

Supervised Learning

Repeating the mistakes of the past

Naive Bayes

$$P(\text{Toss} = \text{heads}) = 0.5.$$

Independence

$$P(\text{Toss} = \text{heads}) = 0.5.$$

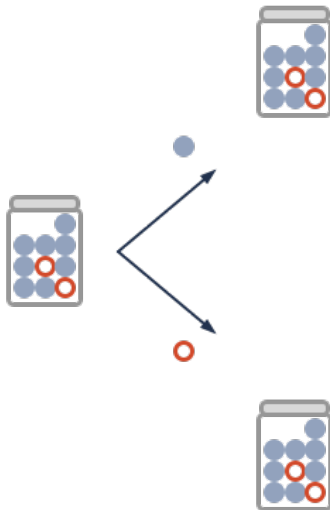
$$P(\text{Toss} = \text{heads} | \text{LastToss} = \text{heads}) = 0.5.$$

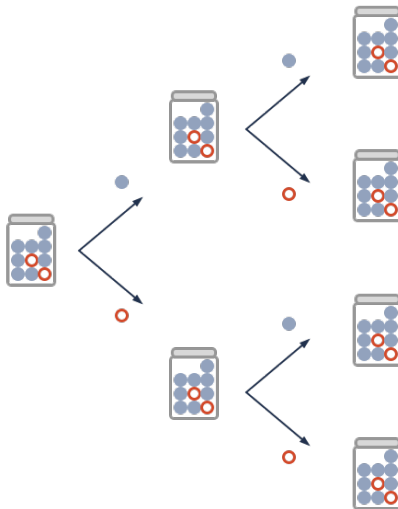
Marbles Problem

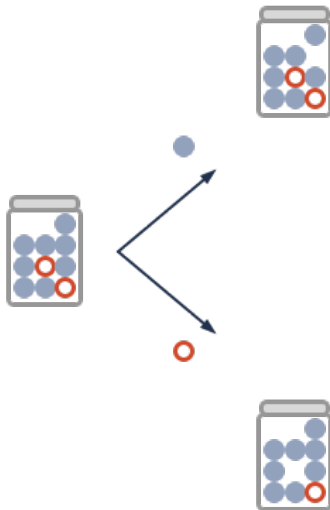
8 Blue, 2 Orange

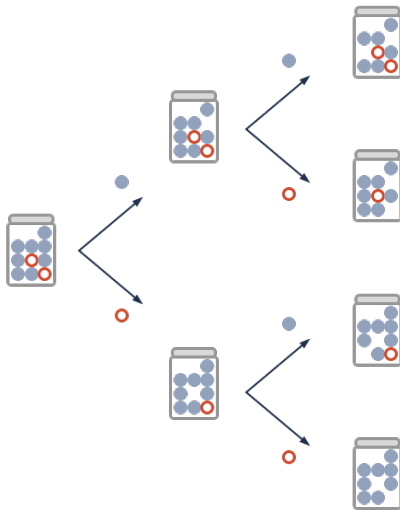
$$P(\text{Select} = \text{blue}) = \frac{8}{8+2} = 0.8$$

$$P(\text{Select} = \text{orange}) = \frac{2}{8+2} = 0.2$$









Marbles Problem

8 Blue, 2 Orange

$P(\text{Select} = \text{blue} \mid \text{Previous} = \text{blue})$

If removing marbles, $\frac{7}{7+2} = 0.78$

If returning marbles, $\frac{8}{8+2} = 0.8$

Bayes' Rule

$$P(A \mid B) = \frac{P(A) P(B|A)}{P(B)}.$$

k classes where each datum has d features

$$P(\text{Class} = i \mid x_1, \dots, x_d) = \frac{P(\text{Class} = i) P(x_1, \dots, x_d \mid \text{Class} = i)}{P(x_1, \dots, x_d)}.$$

Naive because we assume that features are conditionally independent,

$$P(x_a \mid \text{Class} = i, x_1, \dots, x_{a-1}, x_{a+1}, \dots, x_d) = P(x_a \mid \text{Class} = i).$$

$$P(\text{Class} = i \mid x_1, \dots, x_d) = \frac{P(\text{Class} = i) \prod_{\alpha=1}^n P(x_{\alpha} \mid \text{Class} = i)}{P(x_1, \dots, x_d)}.$$

$$\text{Class} = \arg \max_i P(\text{Class} = i) \prod_{\alpha=1}^n P(x_{\alpha} \mid \text{Class} = i)$$

$$P(x_a \mid \text{Class} = i) = \frac{\exp \frac{(x_a - \mu_{a,i})^2}{2\sigma_{a,i}^2}}{\sqrt{2\pi\sigma_{a,i}^2}},$$

Gaussian

where $\sigma_{a,i}$ and $\mu_{a,i}$ are the standard deviation and mean of feature a for class i .

$$P(x_a | \text{Class} = i) = \frac{\sum_{\mathbf{x} \in I} x_a + \alpha}{\sum_{a=1}^n \sum_{\mathbf{x} \in I} x_a + \alpha n},$$

Multinomial

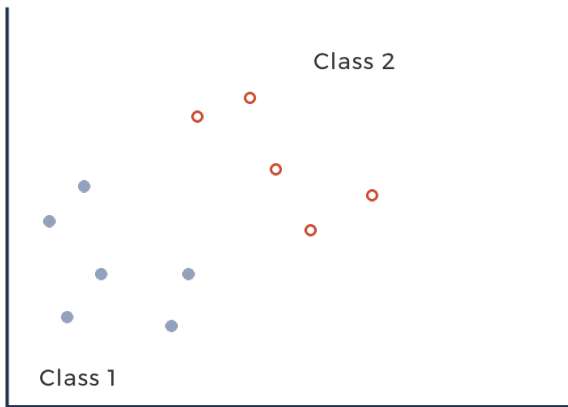
where I is the set of feature vectors of class i , and α is a smoothing parameter.

$$P(x_a \mid \text{Class} = i) = \frac{\sum_{\mathbf{x} \in I} x_a}{m_i} x_a + (1 - \frac{\sum_{\mathbf{x} \in I} x_a}{m_i})(1 - x_a),$$

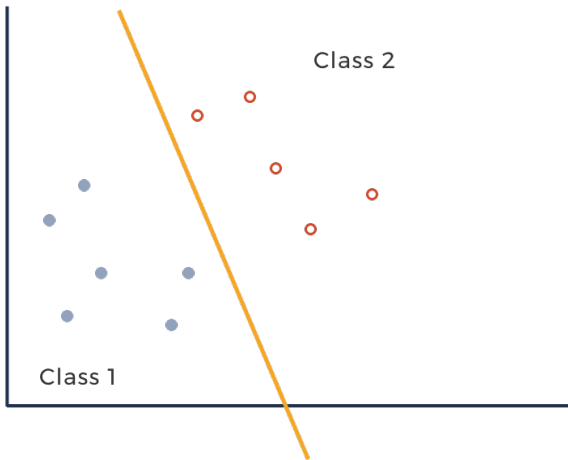
Bernouli

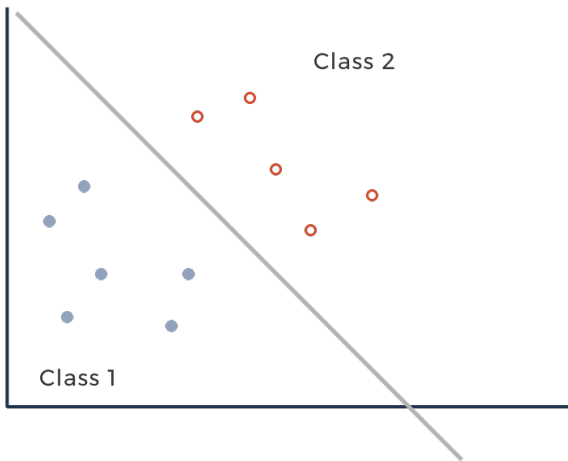
where all n features are boolean values, I is the set of feature vectors of class i , and m_i is the number of vectors in class i .

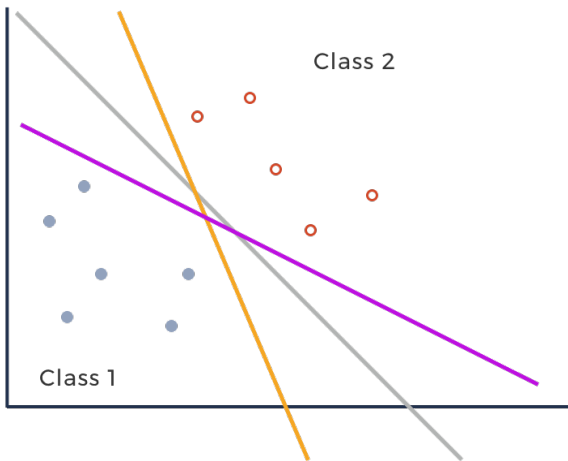
Support Vector Machines

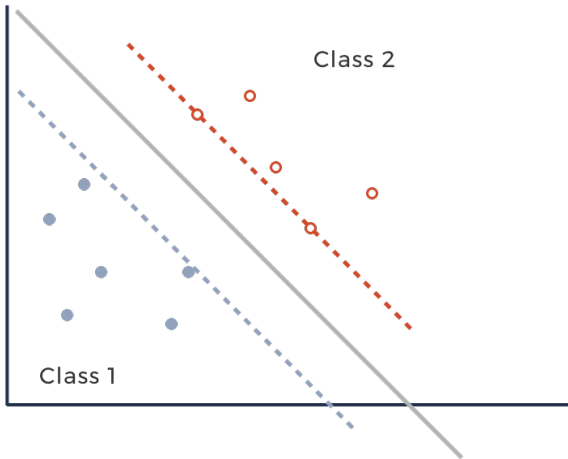


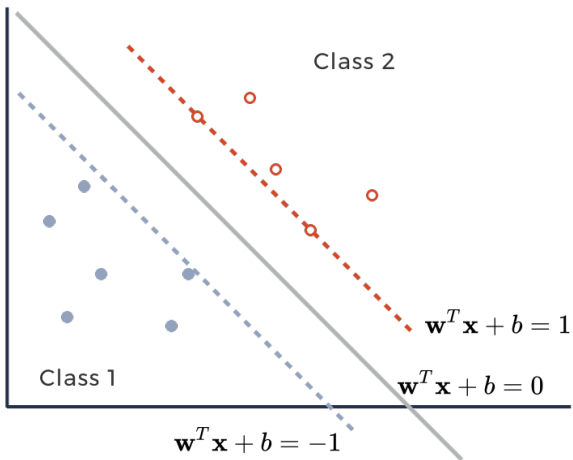


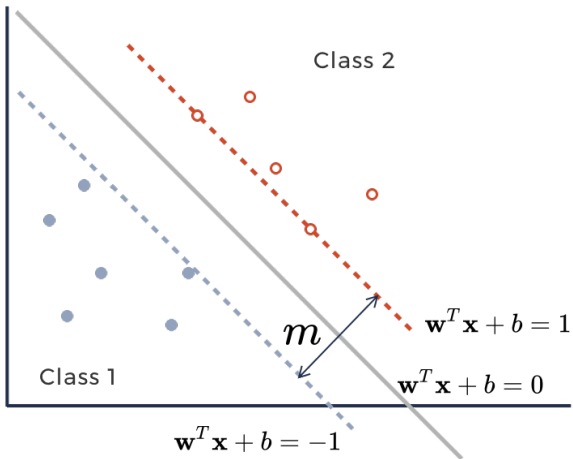




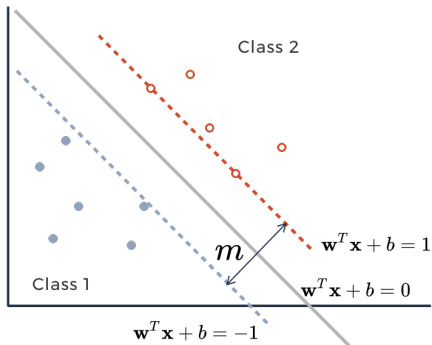








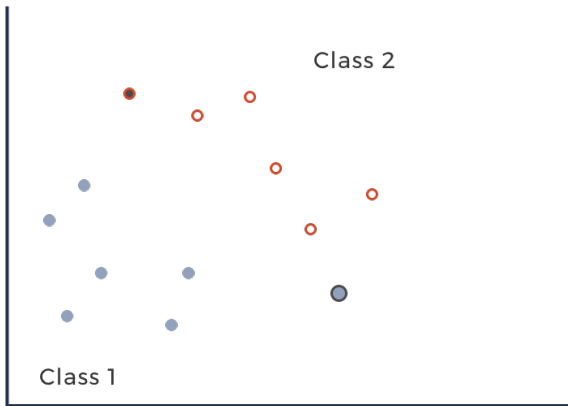
Margin

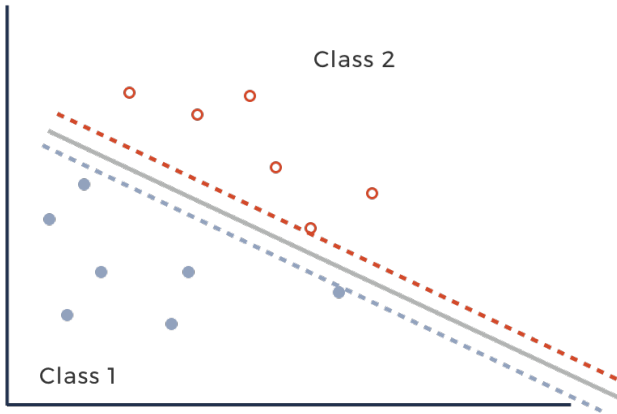


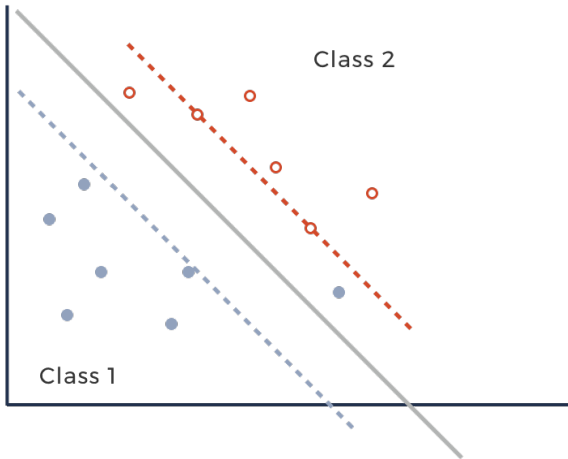
$$\text{margin} = \frac{2}{\|\mathbf{w}\|}.$$

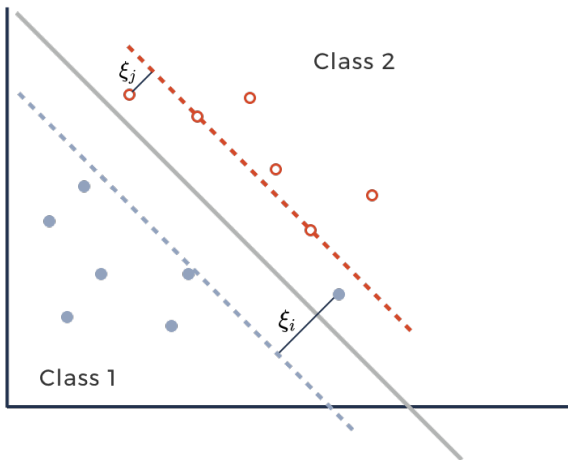
$$\begin{aligned} \min_{\mathbf{w}} \quad & \frac{1}{2} \|\mathbf{w}\|^2 \\ \text{s.t.} \quad & y_i(\mathbf{w}^T \mathbf{x}_i + b) \geq 1 \quad \forall i, \end{aligned}$$

where $y_i \in [1, -1]$ determines the class point i belongs to.







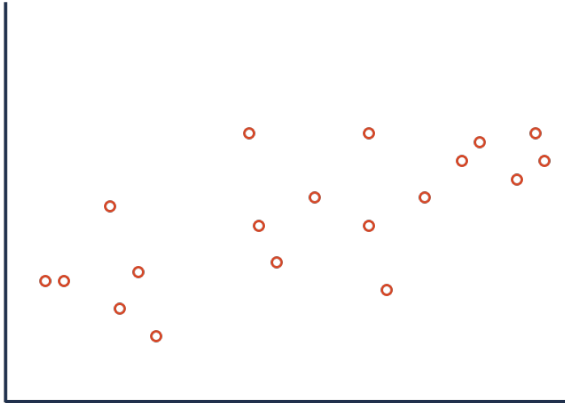


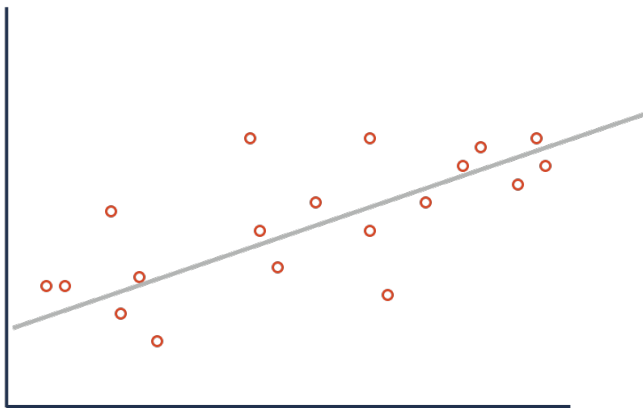
$$\begin{aligned} \min_{\mathbf{w}} \quad & \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^n \xi_i \\ \text{s.t.} \quad & y_i(\mathbf{w}^T \mathbf{x}_i + b) \geq 1 - \xi_i, \xi_i \geq 0 \quad \forall i, \end{aligned}$$

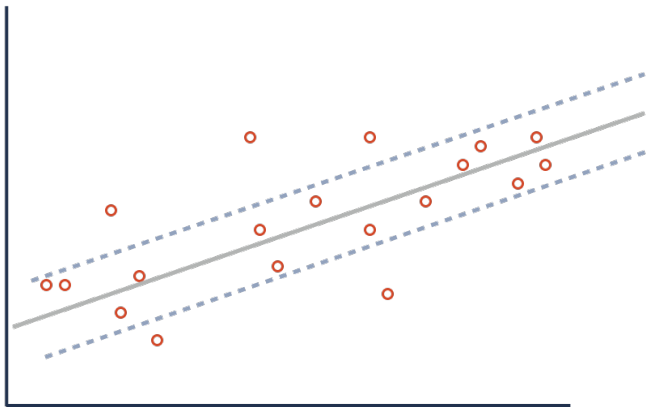
where ξ_i is the error for point i , and C is the penalty tuning parameter.

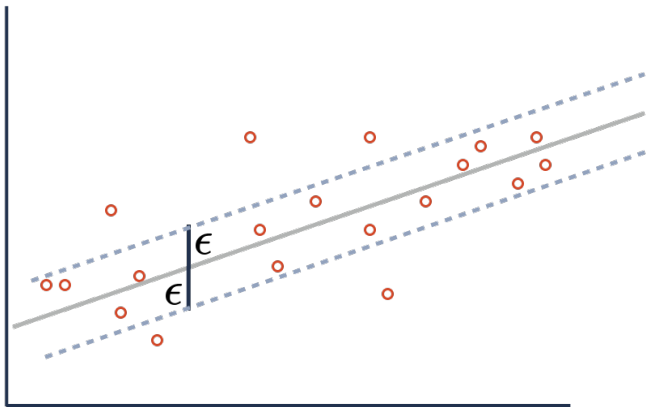
Multi-Class

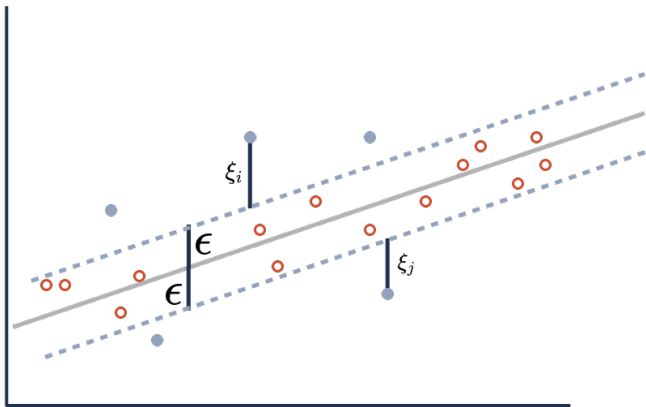
- One vs All / Rest
- One vs One











Support Vector Regression

$$\begin{aligned} \min_{\mathbf{w}} \quad & \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^n \xi_i \\ \text{s.t.} \quad & |y_i - (\langle \mathbf{w}, \mathbf{x}_i \rangle + b)| \leq \epsilon + \xi_i \quad , \epsilon, \xi_i \geq 0 \quad \forall i, \end{aligned}$$

where ϵ is the acceptable error boundary.

Kernel Method

Transform the data into linearly separable features.

$\phi(\mathbf{x})$ defines a transformation that occurs to any given point into the new space

$$K(\mathbf{x}_i, \mathbf{x}_j) = \langle \phi(\mathbf{x}_i), \phi(\mathbf{x}_j) \rangle = \phi(\mathbf{x}_i)^T \phi(\mathbf{x}_j).$$

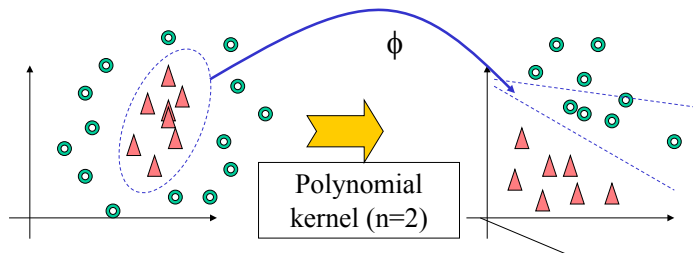
Polynomial Kernel

$$K(\mathbf{x}_i, \mathbf{x}_j) = (\mathbf{x}_i^T \mathbf{x}_j + 1)^d,$$

Polynomial

with degree d ,

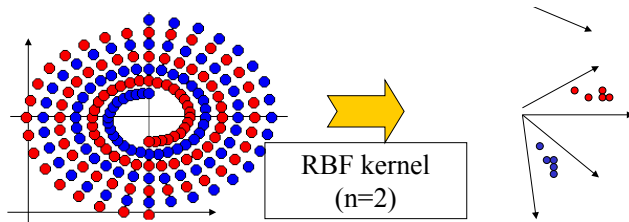
Polynomial Kernel



$$K_{\text{RBF}}(\mathbf{x}_i, \mathbf{x}_j) = \exp\left(\frac{-\|\mathbf{x}_i - \mathbf{x}_j\|^2}{2\sigma^2}\right),$$

with width σ ,

RBF Kernel



Sigmoid Kernel

$$K(\mathbf{x}_i, \mathbf{x}_j) = \tanh(\kappa \mathbf{x}_i^T \mathbf{x}_j + \theta),$$

Sigmoid

with parameters κ and θ .

Questions

These slides are designed for educational purposes, specifically the CSCI-470 Introduction to Machine Learning course at the Colorado School of Mines as part of the Department of Computer Science.

Some content in these slides are obtained from external sources and may be copyright sensitive. Copyright and all rights therein are retained by the respective authors or by other copyright holders. Distributing or reposting the whole or part of these slides not for academic use is HIGHLY prohibited, unless explicit permission from all copyright holders is granted.