

Statistical Machine Learning

5주차

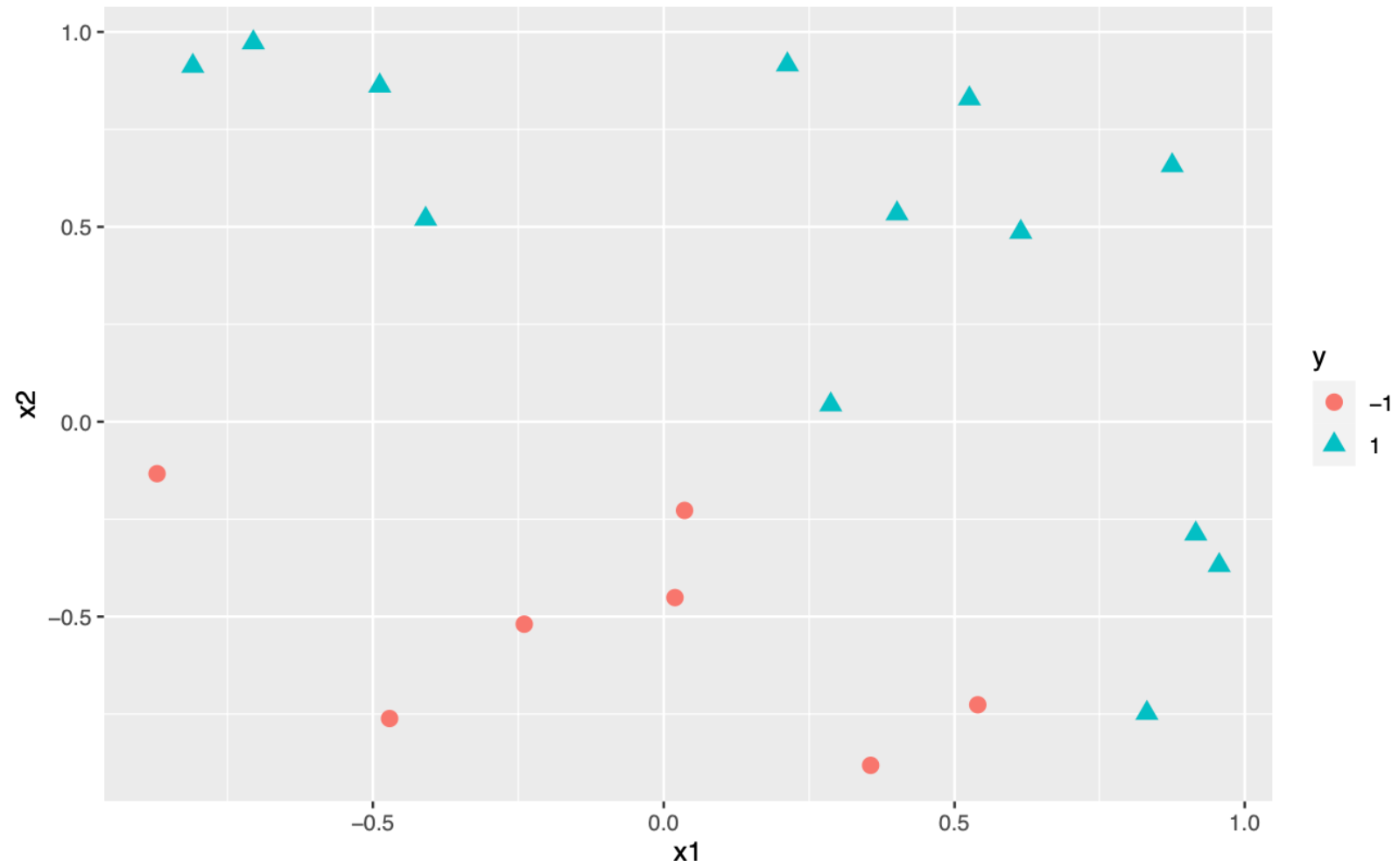
담당: 15기 김지호

- 1. Linear SVM**
- 2. Kernel SVM**
- 3. Decision Tree**

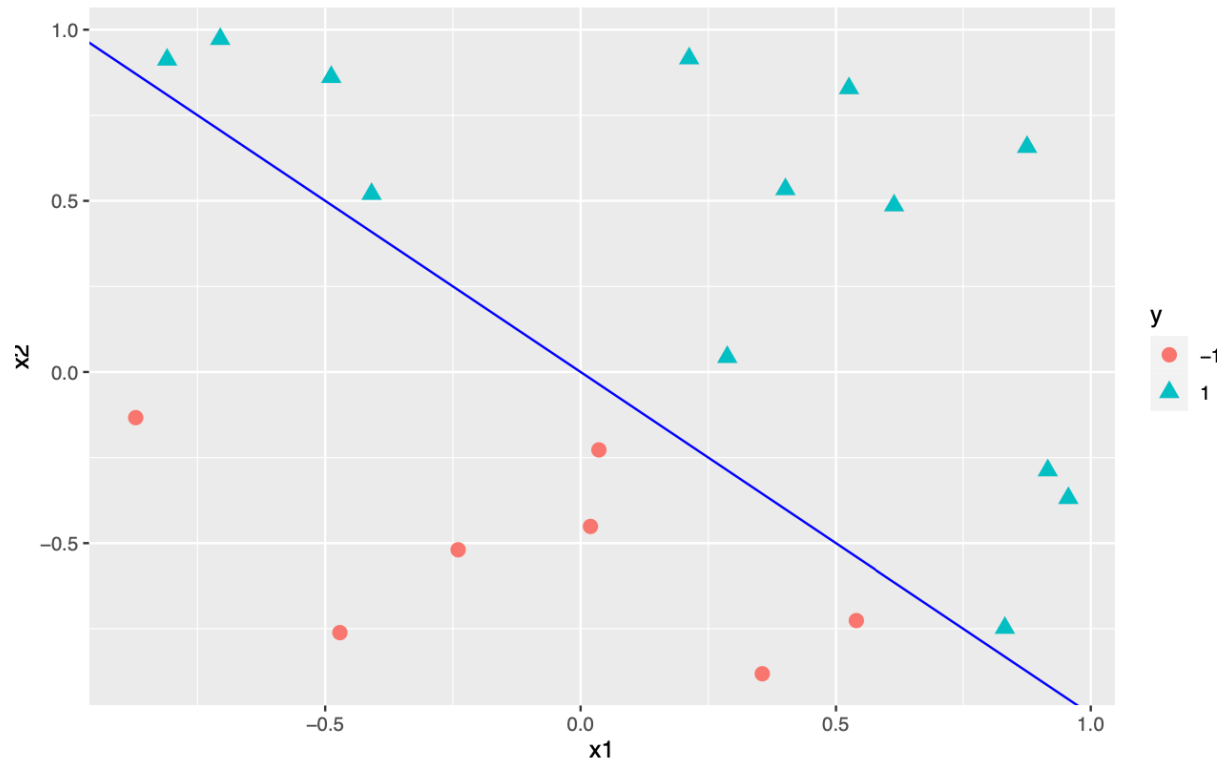
Support Vector Machine

- Classification

- Consider a simple (linearly separable) example.

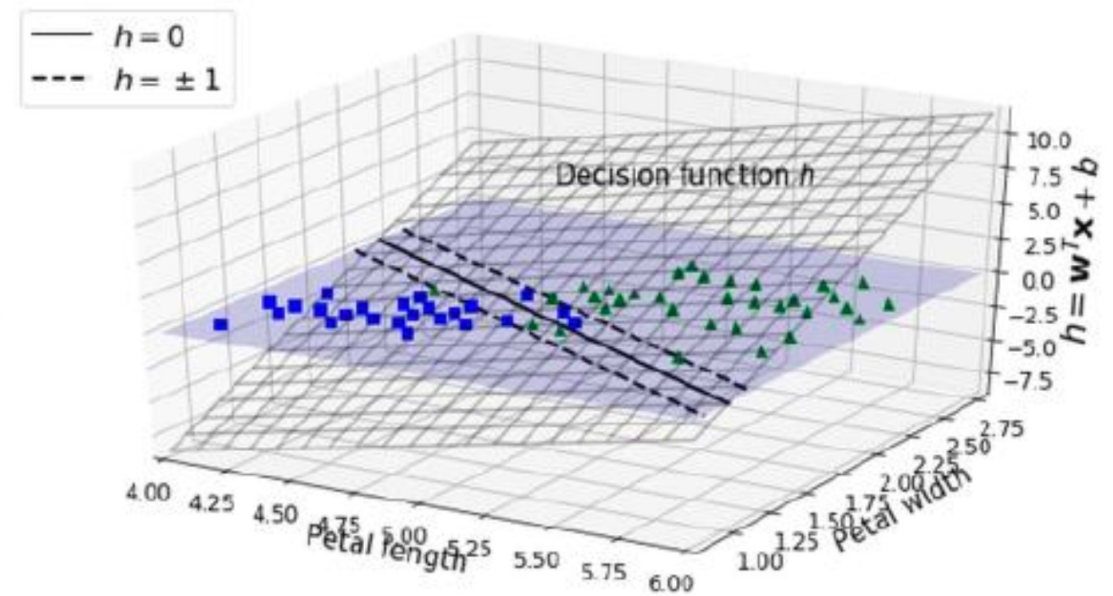
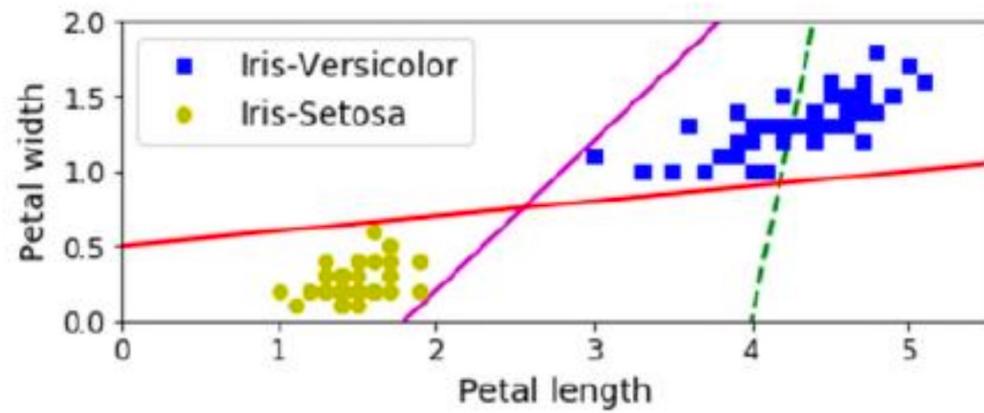


- The line is called the **classification/decision boundary** or **separating hyperplane**.

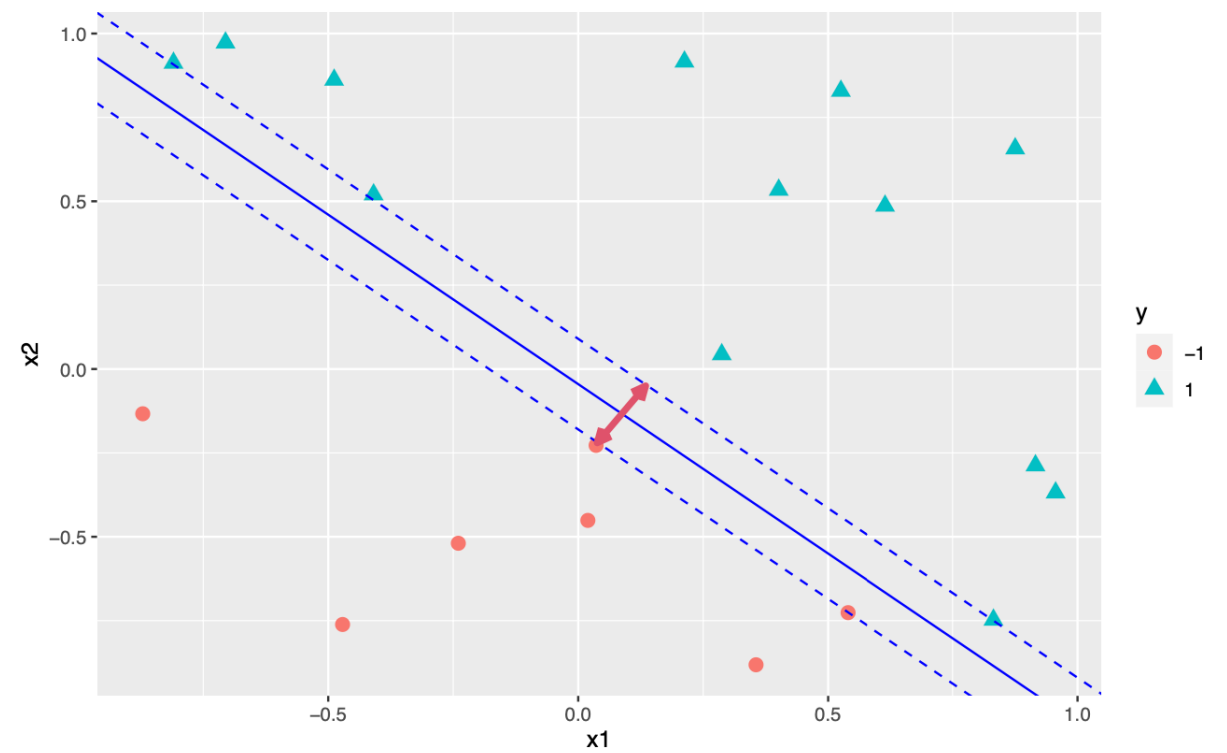
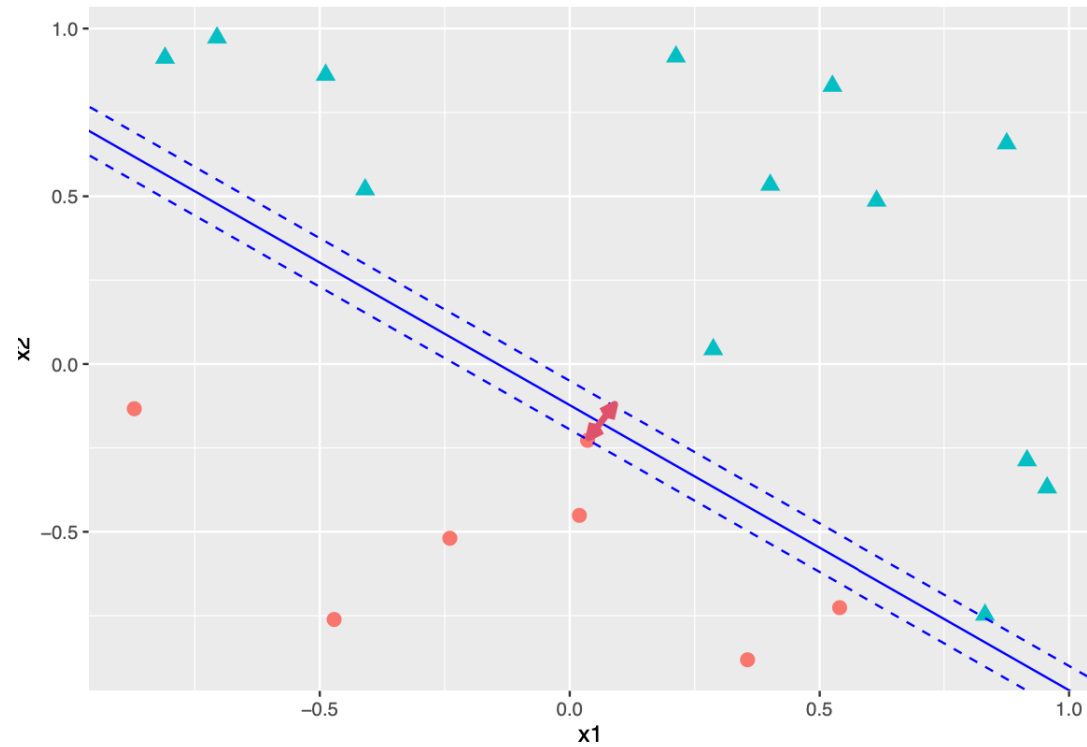


- Decision boundary
$$f(\mathbf{x}) = \beta_0 + \beta^T \mathbf{x} = 0$$
- $y \in \{-1, 1\}$
- Prediction of y given \mathbf{x}
$$\hat{y} = \text{sign}\{f(\mathbf{x})\} = \text{sign}\{\beta_0 + \beta^T \mathbf{x}\}$$

Hyperplane

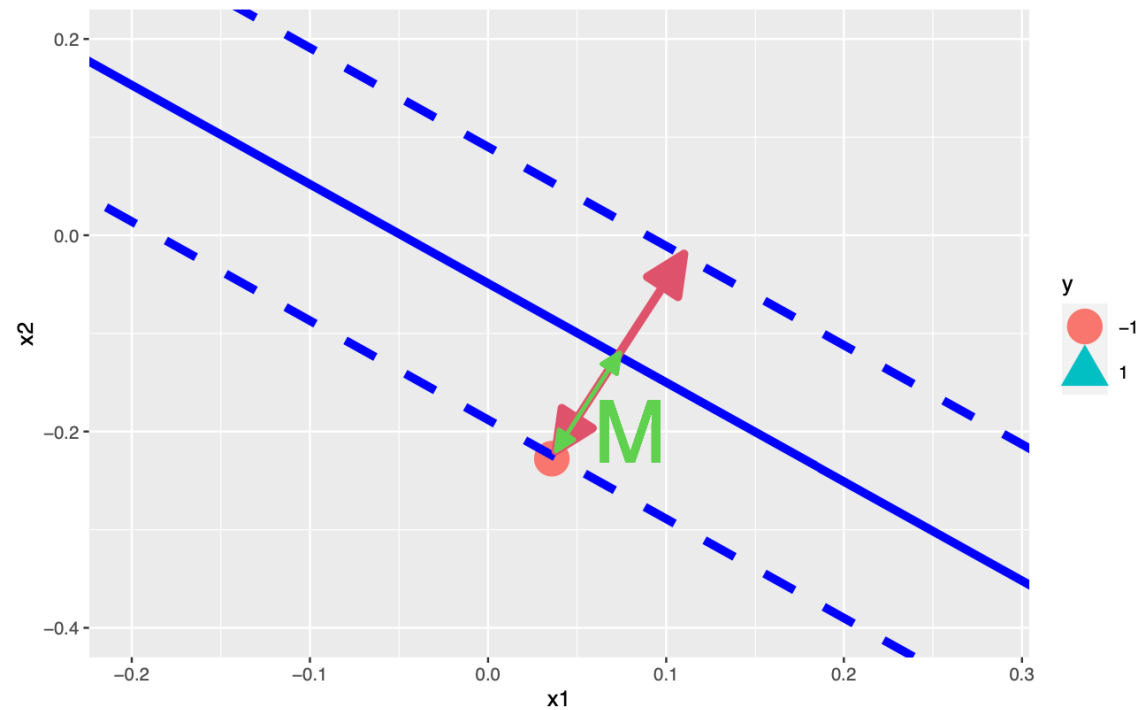


Optimal Separating Hyperplane

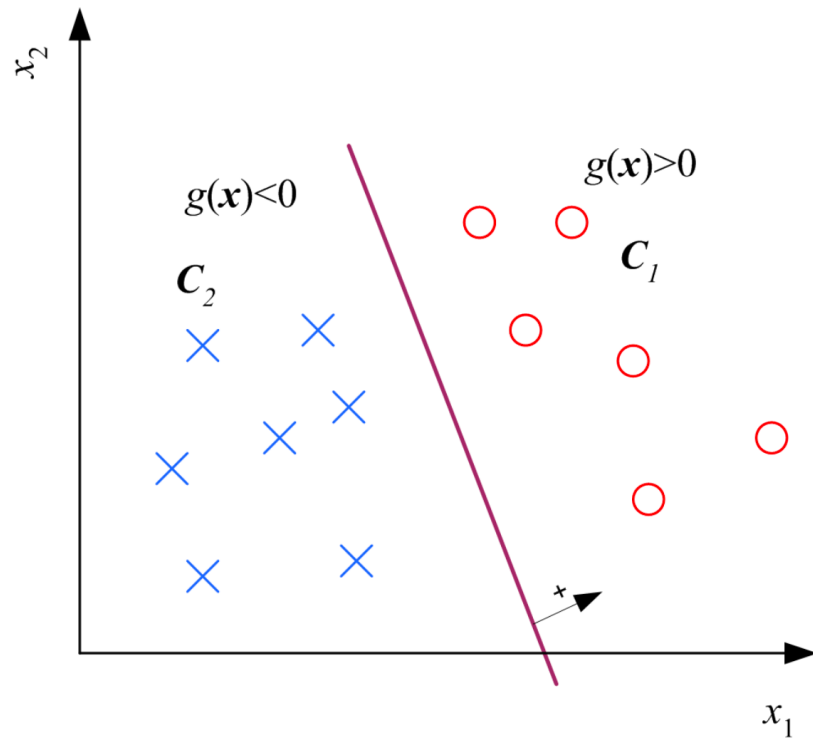


Optimal Separating Hyperplane

- Optimal Separating Hyperplane maximizes **Geometric Margin**, M .



Two Classes

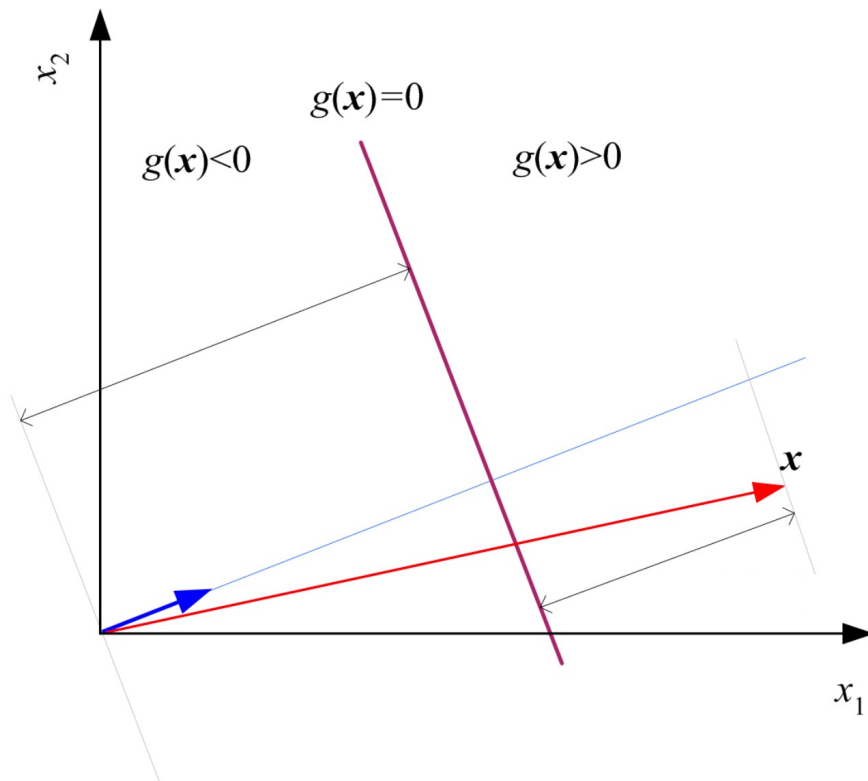


$$g(\mathbf{x}) = g_1(\mathbf{x}) - g_2(\mathbf{x})$$

choose $\begin{cases} C_1 & \text{if } g(\mathbf{x}) > 0 \\ C_2 & \text{otherwise} \end{cases}$

Geometric Margin

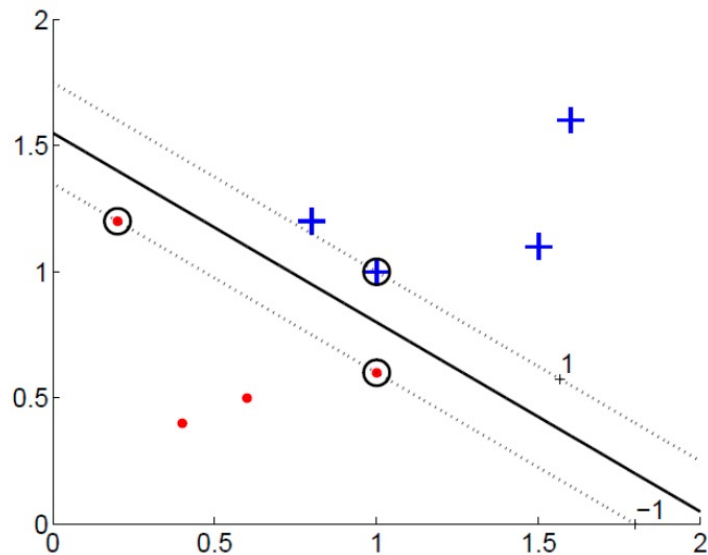
- Geometric margin of x^* $= \frac{|\beta^T x + \beta_0|}{\|\beta\|}$



Geometric Margin

Assume $\|\beta\| = 1$, the geometric margin M of \mathbf{x}_i to the hyperplane $\beta_0 + \beta^T \mathbf{x}$ is

$$M = y_i(\beta_0 + \beta^T \mathbf{x}_i)$$



if \mathbf{x}_i is on right (i.e. $\beta_0 + \beta^T \mathbf{x}_i > 0$ and $y_i = 1$) $M = \beta_0 + \beta^T \mathbf{x}_i$

if \mathbf{x}_i is on left (i.e. $\beta_0 + \beta^T \mathbf{x}_i < 0$ and $y_i = 1$) $M = -(\beta_0 + \beta^T \mathbf{x}_i)$

Maximal Margin Classifier

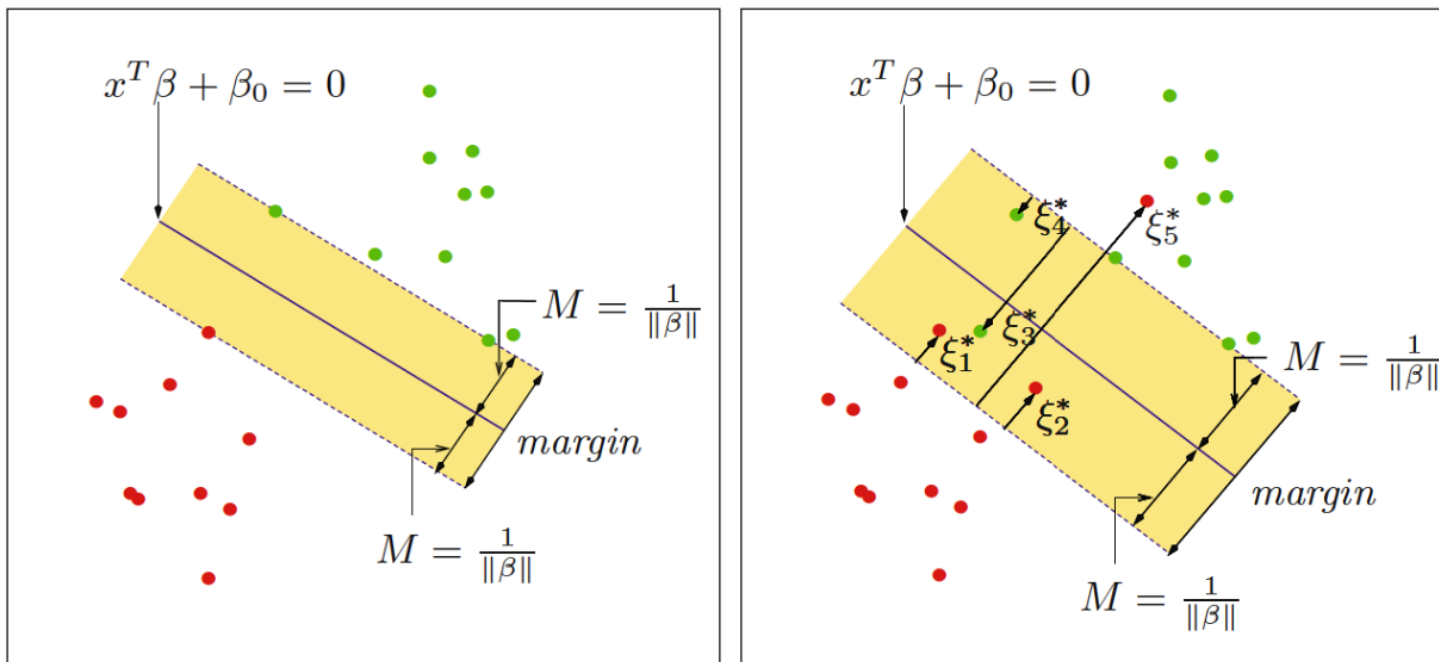
$$\max_{\beta_0, \beta} M \quad \text{subject to} \quad y_i(\beta_0 + \beta^T \mathbf{x}_i) \geq M, i = 1, \dots, n;$$

- For a unique solution, fix $M \parallel \beta \parallel = 1$.
- Maximize $M \Leftrightarrow$ minimize $\parallel \beta \parallel$
 \Leftrightarrow minimize $\parallel \beta \parallel^2$

$$\min_{\beta_0, \beta} \frac{1}{2} \beta^T \beta \quad \text{subject to} \quad y_i(\beta_0 + \beta^T \mathbf{x}_i) \geq M, i = 1, \dots, n;$$

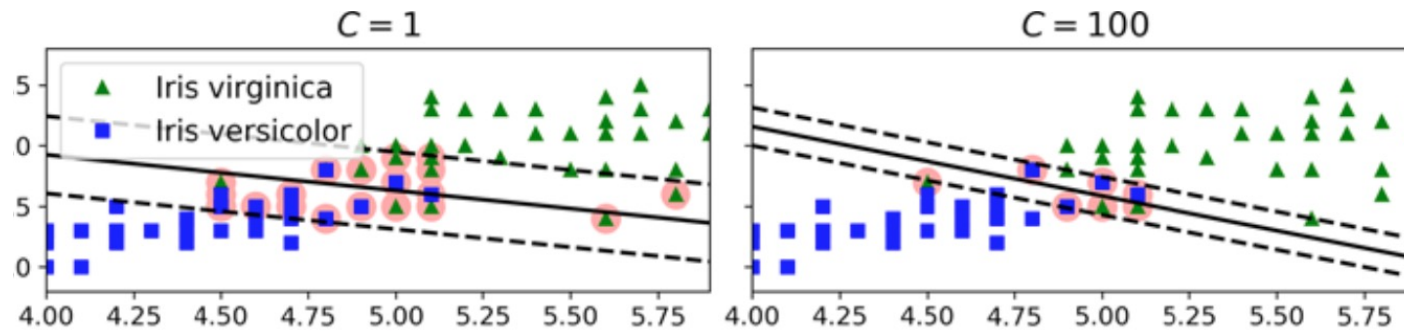
What if Nonseparable?

- Let's relax the constraints by introducing slack variables $\xi_i \geq 0$, and add penalty C for the violations.



Support Vector Machine – Soft Margin Classifier

$$\min_{\beta_0, \beta, \xi_i} \beta^T \beta + C \sum_{i=1}^n \xi_i \quad \text{subject to} \quad y_i(\beta_0 + \beta^T \mathbf{x}_i) \geq 1 - \xi_i, \quad i = 1, \dots, n$$
$$\xi_i \geq 0, \quad i = 1, \dots, n.$$



Computation of SVM

Lagrangian Method

- Constraint Optimization for \mathbf{x} :
$$\begin{aligned} \min_{\mathbf{x}} \quad & f(\mathbf{x}) \\ \text{subject to} \quad & h_i(\mathbf{x}) \leq 0, i = 1, \dots, n \end{aligned} \tag{1}$$
- The **Lagrangian** associated the problem (1) is
$$f(\mathbf{x}) + \sum_{i=1}^n \alpha_i h_i(\mathbf{x})$$
- where $\alpha = (\alpha_1, \dots, \alpha_m)$ are called the **dual variables** or **Lagrange multipliers** associated with the problem (1).

Computation of SVM

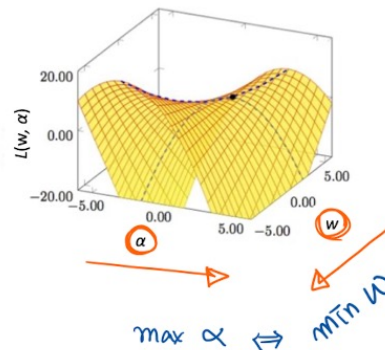
- Lagrangian function of the linear SVM is

$$L_p : \frac{1}{2} \beta^T \beta + C \sum_{i=1}^n \xi_i + \sum_{i=1}^n \alpha_i \{1 - y_i(\beta_0 - \beta^T \mathbf{x}_i) - \xi_i\} - \sum_{i=1}^n \gamma_i \xi_i \quad (2)$$

where α_i and γ_i are (non-negative) Lagrangian multiplier.

- KKT Stationary conditions

Karush-Kuhn-Tucker theorem:
If (w^*, α^*) is a saddle point of $L(w, \alpha)$ in $\alpha \geq 0$, then w^* is an optimal vector.



$$\frac{\partial}{\partial \beta} L_p : \quad \beta = \sum_{i=1}^n \alpha_i y_i \mathbf{x}_i$$

$$\frac{\partial}{\partial \beta_0} L_p : \quad \sum_{i=1}^n \alpha_i y_i = 0$$

$$\frac{\partial}{\partial \xi_i} L_p : \quad \alpha_i = C - \gamma_i$$

Dual Problem of SVM

- Dual function for the linear SVM :
$$g(\alpha) = \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y_i y_j \mathbf{x}_i^T \mathbf{x}_j$$

- Dual problem of the linear SVM :

$$\max_{\alpha} g(\alpha)$$

subject to $0 \leq \alpha_i \leq C, \quad i = 1, \dots, n;$

$$\sum_{i=1}^n \alpha_i y_i = 0.$$

- In Matrix notation :

$$\max_{\alpha} \mathbf{1}^T \alpha - \frac{1}{2} \alpha^T \mathbf{K}^* \alpha$$

subject to $\mathbf{1} \leq \alpha \leq C\mathbf{1},$

$$\mathbf{y}^T \alpha = 0$$

Where $\{\mathbf{K}^*\}_{ij} = y_i y_j \mathbf{x}_i^T \mathbf{x}_j = y_i y_j \langle \mathbf{x}_i, \mathbf{x}_j \rangle$

Computation of SVM

- Primal problem of the linear SVM : $\min_{\beta_0, \beta} \frac{1}{2} \beta^T \beta$ subject to $y_i(\beta_0 + \beta^T \mathbf{x}_i) \geq M, i = 1, \dots, n;$
- Dual problem of the linear SVM : $\max_{\alpha} g(\alpha)$ subject to $0 \leq \alpha_i \leq C, \quad i = 1, \dots, n;$
$$\sum_{i=1}^n \alpha_i y_i = 0.$$
- SVM Solution : $f(\mathbf{x}) = \beta_0 + \sum_{i=1}^n \alpha_i y_i \mathbf{x}_i^T \mathbf{x} = \beta_0 + \sum_{i=1}^n \alpha_i y_i \langle \mathbf{x}_i, \mathbf{x} \rangle$

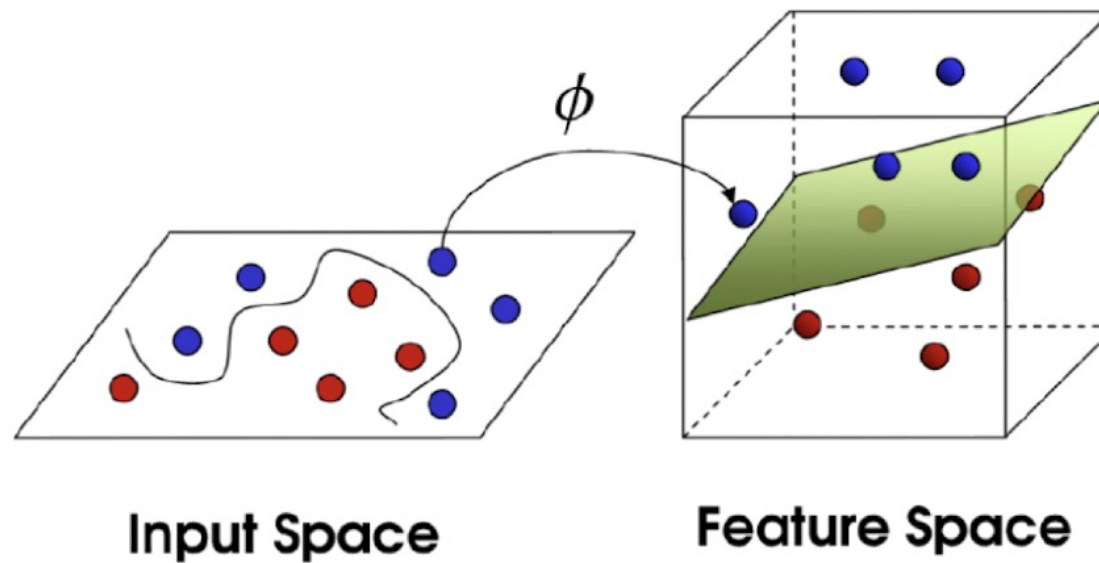
Computation of SVM

- Decision boundary : $f(\mathbf{x}) = \beta_0 + \sum_{i=1}^n \alpha_i y_i \mathbf{x}_i^T \mathbf{x} = \beta_0 + \sum_{i=1}^n \alpha_i y_i \langle \mathbf{x}_i, \mathbf{x} \rangle$
- SVM solution depends on \mathbf{x}_i s only through their inner products.

Extension to Nonlinear Classification

- Using Kernel Trick

Kernel Trick



$$\phi(\mathbf{x}) = (\phi_1(\mathbf{x}), \dots, \phi_n(\mathbf{x}))$$

Kernel Trick

- decision function on space of \mathbf{x} is

$$f(\mathbf{x}) = \beta_0 + \sum_{i=1}^n y_i \alpha_i \langle \mathbf{x}_i, \mathbf{x} \rangle$$

- decision function on feature space is

$$f(\mathbf{x}) = \beta_0 + \sum_{i=1}^n y_i \alpha_i \langle \phi(\mathbf{x}_i), \phi(\mathbf{x}) \rangle$$

- Thus what we need is not the feature ϕ , but its inner product : Kernel Function

$$K(\mathbf{x}, \mathbf{x}') = \langle \phi(\mathbf{x}), \phi(\mathbf{x}') \rangle$$

Kernel Trick

$$f(\mathbf{x}) = \beta_0 + \sum_{i=1}^n y_i \alpha_i K(\mathbf{x}_i, \mathbf{x})$$

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

Linear Kernel

$$K(\mathbf{x}_i, \mathbf{x}_j) = \exp(-\gamma(\mathbf{x}_i - \mathbf{x}_j)^T(\mathbf{x}_i - \mathbf{x}_j))$$

*Gaussian Kernel
(Radial Basis function)*

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

polynomial Kernel

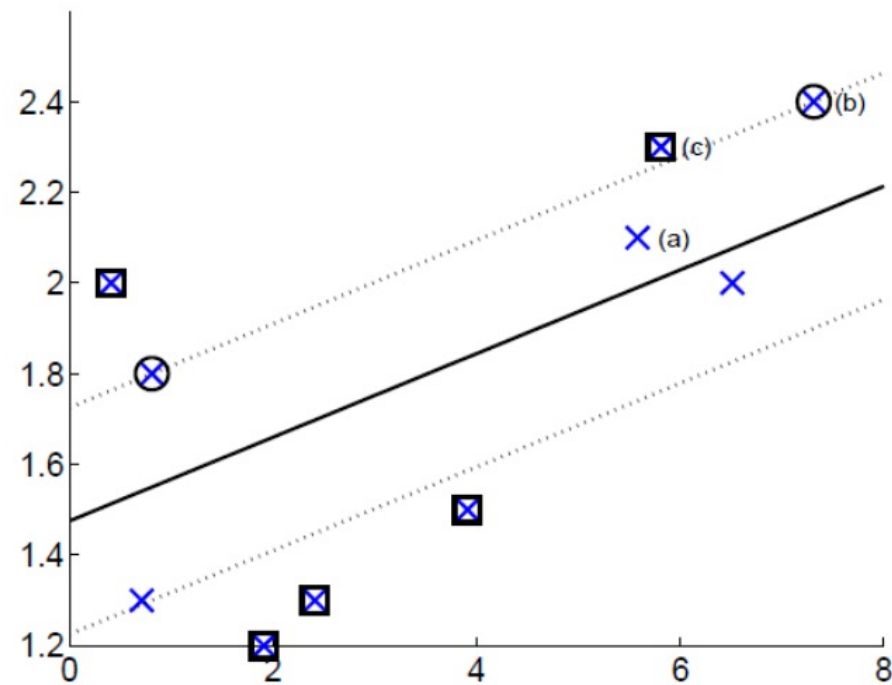
$$K(\mathbf{x}_i, \mathbf{x}_j) = \tanh(k_1 \mathbf{x}_i^T \mathbf{x}_j + k_2)$$

Sigmoid Kernel

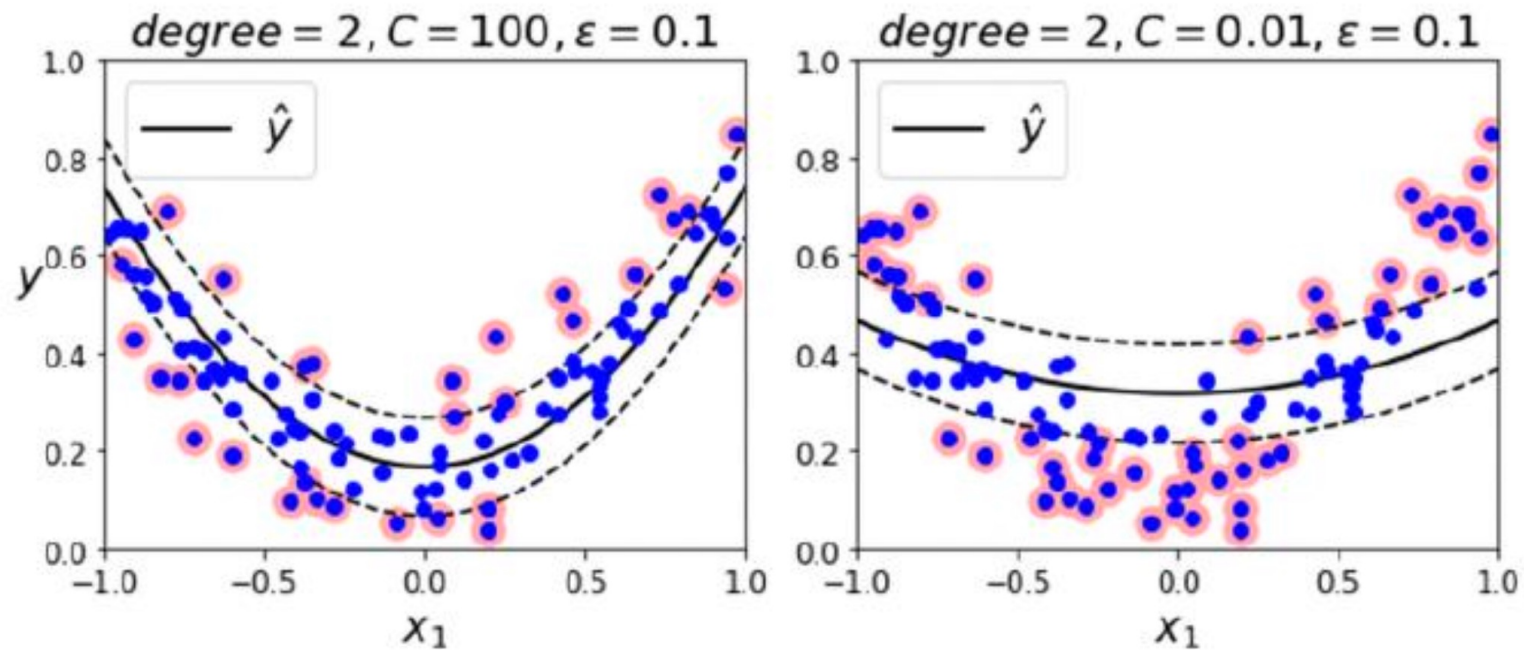
Support Vector Machine

- Regression

SVM - Regression



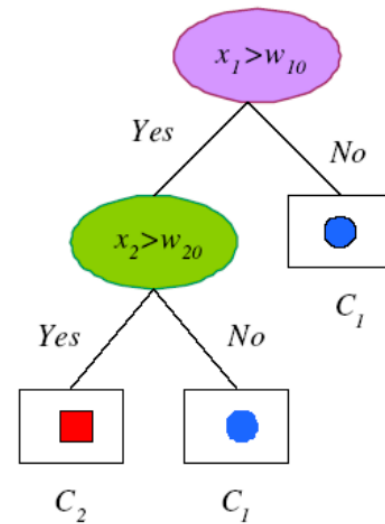
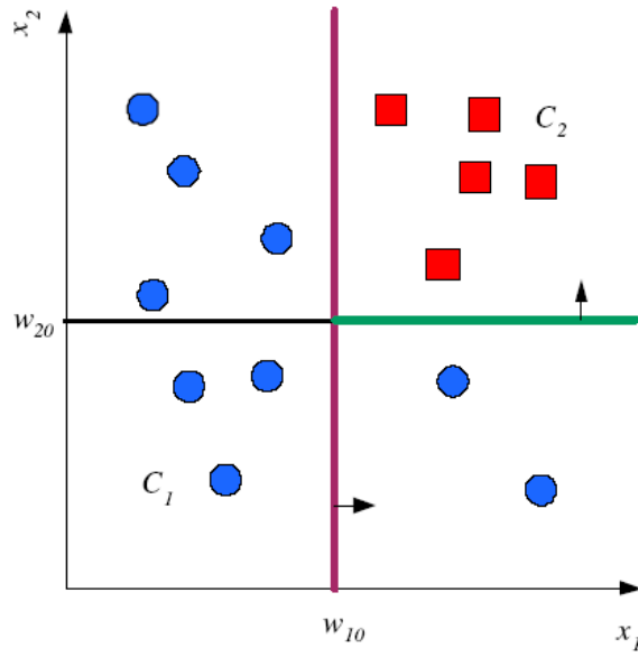
SVM - Regression



Decision Tree

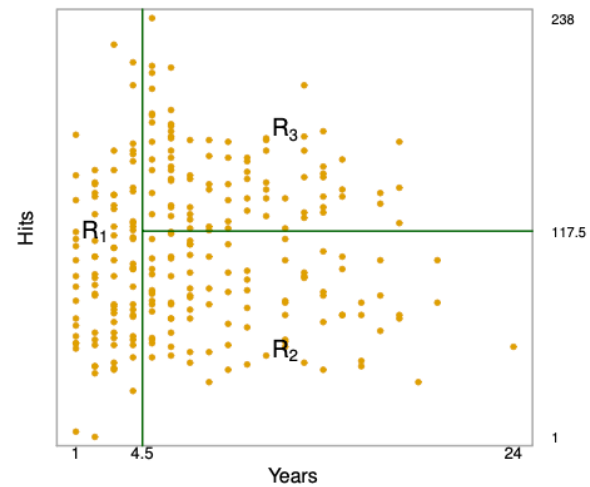
Decision Tree

- Classification

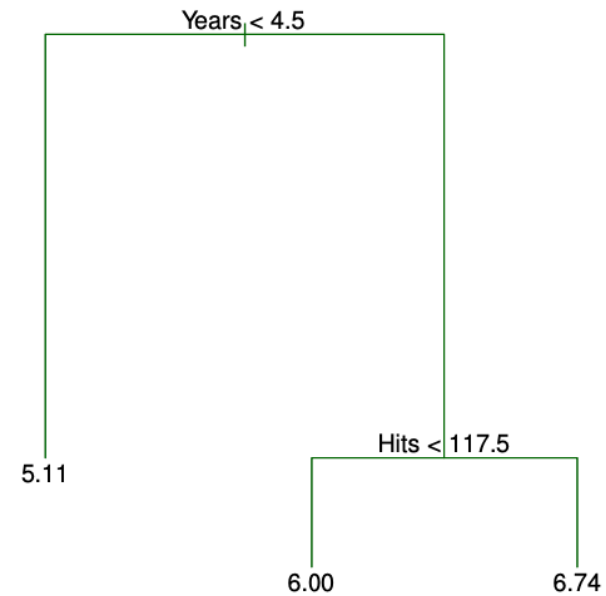


Decision Tree

- Regression



(a) Partitioned Predictor Space



(b) Fitted Model

Impurity

- Minimize the impurity of leaf node
- “Best Split” = Minimize the total impurity
- Measure of Impurity \rightarrow Classification : Entropy $\sum_k p_k^\ell (1 - p_k^\ell)$ or $\sum_k p_k^\ell \log p_k^\ell$
 \rightarrow Regression : MSE $\sum_i (y_i^\ell - \bar{y}^\ell)^2$

Pruning Trees

Size of tree is a tuning parameter.

- Too large Tree: Overfitting (High Variance/Low Bias)
- Too small Tree: Underfitting (Low Variance/High Bias)

Pruning methods

- Prepruning : Early stopping
- Postpruning : Grow the whole tree then prune subtrees