

Course 3

Introduction to Machine Learning

Christophe Eloy

Multivariate linear regression

- Hypothesis: $h_{\theta}(x) = \theta^T x$
- Parameters: $\theta = [\theta_0 \dots \theta_n]^T$
- Cost function (least squares): $J(\theta) = \frac{1}{2N} \sum_{i=1}^N (h_{\theta}(x^{(i)}) - y^{(i)})^2$
- Goal: find $\min_{\theta} J(\theta)$
- Gradient descent: $\theta := \theta - \alpha \nabla_{\theta} J(\theta)$
- Normal equation: $\theta = (X^T X)^{-1} X^T y$

Outline

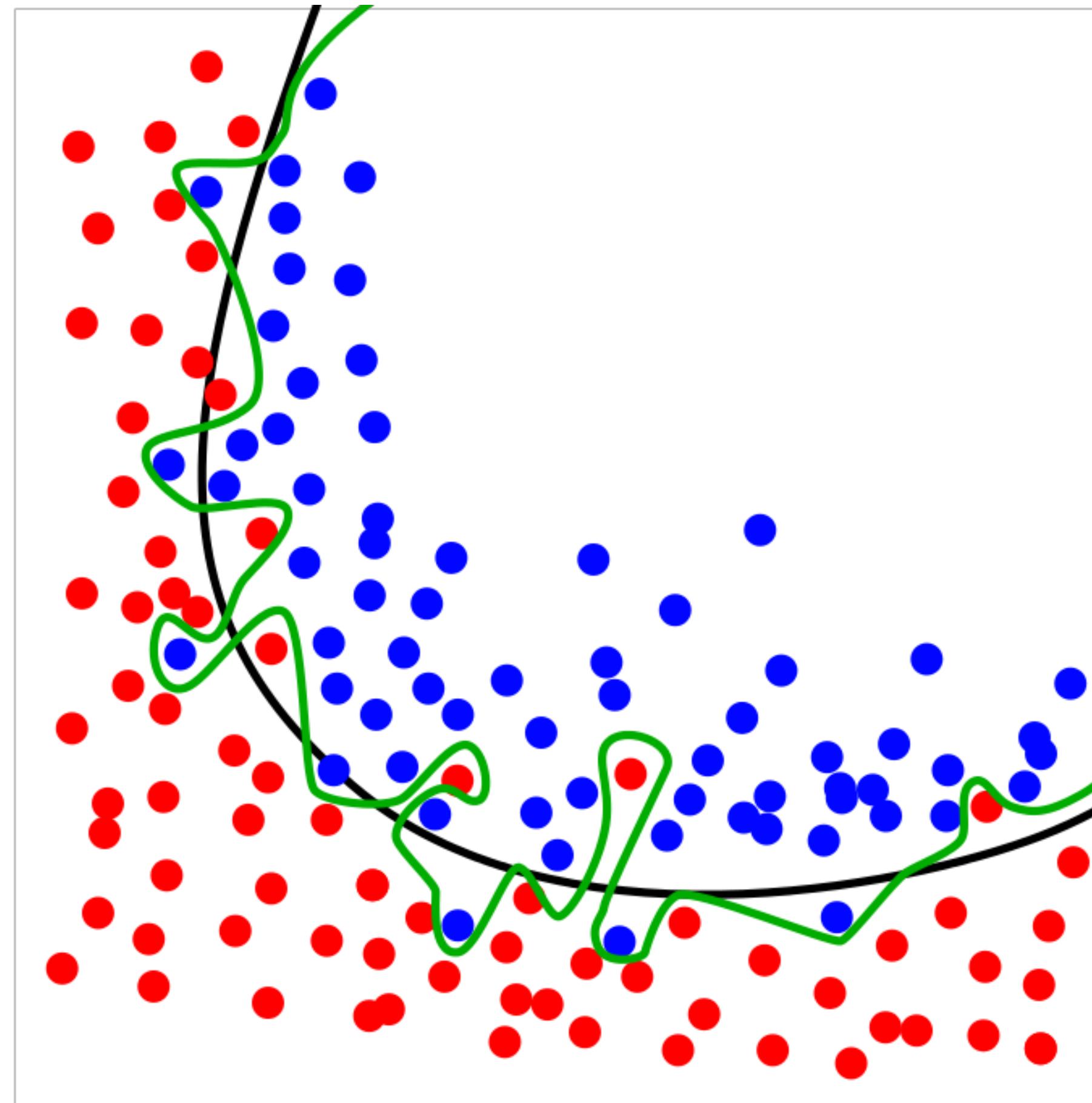
- Logistic regression
- Regularization
- Support vector machines (SVMs)
- Naive Bayes
- Decision trees
- k-means

Logistic regression

- Hypothesis:
$$h_{\theta}(x) = \frac{1}{1 + e^{-\theta^T x}}$$
- Parameters: $\theta = [\theta_0 \dots \theta_n]^T$
- Cost function:
$$J(\theta) = -\frac{1}{N} \sum_{i=1}^N \left[y^{(i)} \log(h_{\theta}(x^{(i)})) + (1 - y^{(i)}) \log(1 - h_{\theta}(x^{(i)})) \right]$$
- Gradient descent:
$$\theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta)$$

$$\theta_j := \theta_j - \alpha \frac{1}{N} \sum_{i=1}^N (h_{\theta}(x^{(i)}) - y^{(i)}) x_j^{(i)}$$

Problem of overfitting



- To avoid overfitting, one method is regularization

Regularization

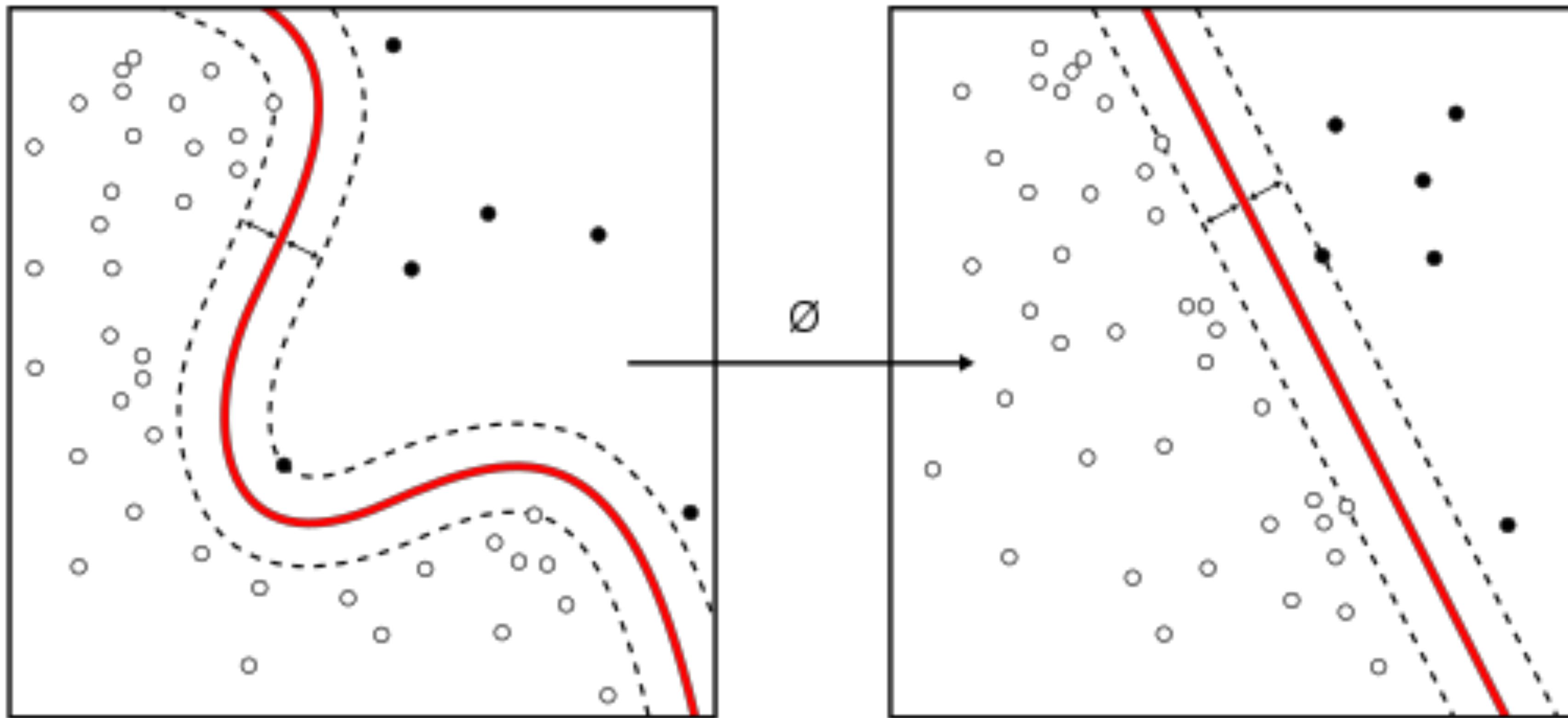
- New cost function:
$$J(\theta) = \frac{1}{2N} \left[\sum_{i=1}^N (h_\theta(\mathbf{x}^{(i)}) - y^{(i)})^2 + \lambda \sum_{j=1}^n \theta_j^2 \right]$$
- Calculation of the gradient

$$\frac{\partial J}{\partial \theta_0} = \frac{1}{N} \sum_{i=1}^N (h_\theta(\mathbf{x}^{(i)}) - y^{(i)}) x_0^{(i)}$$

$$\frac{\partial J}{\partial \theta_j} = \frac{1}{N} \sum_{i=1}^N (h_\theta(\mathbf{x}^{(i)}) - y^{(i)}) x_j^{(i)} + \frac{\lambda}{N} \theta_j$$

- Normal equation
$$\theta = \left(X^T X + \lambda \begin{bmatrix} 0 & & & \\ & 1 & & \\ & & \ddots & \\ & & & 1 \end{bmatrix} \right)^{-1} X^T y$$

Support vector machines



- We want “fat” margins

Support vector machines

- Support vectors (touching the margin) such that $y^{(i)} (\mathbf{a}^T \mathbf{x}_n + b) = 1$
- Size of each margin is $d = 1/\|\mathbf{a}\|$
- Optimization problem: Find $\min(\|\mathbf{a}\|^2)$ subject to $y^{(i)} (\mathbf{a}^T \mathbf{x} + b) \geq 1$
- Soft margin: Find

$$\min \left(\lambda \|\mathbf{a}\|^2 + \frac{1}{n} \sum_{i=1}^N \max \left(0, 1 - y^{(i)} (\mathbf{a}^T \mathbf{x} + b) \right) \right)$$

- Non-linear kernels: \mathbf{x} can be replaced by $\mathbf{f}(\mathbf{x})$

Naive Bayes

- Bayes theorem

$$p(C_k | \mathbf{x}) = \frac{p(C_k) p(\mathbf{x} | C_k)}{p(\mathbf{x})}$$

- “Naive” assumption of independence: $p(C_k, \mathbf{x}) = p(x_1 | C_k) \cdots p(x_n | C_k) p(C_k)$
- Conditional probability

$$p(C_k | \mathbf{x}) = \frac{p(C_k)}{p(\mathbf{x})} \prod_{i=1}^N p(x_i | C_k)$$

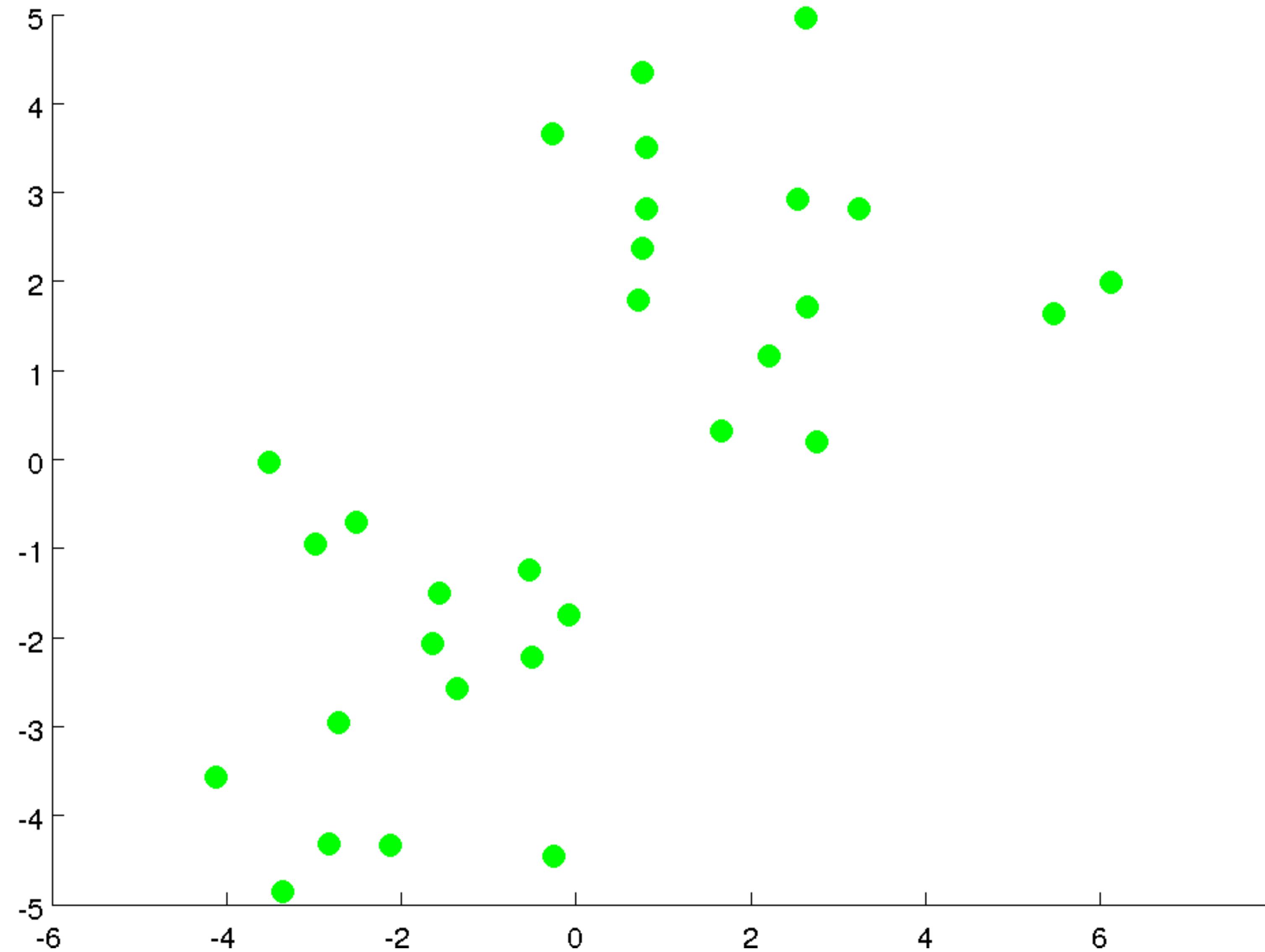
- Naive Bayes classifier: $y = \operatorname{argmax}_k p(C_k) \prod_{i=1}^N p(x_i | C_k)$

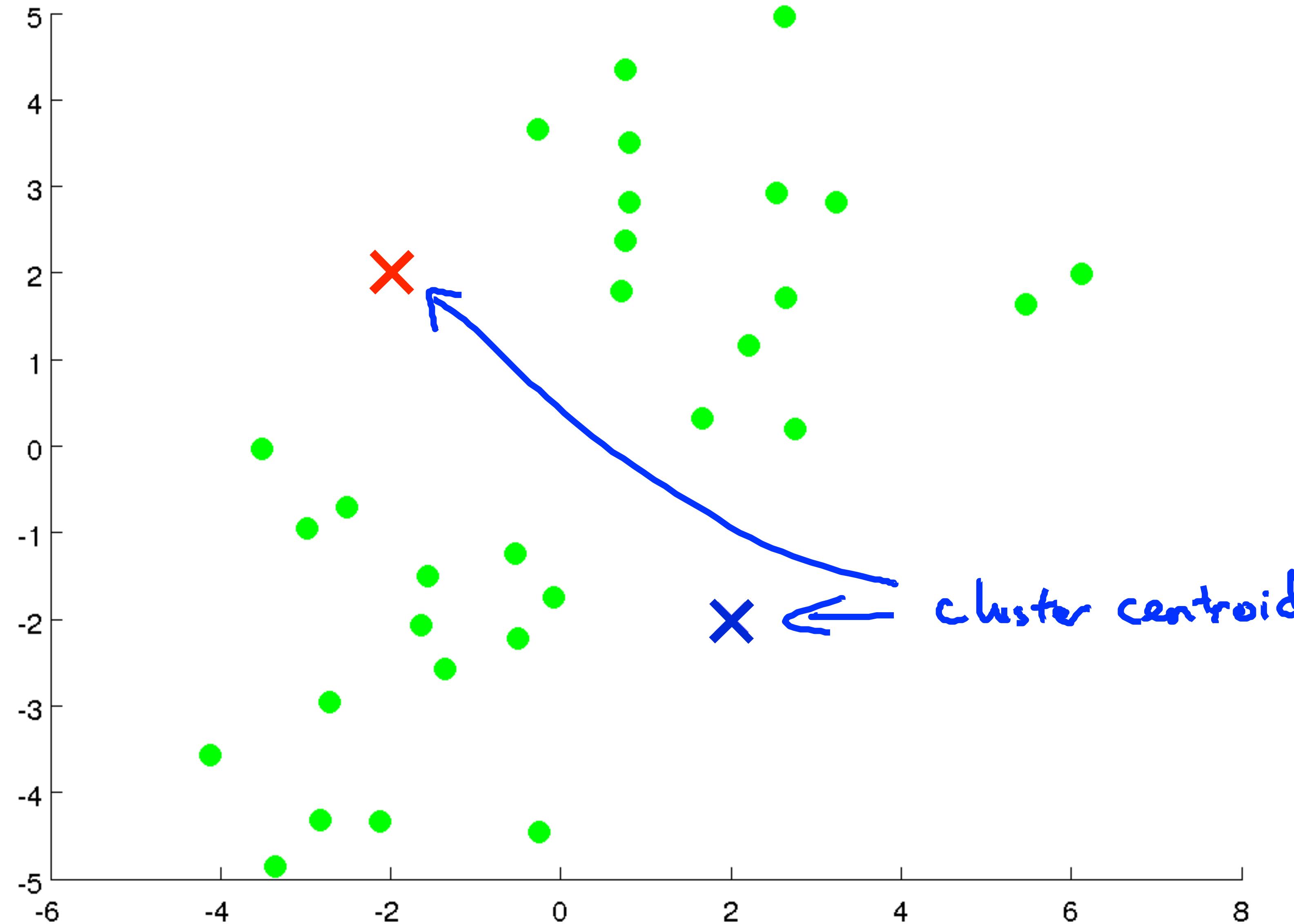
k-means clustering

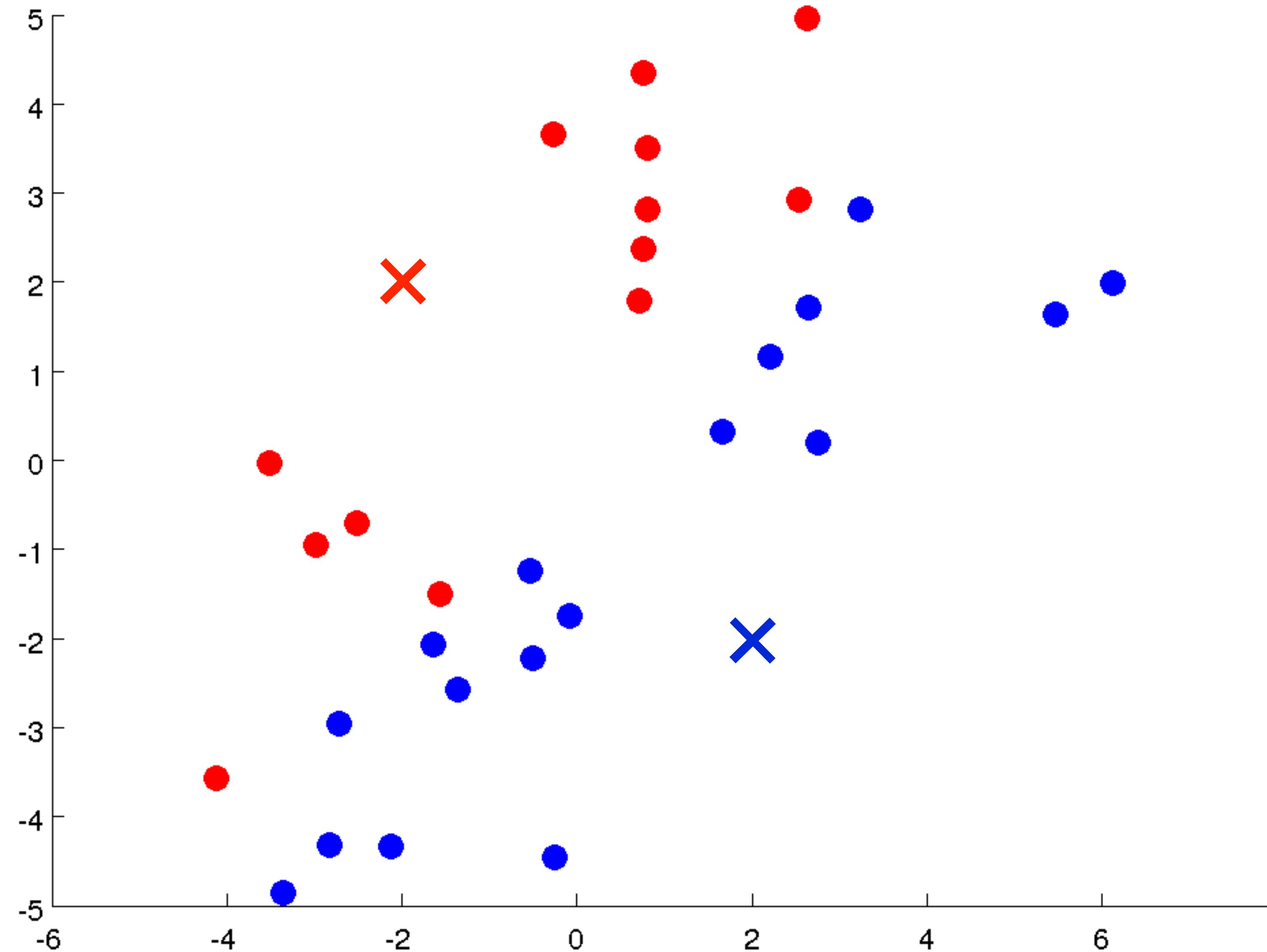
- The goal is to partition the space into K Voronoi cells (clusters)
- Objective function:
$$J = \frac{1}{N} \sum_{i=1}^N \|\mathbf{x}^{(i)} - \boldsymbol{\mu}_i\|^2$$
- Where $\boldsymbol{\mu}_i$ is the mean of points in cluster S_i
- Iterative algorithm:
 - Randomly choose K centroids
 - Compute for each \mathbf{x} the cluster to which it belongs
 - Calculate the average of \mathbf{x} 's in each cluster and update centroid position

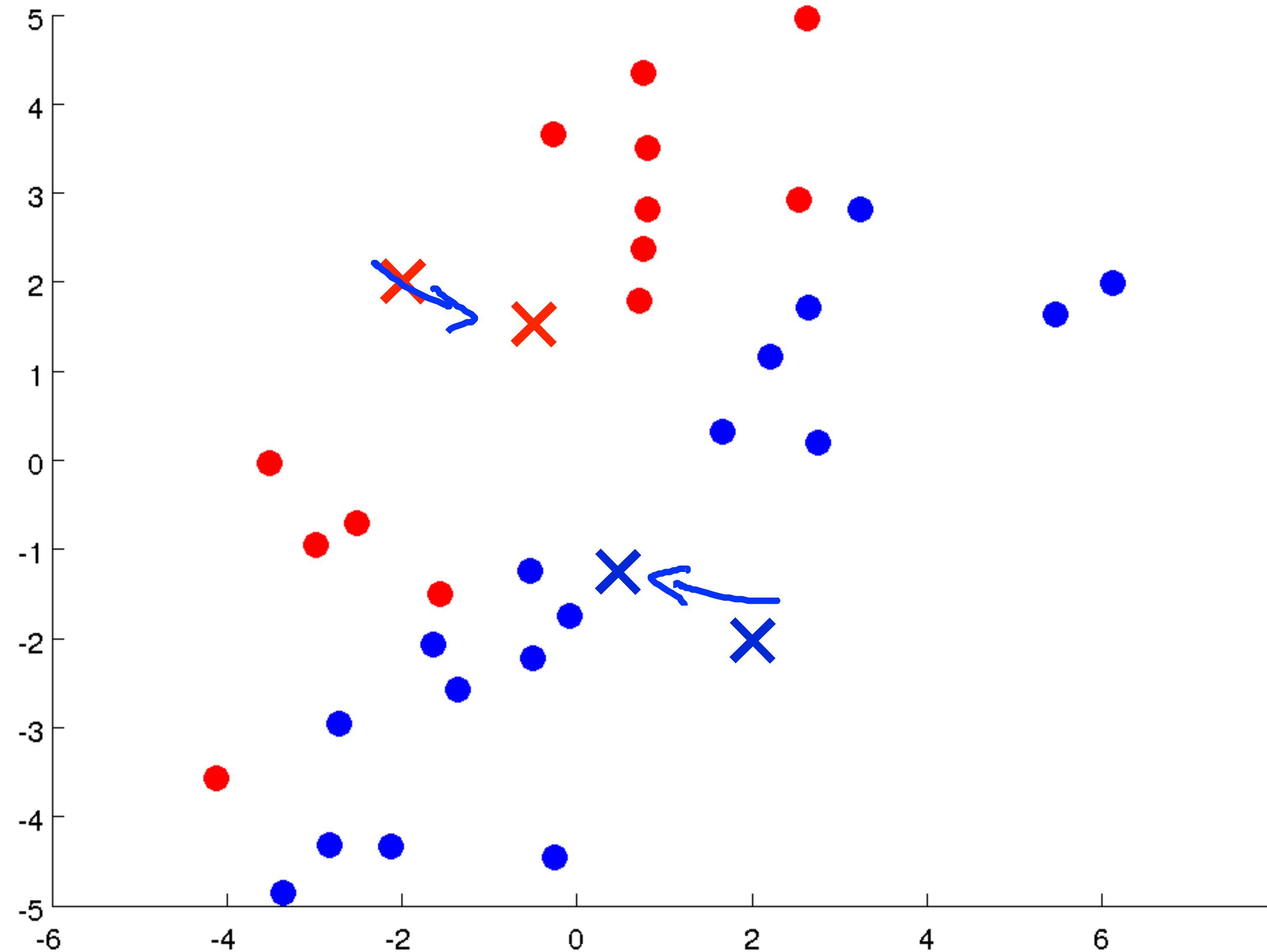
Example of k -means

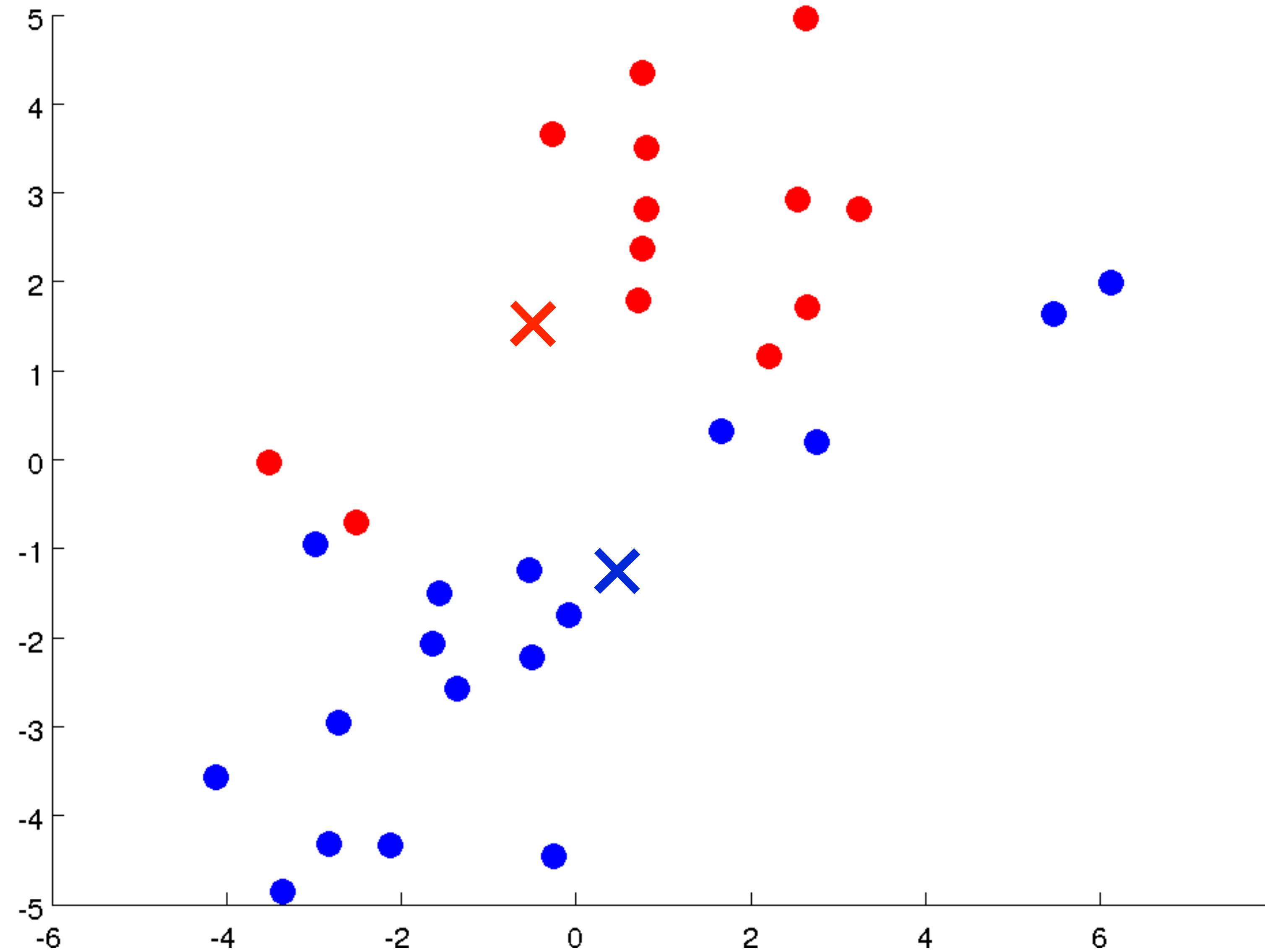
From MOOC Machine Learning, Andrew Ng (Stanford)

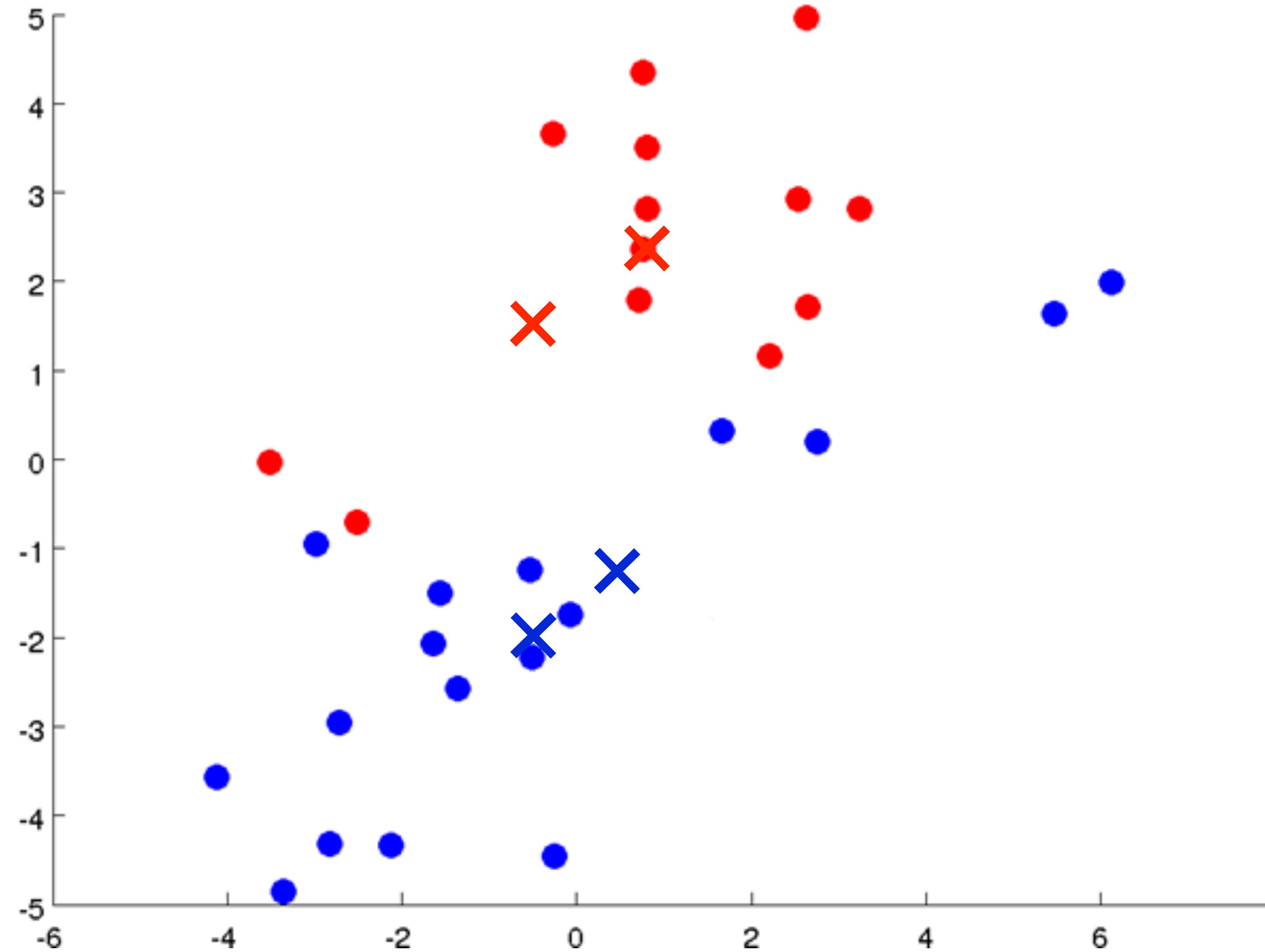


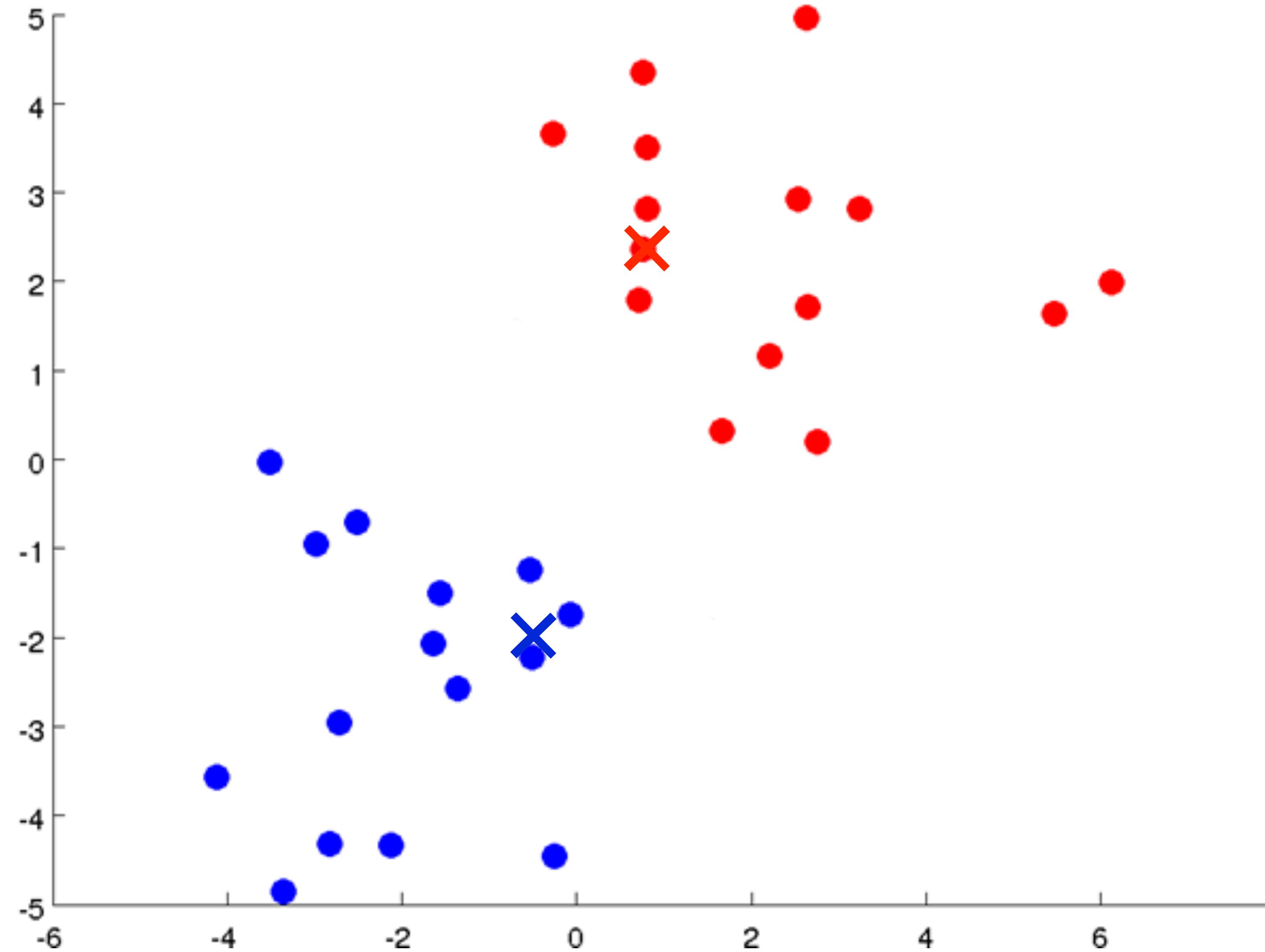


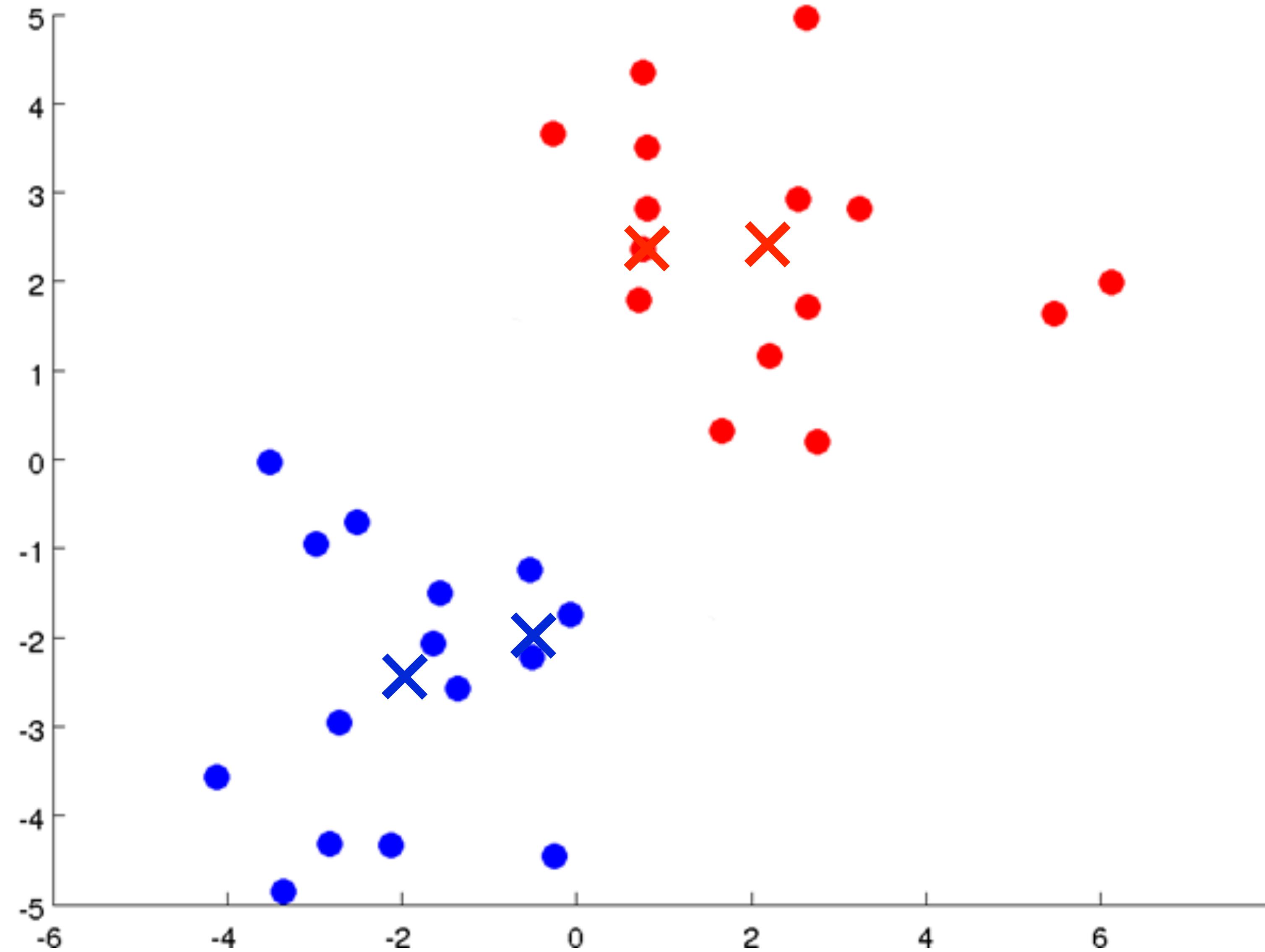


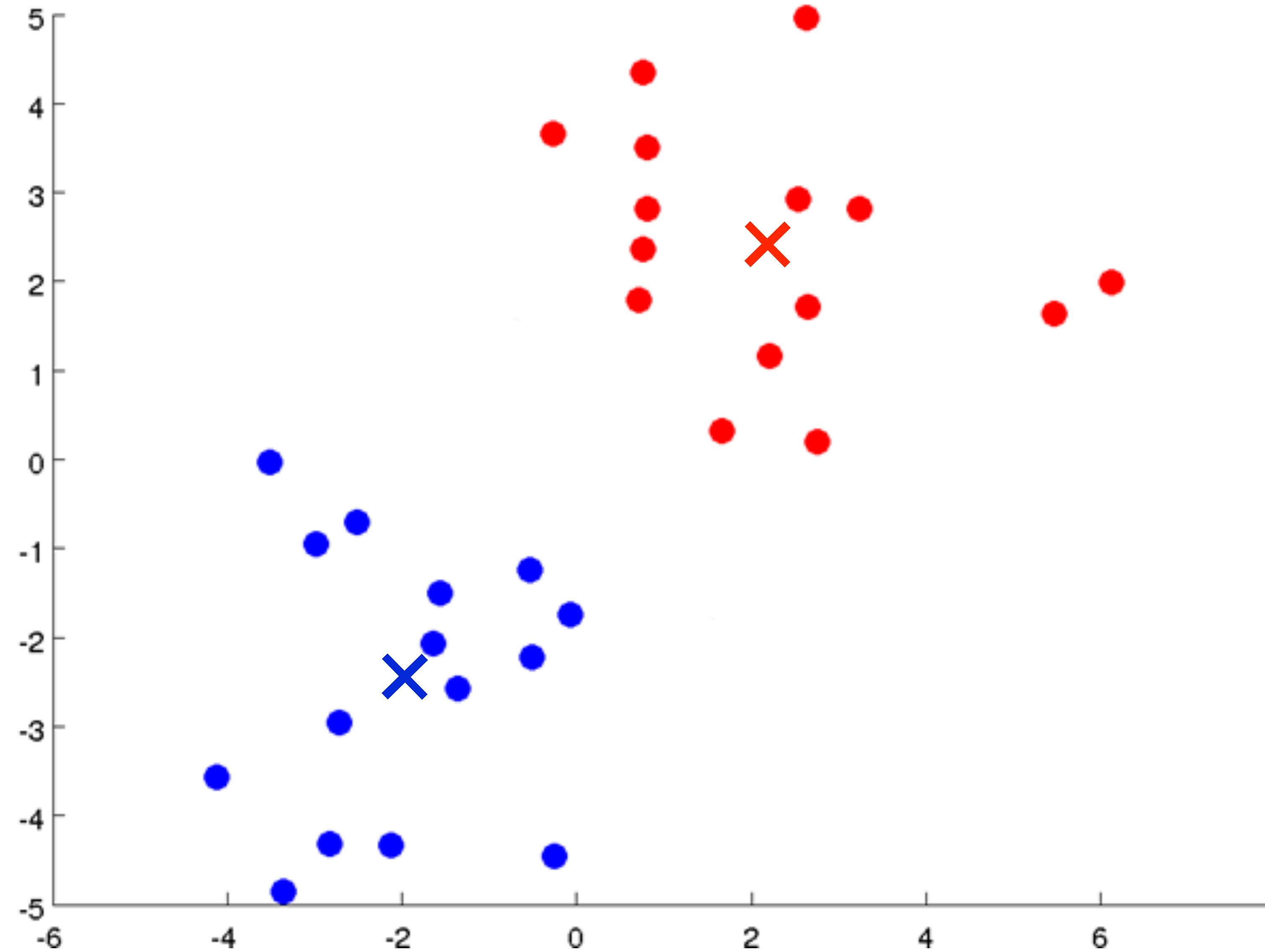








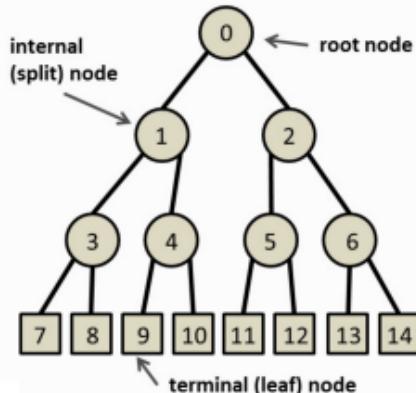




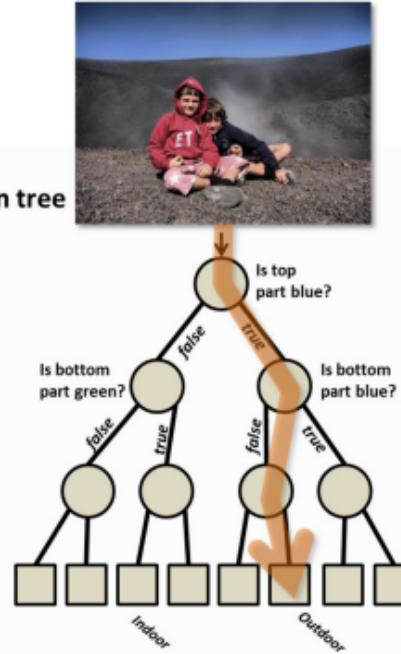
Decision Trees

Image classification example

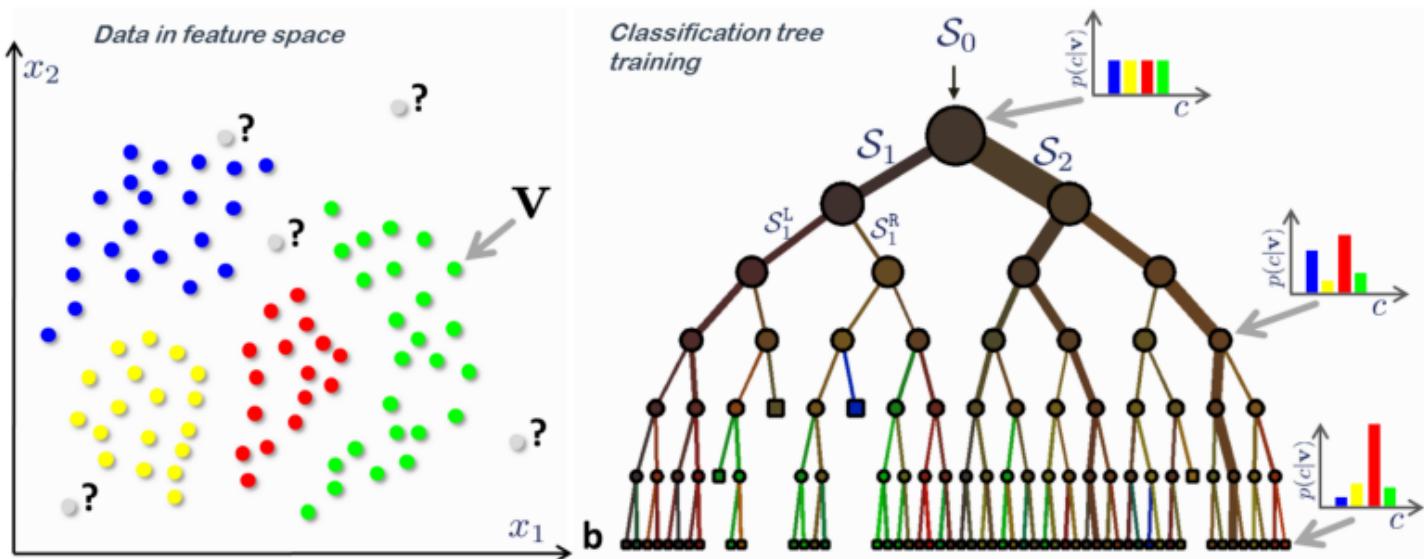
A general tree structure



A decision tree



Another classification tree



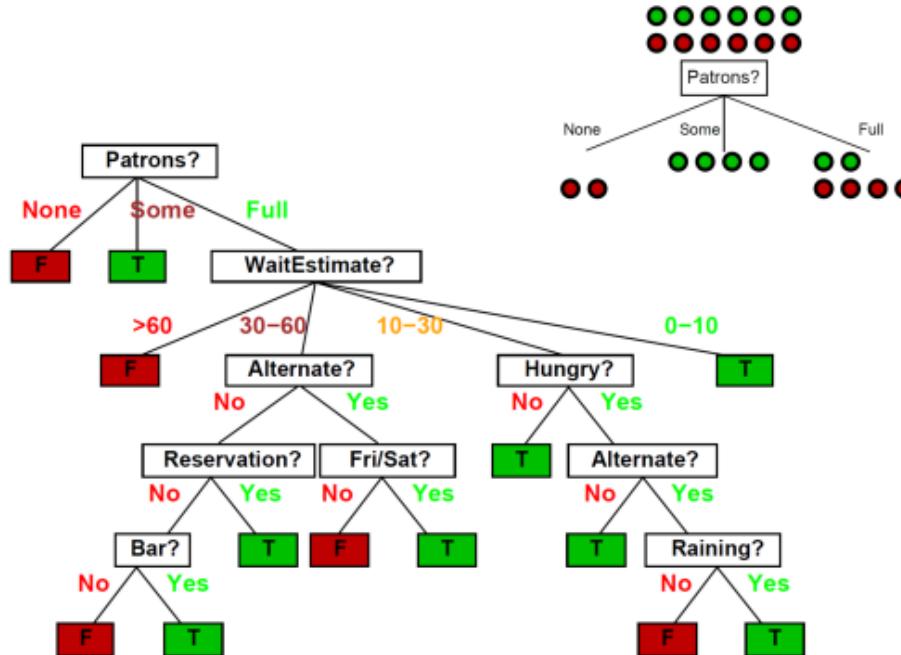
How do we build a tree

Building a node from data

Example	Input Attributes										Goal <i>WillWait</i>
	Alt	Bar	Fri	Hun	Pat	Price	Rain	Res	Type	Est	
x₁	Yes	No	No	Yes	Some	\$\$\$	No	Yes	French	0–10	$y_1 = Yes$
x₂	Yes	No	No	Yes	Full	\$	No	No	Thai	30–60	$y_2 = No$
x₃	No	Yes	No	No	Some	\$	No	No	Burger	0–10	$y_3 = Yes$
x₄	Yes	No	Yes	Yes	Full	\$	Yes	No	Thai	10–30	$y_4 = Yes$
x₅	Yes	No	Yes	No	Full	\$\$\$	No	Yes	French	>60	$y_5 = No$
x₆	No	Yes	No	Yes	Some	\$\$	Yes	Yes	Italian	0–10	$y_6 = Yes$
x₇	No	Yes	No	No	None	\$	Yes	No	Burger	0–10	$y_7 = No$
x₈	No	No	No	Yes	Some	\$\$	Yes	Yes	Thai	0–10	$y_8 = Yes$
x₉	No	Yes	Yes	No	Full	\$	Yes	No	Burger	>60	$y_9 = No$
x₁₀	Yes	Yes	Yes	Yes	Full	\$\$\$	No	Yes	Italian	10–30	$y_{10} = No$
x₁₁	No	No	No	No	None	\$	No	No	Thai	0–10	$y_{11} = No$
x₁₂	Yes	Yes	Yes	Yes	Full	\$	No	No	Burger	30–60	$y_{12} = Yes$

How do we build a tree

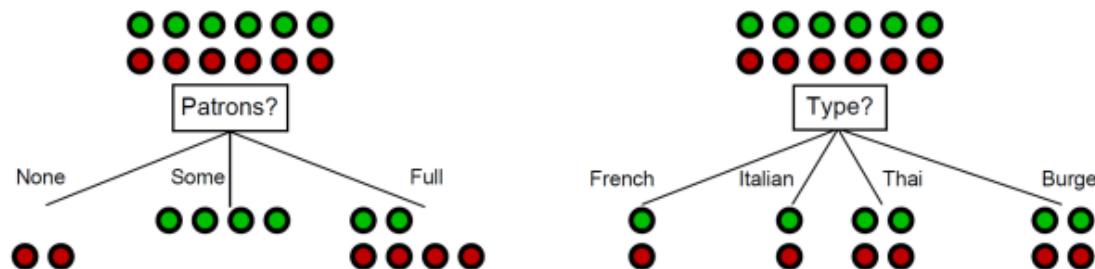
A learned tree



How do we build a tree

Which split is better?

Ideally we want to separate negative and positive examples



How do we build a tree

- Shannon Entropy

$$H = - \sum_i p_i \log_2(p_i)$$

- Expected Entropy (for a feature F with K values)

$$EH(F) = - \sum_{i=1}^K \frac{n_i}{N} H_i$$

- Information Gain $I(F)$

$$I(F) = H(Data) - EH(F)$$

How do we build a tree

The patron vs type example

