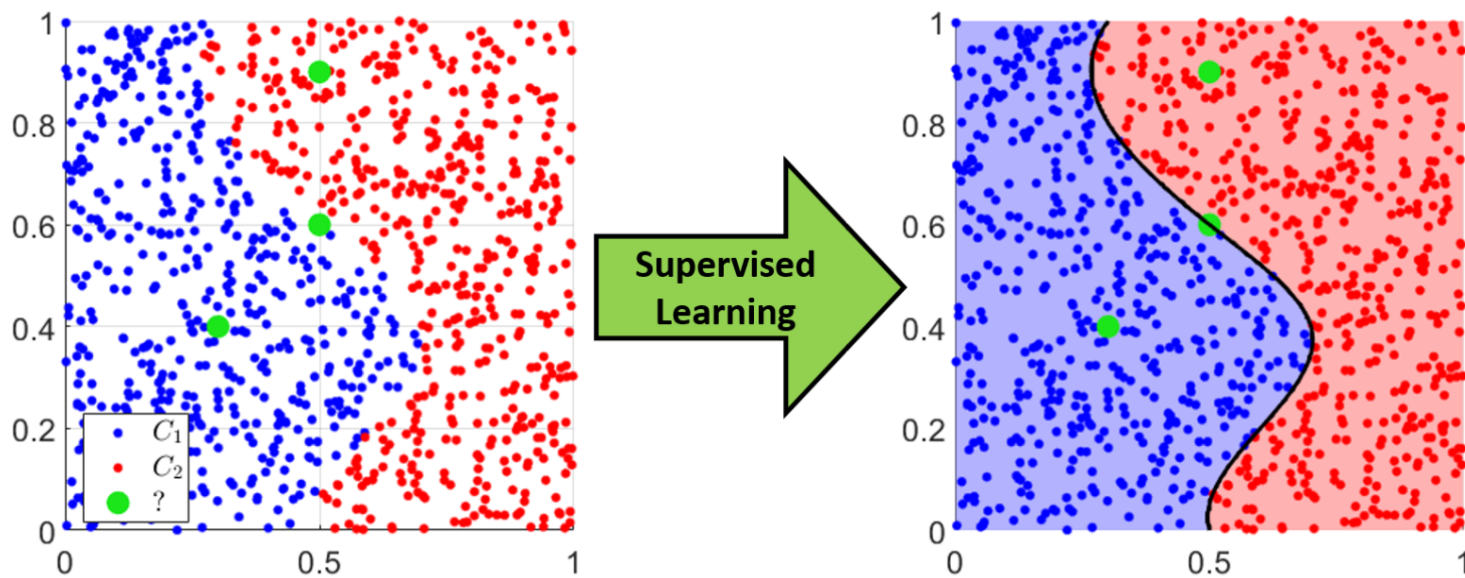
Consider a new pair (\mathbf{x}_0, y_0) where \mathbf{x}_0 is known and y_0 is unknown.

Given the set \mathcal{D} , our goal is to derive a classifier function $f: \mathcal{X} \rightarrow \mathcal{Y}$, such that:

$$y_0 = f(\mathbf{x}_0)$$

For example, \mathcal{D} is the set of all **red** and **blue** points (left figure).

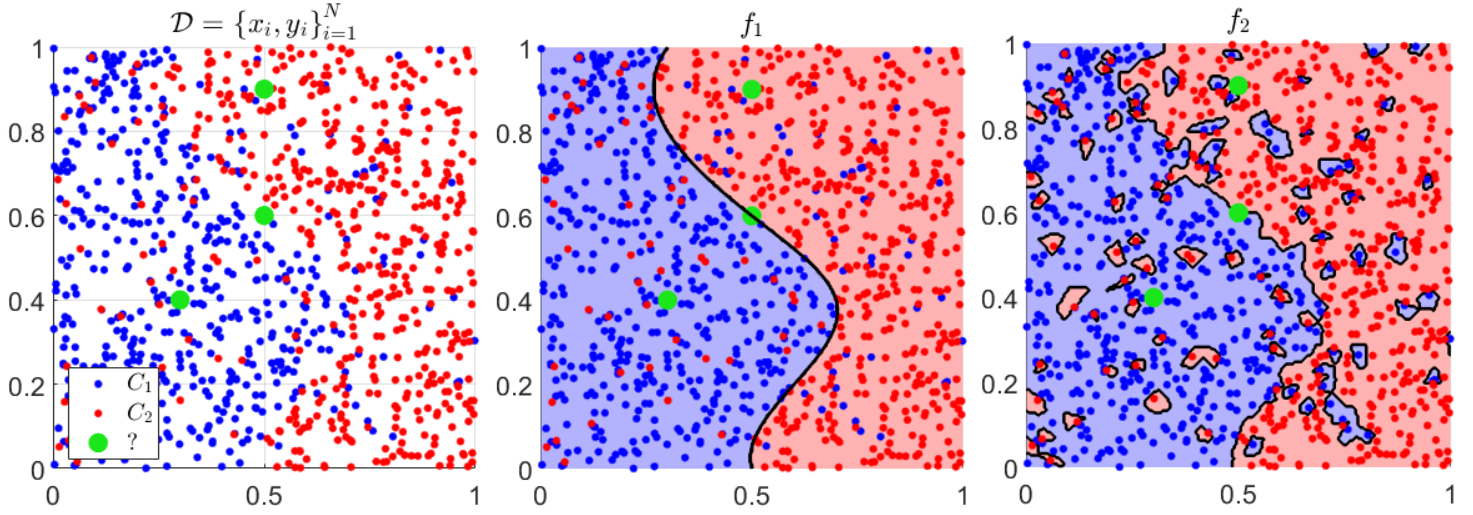
Can we guess the classes of the three **green** points?



Notes:

- For $\mathcal{Y} = \{0, 1\}$, the problem is known as **binary classification**.
- If \mathcal{Y} is a continuous set we called this problem **regression**.

In actual problems the data might not be perfectly separable:



and learning a perfect f is impossible.

2 Quality Index

2.1 Loss Function

A loss function $\ell : \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}$ satisfies:

1.

$$\ell(y, y) = 0, \quad \forall y \in \mathcal{Y}$$

2.

$$\ell(\hat{y}, y) \geq 0$$

Examples

1. Hamming loss:

$$\ell(\hat{y}, y) = \mathbf{I}\{\hat{y} \neq y\} \triangleq \begin{cases} 1 & \hat{y} \neq y \\ 0 & \text{else} \end{cases}$$

2. Squared loss (for continuous \mathcal{Y}):

$$\ell(\hat{y}, y) = (\hat{y} - y)^2$$

2.2 Empirical Risk (Training and Test Loss)

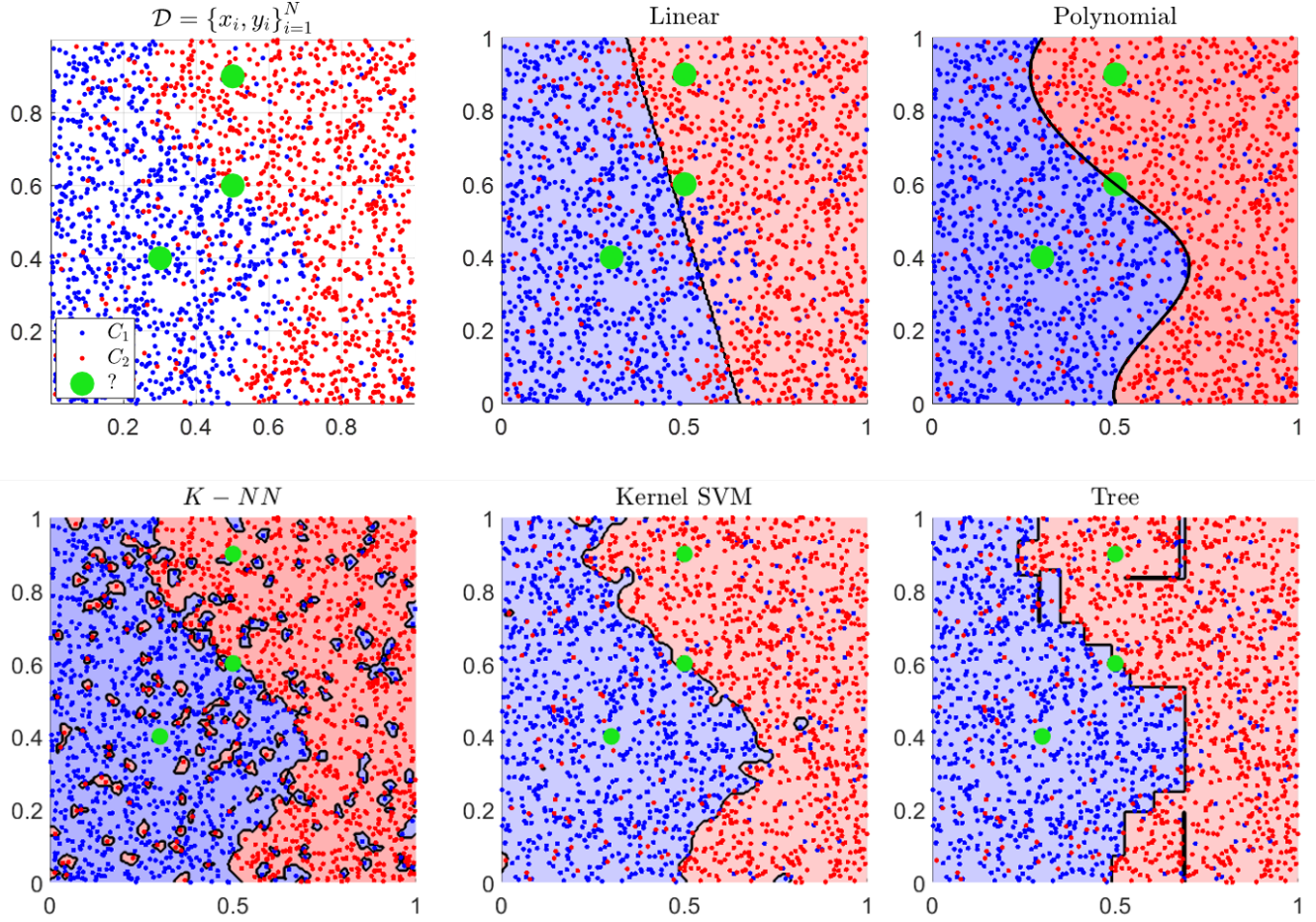
Let $f : \mathcal{X} \rightarrow \mathcal{Y}$ be a classifier function.

Based on a loss function ℓ ,

we can evaluate the performance of f on a set $\mathcal{D} = \{(\mathbf{x}_i, y_i)\}$ by:

$$L(f) \triangleq \frac{1}{N} \sum_{i=1}^N \ell(f(\mathbf{x}_i), y_i)$$

Example – Models and decision regions



3 Linear Binary Classification Using Perceptron

Let $\mathcal{D} = \{x_i, y_i\}_{i=1}^N$ be a training set where:

- $x_i \in \mathcal{X} = \mathbb{R}^D$.
- $y_i \in \mathcal{Y} = \{-1, 1\}$ (binary classification)

We search for a linear classifier:

$$f : \mathcal{X} \rightarrow \mathcal{Y}$$

such that:

$$f(x) = \text{sign}(w^T x - b) = \begin{cases} 1 & w^T x - b \geq 0 \\ -1 & w^T x - b < 0 \end{cases}$$

where $w \in \mathbb{R}^D$ and $b \in \mathbb{R}$ are the classifier model parameters.

- Notice that for any value of C :

$$\text{sign}(w^T x - b) = \text{sign}\left(\frac{w^T x}{C} - \frac{b}{C}\right), \quad \forall x \in \mathbb{R}^D$$

Thus, we can assume (without loss of generality) that $\|w\| = 1$.

3.1 Analysis

Consider the function:

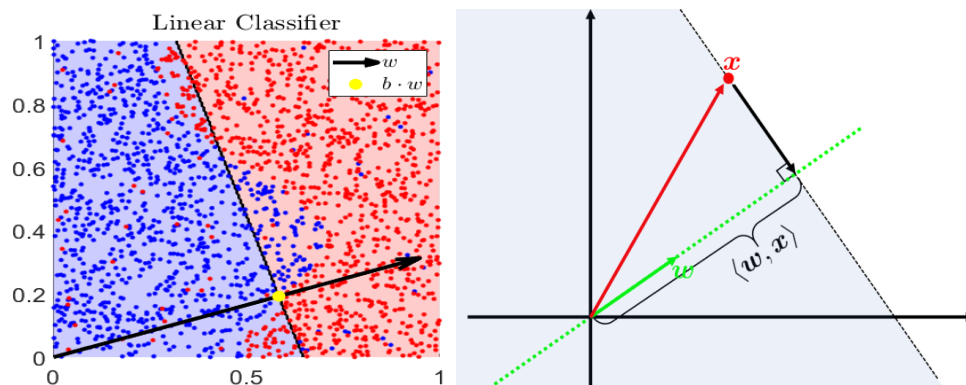
$$f(\mathbf{x}) = \text{sign}(\mathbf{w}^T \mathbf{x} - b), \quad \|\mathbf{w}\| = 1$$

This function splits the input space \mathcal{X} into two half-planes.

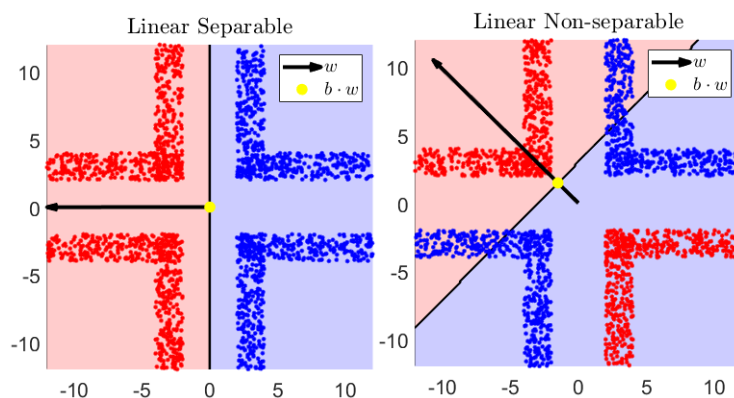
The decision boundary is given by:

$$\mathbf{w}^T \mathbf{x} - b = 0$$

$$\langle \mathbf{w}, \mathbf{x} \rangle = b$$



Linear separable and linear non-separable



3.2 The Perceptron Algorithm

For convenience purposes, we denote:

$$\tilde{\mathbf{x}}_i \triangleq \begin{bmatrix} 1 \\ \mathbf{x}_i \end{bmatrix}, \quad \tilde{\mathbf{w}} = \begin{bmatrix} -b \\ \mathbf{w} \end{bmatrix}$$

So, we have:

$$\Rightarrow \boxed{\mathbf{w}^T \mathbf{x} - b = \tilde{\mathbf{w}}^T \tilde{\mathbf{x}}}$$

Algorithm 1 The Perceptron Algorithm

Input: Training set $\mathcal{D} = \{\tilde{\mathbf{x}}_i, y_i\}$

Output: The linear classifiers parameters: $\tilde{\mathbf{w}}$ (that is, \mathbf{w} and b).

1. Set $\tilde{\mathbf{w}}_1$ with some initial guess.

2. **for** $k = 1, 2, 3, \dots$

(a) Choose some $(\tilde{\mathbf{x}}_k, y_k) \in \mathcal{D}$

(b) Compute:

$$\hat{y}_k = \text{sign}(\tilde{\mathbf{w}}_k^T \tilde{\mathbf{x}}_k)$$

(c) Update:

$$\tilde{\mathbf{w}}_{k+1} = \tilde{\mathbf{w}}_k + \frac{1}{2} (y_k - \hat{y}_k) \tilde{\mathbf{x}}_k$$

Notes:

- If $\hat{y}_k = y_k$, there is no update: $\tilde{\mathbf{w}}_{k+1} = \tilde{\mathbf{w}}_k$
- If $\hat{y}_k = -1$ and $y_k = 1$, then, the value $\tilde{\mathbf{w}}_k^T \tilde{\mathbf{x}}_k < 0$ is too small (and we should increase it).
 - Step (c) increases the value of $\tilde{\mathbf{w}}_k^T \tilde{\mathbf{x}}_k$:

$$\tilde{\mathbf{w}}_{k+1} = \tilde{\mathbf{w}}_k + \frac{1}{2} (y_k - \hat{y}_k) \tilde{\mathbf{x}}_k = \tilde{\mathbf{w}}_k + \tilde{\mathbf{x}}_k$$

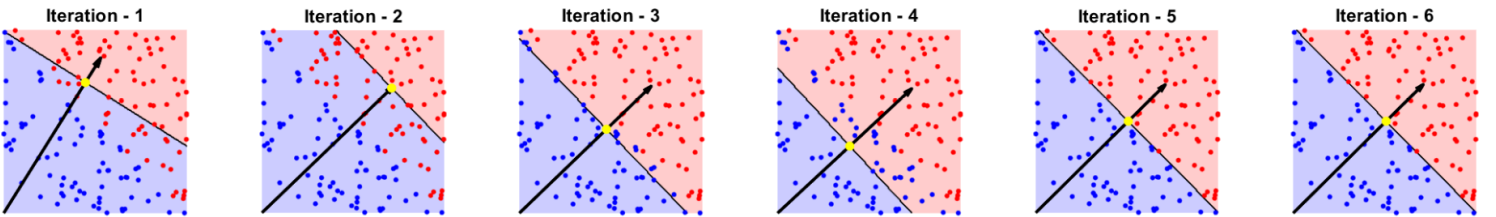
$$\Rightarrow \boxed{\tilde{\mathbf{w}}_{k+1}^T \tilde{\mathbf{x}}_k = (\tilde{\mathbf{w}}_k + \tilde{\mathbf{x}}_k)^T \tilde{\mathbf{x}}_k = \tilde{\mathbf{w}}_k^T \tilde{\mathbf{x}}_k + \|\tilde{\mathbf{x}}_k\|_2^2 > \tilde{\mathbf{w}}_k^T \tilde{\mathbf{x}}_k}$$

- In the same manner, the algorithm updates the miss-classification $\hat{y}_k = 1$ and $y_k = -1$.
- The algorithm converge in finite number of iterations if the problem is linear separable.
- In the linear non-separable case, there is no guarantee for convergence.

Example

Let $\mathcal{D} = \{\mathbf{x}_i, y_i\}_{i=1}^N$ with $N = 150$.

We apply the algorithm for $5N$ iterations, namely, we run over \mathcal{D} , 5 times:



4 Quadratic Classifier Using Bayesian Approach

4.1 Introduction

Given the statistic $p_{X,Y}$ we can use the MAP classifier:

$$f_{MAP}(\mathbf{x}) = \arg \max_{C_k \in \mathcal{Y}} p(\mathbf{x}|C_k) P_Y(C_k)$$

In practice, the statistic is usually unknown.

However, we can use the training set $\mathcal{D} = \{\mathbf{x}_i, y_i\}_{i=1}^N$ to estimate it.

Prior Estimation The estimation of P_Y is given by:

$$\hat{P}_Y(C_k) = \frac{1}{N} \sum_{i=1}^N \mathbf{I}\{y_i = C_k\}$$

Likelihood Estimation Lets assume that given C_k :

$$\mathbf{x}|C_k \sim \mathcal{N}(\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k), \quad \forall C_k \in \mathcal{Y}$$

Thus, we can estimate the mean and covariance of each class:

$$\begin{aligned} \hat{\boldsymbol{\mu}}_k &= \frac{1}{N_k} \sum_{\mathbf{x}_i \in C_k} \mathbf{x}_i, \quad N_k \triangleq \sum_{i=1}^N \mathbf{I}\{y_i = C_k\} \\ \hat{\boldsymbol{\Sigma}}_k &= \frac{1}{N_k} \sum_{\mathbf{x}_i \in C_k} (\mathbf{x}_i - \hat{\boldsymbol{\mu}}_k)(\mathbf{x}_i - \hat{\boldsymbol{\mu}}_k)^T \\ \Rightarrow p(\mathbf{x}|C_k) &= |2\pi\boldsymbol{\Sigma}_k|^{-\frac{1}{2}} \exp\left(-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu}_k)^T \boldsymbol{\Sigma}_k^{-1}(\mathbf{x} - \boldsymbol{\mu}_k)\right) \end{aligned}$$

4.2 Binary Case Analysis

Consider the binary case: $\mathcal{Y} = \{C_1, C_2\}$.

The MAP classifier is given by:

$$\begin{aligned} p(\mathbf{x}|C_1) P_Y(C_1) &\underset{C_2}{\overset{C_1}{\gtrless}} p(\mathbf{x}|C_2) P_Y(C_2) \\ p_1 p(\mathbf{x}|C_1) &\underset{C_2}{\overset{C_1}{\gtrless}} (1 - p_1) p(\mathbf{x}|C_2), \quad p_1 \triangleq P_Y(C_1) \end{aligned}$$

In other words, the boundary decision is given by:

$$\begin{aligned}
 p_1 p(\mathbf{x}|C_1) &= (1 - p_1) p(\mathbf{x}|C_2) \\
 p_1 |2\pi \Sigma_1|^{-\frac{1}{2}} \exp\left(-\frac{1}{2}(\mathbf{x} - \mu_1)^T \Sigma_1^{-1}(\mathbf{x} - \mu_1)\right) &= (1 - p_1) |2\pi \Sigma_2|^{-\frac{1}{2}} \exp\left(-\frac{1}{2}(\mathbf{x} - \mu_2)^T \Sigma_2^{-1}(\mathbf{x} - \mu_2)\right) \\
 p_1 d_1 \exp\left(-\frac{1}{2}(\mathbf{x} - \mu_1)^T \Sigma_1^{-1}(\mathbf{x} - \mu_1)\right) &= (1 - p_1) d_2 \exp\left(-\frac{1}{2}(\mathbf{x} - \mu_2)^T \Sigma_2^{-1}(\mathbf{x} - \mu_2)\right), \quad d_{1,2} \triangleq |\Sigma_{1,2}|^{-\frac{1}{2}} \\
 \log(p_1 d_1) - \frac{1}{2}(\mathbf{x} - \mu_1)^T \Sigma_1^{-1}(\mathbf{x} - \mu_1) &= \log((1 - p_1) d_2) - \frac{1}{2}(\mathbf{x} - \mu_2)^T \Sigma_2^{-1}(\mathbf{x} - \mu_2) \\
 \underbrace{2 \log\left(\frac{p_1 d_1}{(1 - p_1) d_2}\right)}_{\triangleq b} - (\mathbf{x} - \mu_1)^T \Sigma_1^{-1}(\mathbf{x} - \mu_1) &= -(\mathbf{x} - \mu_2)^T \Sigma_2^{-1}(\mathbf{x} - \mu_2) \\
 b - (\mathbf{x} - \mu_1)^T \Sigma_1^{-1}(\mathbf{x} - \mu_1) + (\mathbf{x} - \mu_2)^T \Sigma_2^{-1}(\mathbf{x} - \mu_2) &= 0 \\
 \underbrace{b - \mu_1^T \Sigma_1^{-1} \mu_1 + \mu_2^T \Sigma_2^{-1} \mu_2}_{\triangleq \tilde{b}} - \underbrace{2(\mu_2^T \Sigma_2^{-1} - \mu_1^T \Sigma_1^{-1}) \mathbf{x}}_{\triangleq \mathbf{w}^T} + \underbrace{\mathbf{x}^T (\Sigma_2^{-1} - \Sigma_1^{-1}) \mathbf{x}}_{\triangleq \mathbf{M}} &= 0 \\
 \Rightarrow \boxed{\tilde{b} - \mathbf{w}^T \mathbf{x} + \mathbf{x}^T \mathbf{M} \mathbf{x} = 0}
 \end{aligned}$$

This is a **quadratic** decision boundary.

Examples Consider two states with equal a-priori probability:

$$\mathcal{Y} = \{C_1, C_2\}$$

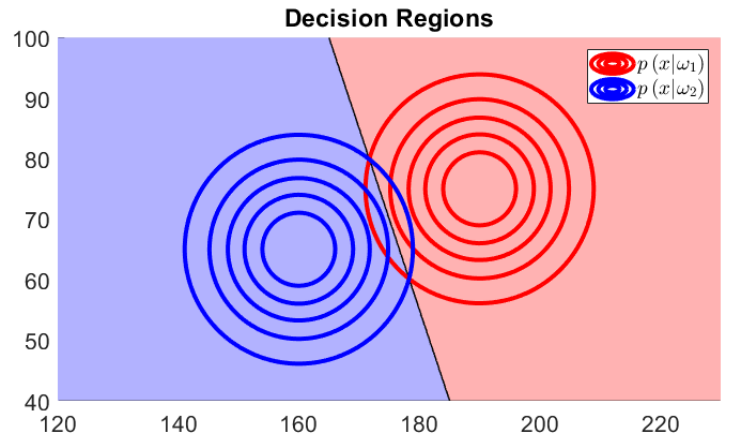
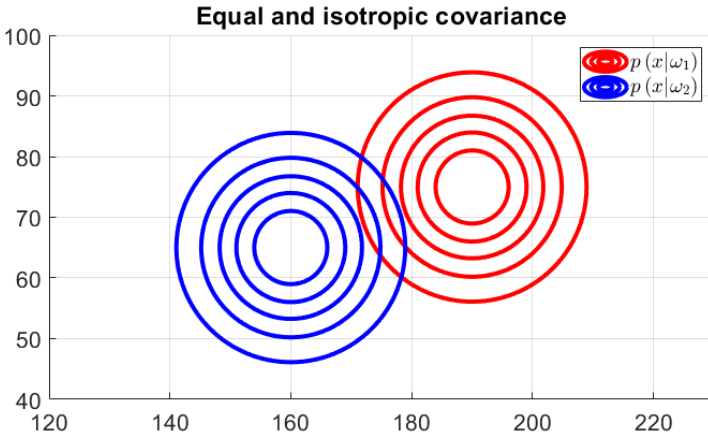
$$P_Y(C_1) = P_Y(C_2) = \frac{1}{2}$$

Given each state C_i , the observation X is a random Gaussian vector (with Σ_i and μ_i):

$$p_{X|Y}(\mathbf{x}|C_i) = \frac{1}{2\pi} |\Sigma_i|^{-1} \exp\left(-\frac{1}{2}(\mathbf{x} - \mu_i)^T \Sigma_i^{-1}(\mathbf{x} - \mu_i)\right)$$

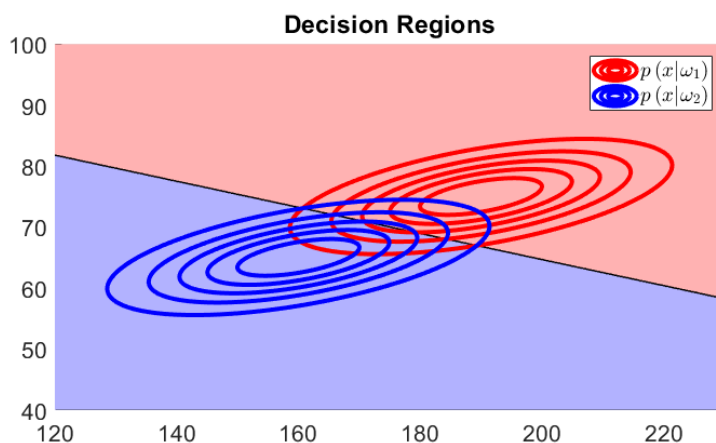
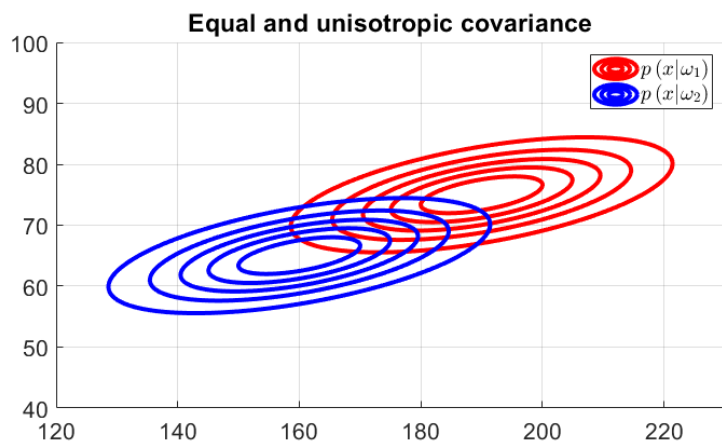
Equal isotropic covariance

$$\begin{cases} \Sigma_1 = \Sigma_2 = \mathbf{I} \\ \mu_1 \neq \mu_2 \end{cases}$$



Equal covariances

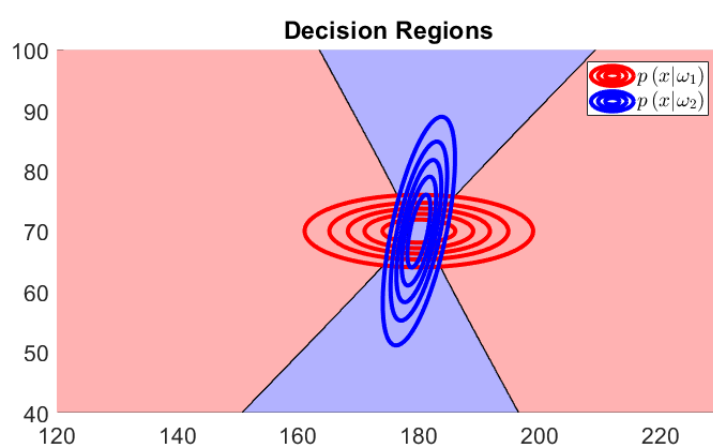
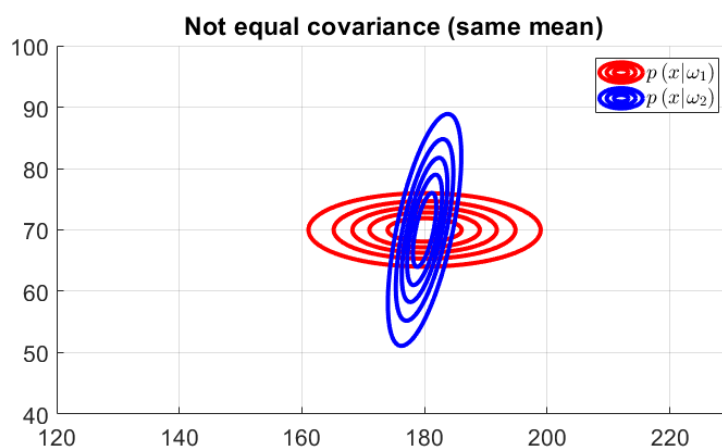
$$\begin{cases} \Sigma_1 = \Sigma_2 \\ \mu_1 \neq \mu_2 \end{cases}$$



Notice that when $\Sigma_1 = \Sigma_2$ the decision boundary is linear (a single straight line).

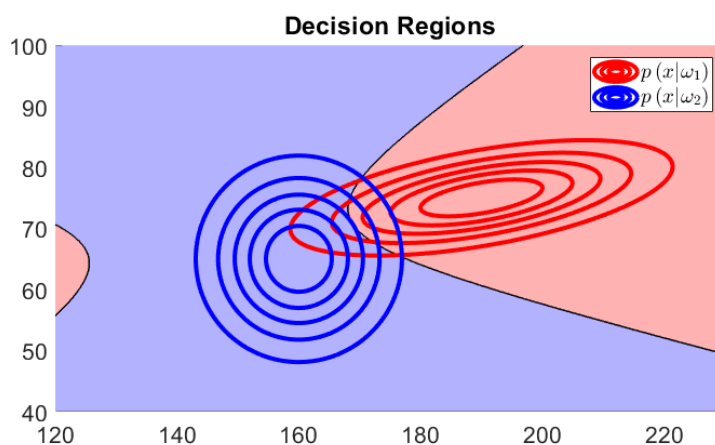
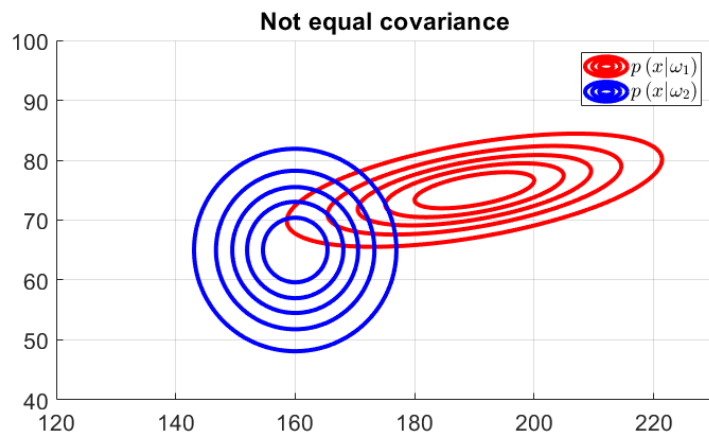
Equal mean

$$\begin{cases} \Sigma_1 \neq \Sigma_2 \\ \mu_1 = \mu_2 \end{cases}$$

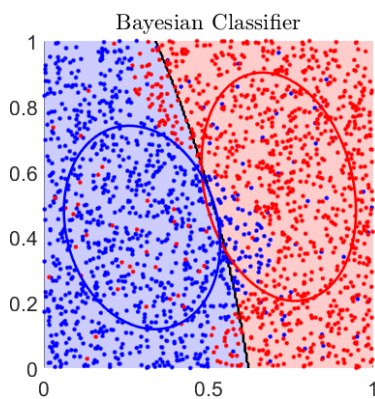


Non-equal covariance

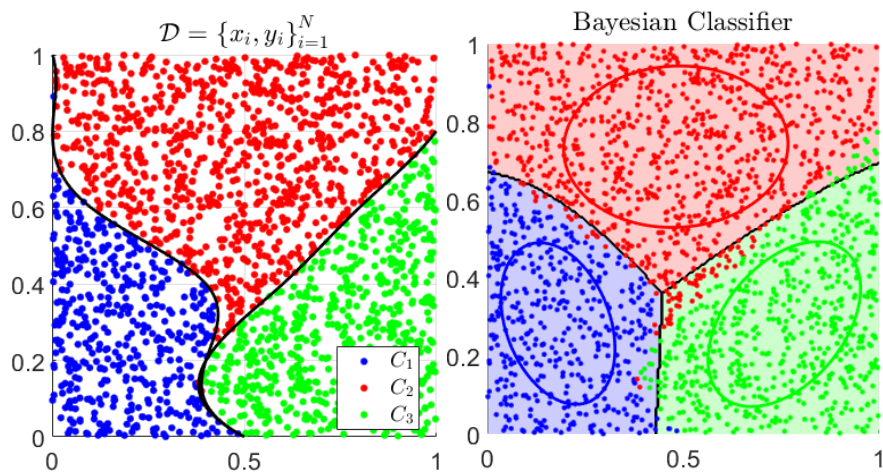
$$\begin{cases} \Sigma_1 \neq \Sigma_2 \\ \mu_1 \neq \mu_2 \end{cases}$$



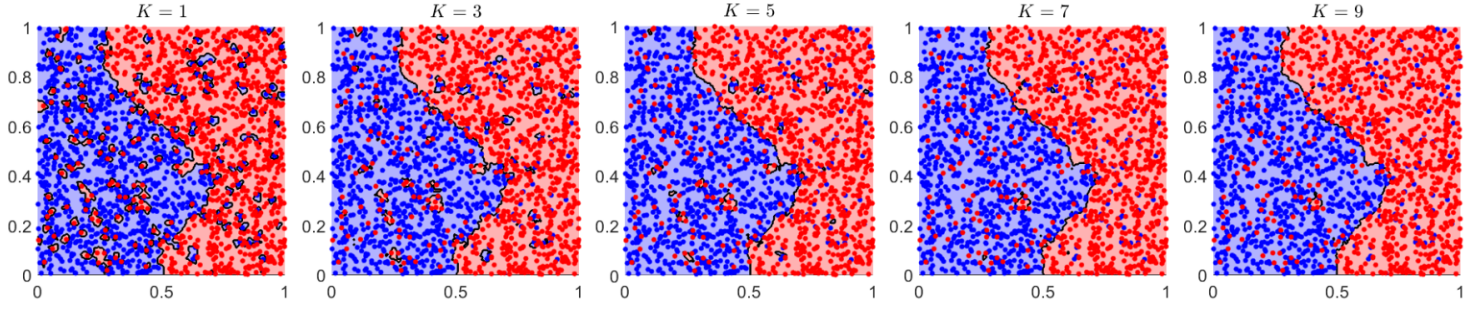
Binary case – with estimation



Trinary case – with estimation



5 K Nearest Neighbors (K -NN)



Let $\mathcal{D} = \{\mathbf{x}_i, y_i\}_{i=1}^N$ be a training set.

5.1 1-NN ($K = 1$)

The 1-NN classifier is given by:

$$\hat{y} = f_{NN}(\mathbf{x}) = y_{k(\mathbf{x})}$$

such that:

$$k(\mathbf{x}) = \arg \min_{i \in \{1, 2, \dots, N\}} d(\mathbf{x}, \mathbf{x}_i)$$

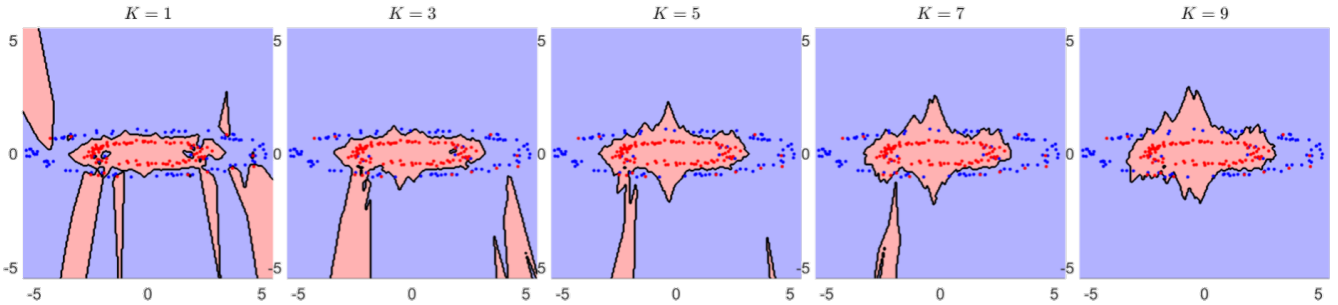
where $d: \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ is a metric function (for example: $d(\mathbf{x}, \mathbf{x}_i) = \|\mathbf{x} - \mathbf{x}_i\|_2$).
In words, a new sample \mathbf{x} will be labeled by the class of the nearest neighbor $\mathbf{x}_i \in \mathcal{D}$.

5.2 K -NN

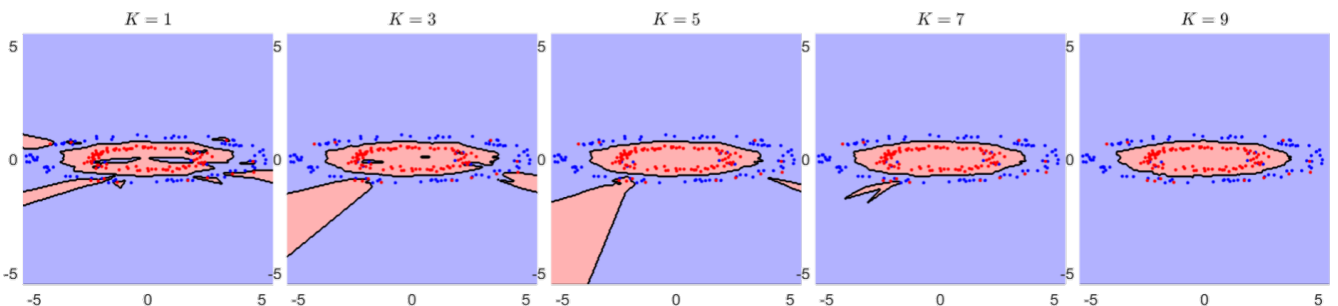
For $K > 1$, a new sample \mathbf{x} will be labeled by the class most common among its K nearest neighbors.

K -NN with different metric (example)

Euclidean distance: $\|\mathbf{x}_i - \mathbf{x}_j\|_2$



Mahalanobis distance: $\|\mathbf{x}_i - \mathbf{x}_j\|_{\Sigma_x} \triangleq (\mathbf{x}_i - \boldsymbol{\mu}_x)^T \Sigma_x^{-1} (\mathbf{x}_i - \boldsymbol{\mu}_x)$



5.3 Lower and Upper Performance Bounds (Extra)

Consider the case of 1-Nearest Neighbors.

The optimal MAP estimator is given by:

$$C_m \triangleq \arg \max_{C_i \in \{1,2,\dots,K\}} P(C_i|\mathbf{x}), \quad K\text{-number of classes}$$

The error of the MAP classifier is given by:

$$P^*(\text{error}|\mathbf{x}) = 1 - P(C_m|\mathbf{x})$$

We denote the error of the 1-NN classifier by:

$$P(\text{error}|\mathbf{x}) = 1 - P(f_{NN}(\mathbf{x})|\mathbf{x})$$

Lemma 1. *In the limit $N \rightarrow \infty$:*

(We also assume that p_X is continuous and non-zero everywhere)

$$P^*(\text{error}|\mathbf{x}) \leq P(\text{error}|\mathbf{x}) \leq 2P^*(\text{error}|\mathbf{x})$$

Proof. Given a new point \mathbf{x} with the (unknown) label $y \in \mathcal{Y} = \{C_1, C_2, \dots, C_K\}$.

By the optimality of the MAP classifier:

$$\boxed{P^*(\text{error}|\mathbf{x}) \leq P(\text{error}|\mathbf{x})}$$

Since $N \rightarrow \infty$ and p_X is non-zeros everywhere there is some pair $(\mathbf{x}_0, y_0) \in \mathcal{D}$ such that $\mathbf{x}_0 = \mathbf{x}$.

Thus:

$$f_{NN}(\mathbf{x}) = y_0$$

Thus, the error is given by:

$$\begin{aligned} P(\text{error}|\mathbf{x}) &= P(y \neq y_0|\mathbf{x}, \mathbf{x}_0) \\ &= P(y \neq y_0|\mathbf{x}) \\ &= 1 - \sum_{i=1}^C P(y = C_i, y_0 = C_i|\mathbf{x}) \\ &= 1 - \sum_{i=1}^C P_y^2(C_i|\mathbf{x}) \\ &\stackrel{(*)}{\leq} 2P^*(\text{error}|\mathbf{x}) \end{aligned}$$

where (*):

$$\begin{aligned} (1 - P(C_i|\mathbf{x}))^2 &\geq 0 \geq -\sum_{j \neq i} P^2(C_j|\mathbf{x}) \\ \Rightarrow (1 - P(C_i|\mathbf{x}))^2 &\geq -\sum_{j \neq i} P^2(C_j|\mathbf{x}) \\ 1 - 2P(C_i|\mathbf{x}) + P^2(C_i|\mathbf{x}) &\geq -\sum_{j \neq i} P^2(C_j|\mathbf{x}) \\ 1 - 2P(C_i|\mathbf{x}) &\geq -\sum_j P^2(C_j|\mathbf{x}) \\ 2 - 2P(C_i|\mathbf{x}) &\geq 1 - \sum_j P^2(C_j|\mathbf{x}) \\ 2P^*(\text{error}|\mathbf{x}) &\geq 1 - \sum_j P^2(C_j|\mathbf{x}) \end{aligned}$$

Overall:

$$\Rightarrow \boxed{P(\text{error}|\mathbf{x}) \leq 2P^*(\text{error}|\mathbf{x})}$$

□