

→ SVM with quadratic penalty.

$$\text{Minimize } \frac{1}{2} \|w\|^2 + c \sum_i \xi_i^2$$

$$\text{subject to } y_i (w^T x^{(i)} + b) \geq 1 - \xi_i \quad \forall i \quad \xi_i \geq 0$$

$$\Rightarrow 1 - y_i (w^T x^{(i)} + b) - \xi_i \leq 0, \quad -\xi_i \leq 0$$

Here w is a vector $\in \mathbb{R}^{(n)}$ same as $x^{(i)}$

ξ is a vector $\in \mathbb{R}^{(n)}$ (No. of data samples)

λ_1, λ_2 are vectors corresponding to $1 - y_i (w^T x^{(i)} + b) - \xi_i \leq 0, \quad -\xi_i \leq 0$

$$L(w, b, \xi, \lambda_1, \lambda_2) = \frac{1}{2} \|w\|^2 + c \sum_{i=1}^n \xi_i^2 + \sum_{i=1}^n \lambda_1^{(i)} (1 - y_i (w^T x^{(i)} + b) - \xi_i) - \sum_{i=1}^n \lambda_2^{(i)} \xi_i$$

$$\begin{aligned} \frac{\partial L}{\partial w_1} &= w_1 + \sum_{i=1}^n \frac{\partial}{\partial w_1} \left(\lambda_1^{(i)} (1 - y_i (w^T x^{(i)} + b) - \xi_i) \right) \\ &= w_1 + \sum_{i=1}^n \lambda_1^{(i)} \left(-y_i \frac{\partial}{\partial w_1} (w_1 x_1^{(i)} + w_2 x_2^{(i)} + \dots + w_n x_n^{(i)}) \right) \end{aligned}$$

$$= w_1 + \sum_{i=1}^n \lambda_1^{(i)} (-y_i x_1^{(i)}) = 0 \Rightarrow w_1 = - \sum_{i=1}^n \lambda_1^{(i)} y_i x_1^{(i)}$$

$$\text{If for other } w_j = - \sum_{i=1}^n \lambda_1^{(i)} y_i x_j^{(i)} \leftarrow [j\text{-DIMENSION}] \quad \text{--- (1)}$$

$$\frac{\partial L}{\partial \xi_j} = 2c \xi_j + \frac{\partial}{\partial \xi_j} \left(\sum_{i=1}^n \lambda_1^{(i)} (1 - y_i (w^T x^{(i)} + b) - \xi_i) \right) + \frac{\partial}{\partial \xi_j} \left(\sum_{i=1}^n \lambda_2^{(i)} \xi_i \right)$$

$$\frac{\partial L}{\partial \hat{\epsilon}_j} = 2c \hat{\epsilon}_j + \frac{\partial}{\partial \hat{\epsilon}_j} \left(\sum_{i=1}^m \lambda_1^{(i)} (1 - y_i (w^T x^{(i)} + b) - \hat{\epsilon}_i) \right) - \frac{\partial}{\partial \hat{\epsilon}_j} \left(\sum_{i=1}^m \lambda_2^{(i)} \hat{\epsilon}_i \right)$$

$$\Rightarrow 2c \hat{\epsilon}_j + (-\lambda_1^{(j)}) - \lambda_2^{(j)} = 0 \Rightarrow \hat{\epsilon}_j = \frac{\lambda_1^{(j)} + \lambda_2^{(j)}}{(2c)} \quad \text{--- (2)} \quad [j - \text{DATA POINTS}]$$

$$\frac{\partial L}{\partial b} = \frac{\partial}{\partial b} \left(\sum_{i=1}^m \lambda_1^{(i)} (1 - y_i (w^T x^{(i)} + b) - \hat{\epsilon}_i) \right) = 0$$

$$= - \sum_{i=1}^m \lambda_1^{(i)} y_i = 0 \Rightarrow \lambda_1^{(1)} y_1 + \lambda_1^{(2)} y_2 + \dots + \lambda_1^{(m)} y_m = 0 \quad \text{--- (3)}$$

$$\vec{w} = \sum_{i=1}^m \lambda_1^{(i)} y_i \vec{x}^{(i)} \quad \text{from (1)}$$

$$\Rightarrow \min \frac{1}{2} \|\vec{w}\|^2 + c \sum_i \hat{\epsilon}_i^2 + \sum_{i=1}^m \lambda_1^{(i)} (1 - y_i (w^T x^{(i)} + b) - \hat{\epsilon}_i) - \sum_{i=1}^m \lambda_2^{(i)} (\hat{\epsilon}_i)$$

$$\Rightarrow \frac{1}{2} \left(\sum_{i=1}^m \sum_{j=1}^m \lambda_1^{(i)} \lambda_1^{(j)} y_i y_j \vec{x}^{(i)} \cdot \vec{x}^{(j)} \right) + c \sum_{i=1}^m \left(\frac{\lambda_1^{(i)} + \lambda_2^{(i)}}{2c} \right)^2 \quad \left[\lambda_1^{(i)} = -\lambda_2^{(i)} + 2c \hat{\epsilon}_i \right]$$

$$+ \sum_{i=1}^m \lambda_1^{(i)} \left(1 - y_i \left(\sum_{j=1}^m \lambda_1^{(j)} y_j \vec{x}^{(j)} \right) x^{(i)} - b y_i - \left(\frac{\lambda_1^{(i)} + \lambda_2^{(i)}}{2c} \right) \right) - \sum_{i=1}^m \lambda_2^{(i)} \left(\frac{\lambda_1^{(i)} + \lambda_2^{(i)}}{2c} \right)$$

$$\Rightarrow \frac{1}{2} \left(\sum_{i=1}^m \sum_{j=1}^m \lambda_1^{(i)} \lambda_1^{(j)} y_i y_j \vec{x}^{(i)} \cdot \vec{x}^{(j)} \right) + \left(\sum_{i=1}^m \lambda_1^{(i)} \right) - b \left(\sum_{i=1}^m \lambda_1^{(i)} y_i \right) - \sum_{i=1}^m \lambda_1^{(i)} y_i \left(\sum_{j=1}^m \lambda_1^{(j)} y_j \vec{x}^{(j)} \right) x^{(i)} - \sum_{i=1}^m \left(\frac{\lambda_1^{(i)2} + 2\lambda_1^{(i)} \lambda_2^{(i)} + \lambda_2^{(i)2}}{2c} \right)$$

$$= -\frac{1}{2} \sum_{i=1}^m \sum_{j=1}^m \lambda_1^{(i)} \lambda_1^{(j)} y_i y_j \vec{x}^{(i)} \cdot \vec{x}^{(j)} + \sum_{i=1}^m \lambda_1^{(i)} - \sum_{i=1}^m \left(\frac{(\lambda_1^{(i)} + \lambda_2^{(i)})^2}{4c} \right)$$

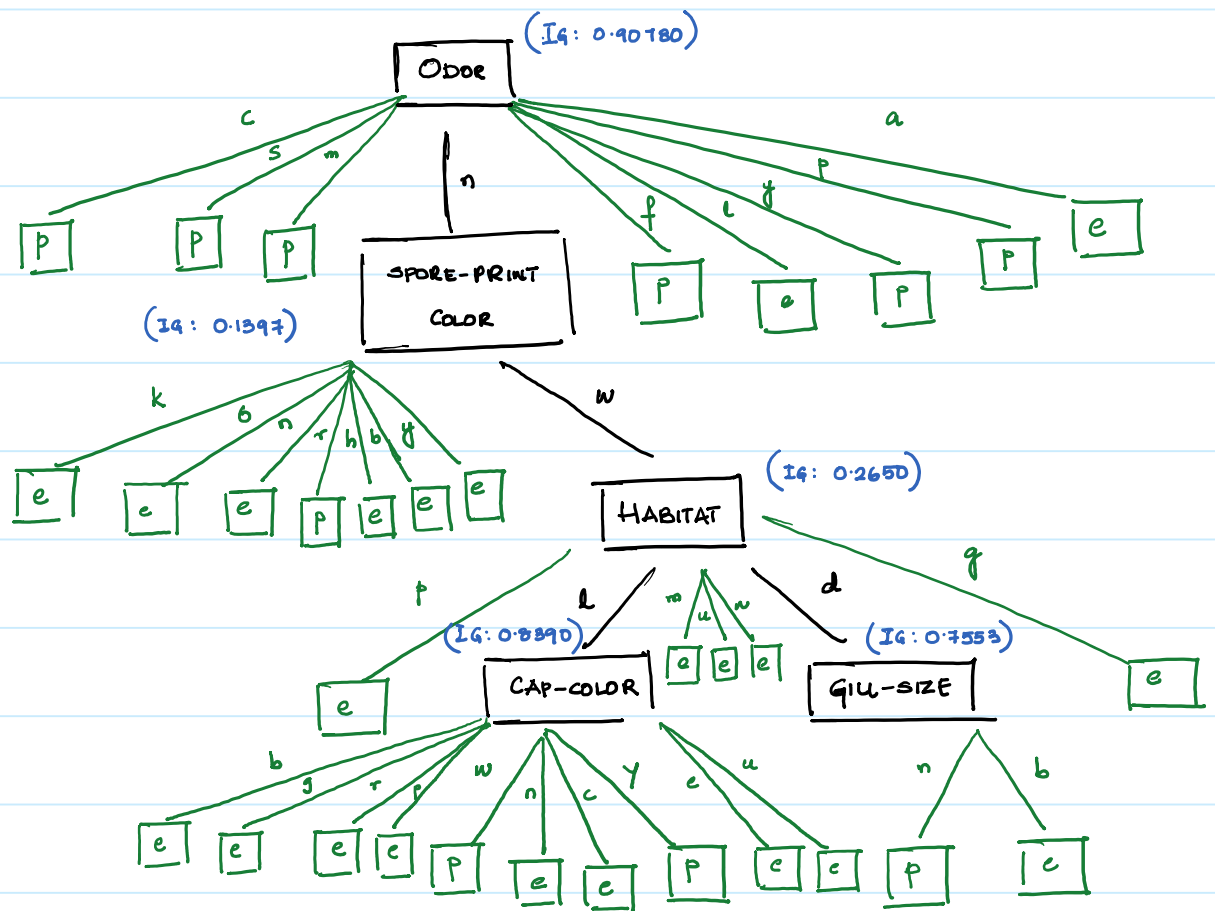
Such that $\sum \lambda_1^{(i)} y_i = 0$ & $\lambda_1^{(i)} \geq 0, \lambda_2^{(i)} \geq 0$

This does not have an upper bound on $\lambda_1^{(i)}$. KERNEL can be applied.

For this term $\vec{x}^{(i)} \cdot \vec{x}^{(j)}$ can be written as $k(\vec{x}^{(i)}, \vec{x}^{(j)})$

3-1)

Decision Tree :



Information gain : ODOR \rightarrow 0.907

Spore Print Color \rightarrow 0.139

HABITAT \rightarrow 0.2650

CAP-COLOR \rightarrow 0.8390

GILL SIZE \rightarrow 0.7553

3-2)

Accuracy of the decision tree on test data is 100%

3-3)

Considering decision trees of height 1 i.e. one split

Yes it is optimal on the training data. As INFORMATION GAIN

is a GREEDY HEURISTIC and for trees of height 1 its result is optimal.

For multilevel trees, tree produced might not be optimal.