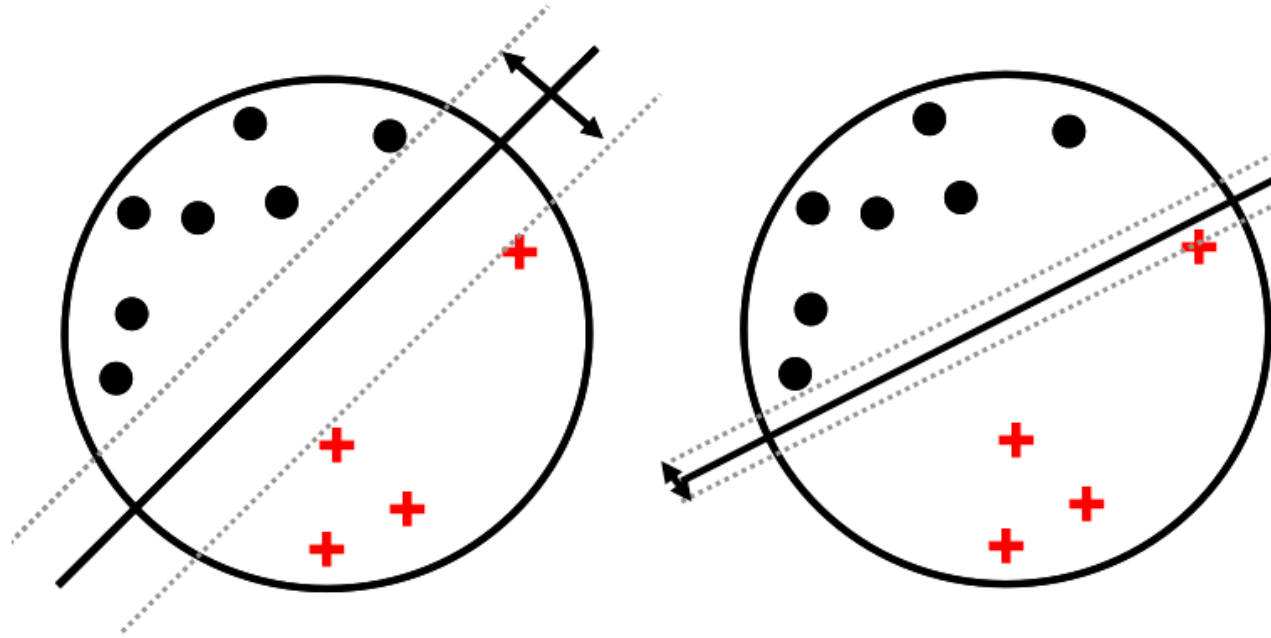# CS 461: Machine Learning Principles

Class 13: Oct. 17

SVM (SMO Algorithm)

& Decision Tree

Instructor: Diana Kim

Computing a Large Margin Classifier



There are many hyperplanes to sperate the training data points.
But a maximum margin classifier on training set is desirable.
- high confident separation
- robust to perturbation of data (generalization)

# SVM Problem [Hard Margin SVM]

$$w*, b* = \arg\max_{w,b} \frac{\Delta}{||w||}$$

$$\text{subject to} \quad t_n(w^t x_n + b) \geq \Delta \quad \forall n$$

$$w*, b* = \arg\min_{w,b} \frac{1}{2}||w||^2$$

$$\text{subject to} \quad t_n(w^t x_n + b) \geq 1 \quad \forall n$$

All data points are separated by the minimum  distance of $\Delta/||W||$ hyperplane , we want to maximize the minimum distance.

# SVM Primal and Dual Problem

- Primal

$$w*, b* = \arg\min_{w,b} \frac{1}{2}||w||^2$$

$$\text{subject to} \quad t_n(w^t x_n + b) \geq 1 \quad \forall n$$

- Dual

$$\lambda*_{n=1}^{N} = \arg\max_{\lambda*} \sum_{n=1}^{N} \lambda*_n - \frac{1}{2}\sum_{n=1}^{N}\sum_{m=1}^{N} \lambda_n^* \cdot t_n \lambda_m^* \cdot t_m \cdot \kappa(x_n, x_m)$$

$$\text{subject to} \quad \lambda_n^* \geq 0 \quad \forall n$$

$$\sum_{n=1}^{N} \lambda_n^* \cdot t_n = 0$$

# Advantage of Dual Formulation

- Dual

$$\lambda *_{n=1}^{N} = \arg\max_{\lambda*} \sum_{n=1}^{N} \lambda *_n - \frac{1}{2} \sum_{n=1}^{N} \sum_{m=1}^{N} \lambda_n^* \cdot t_n \lambda_m^* \cdot t_m \cdot \kappa(x_n, x_m)$$

$$\text{subject to} \quad \lambda_n^* \geq 0 \quad \forall n$$

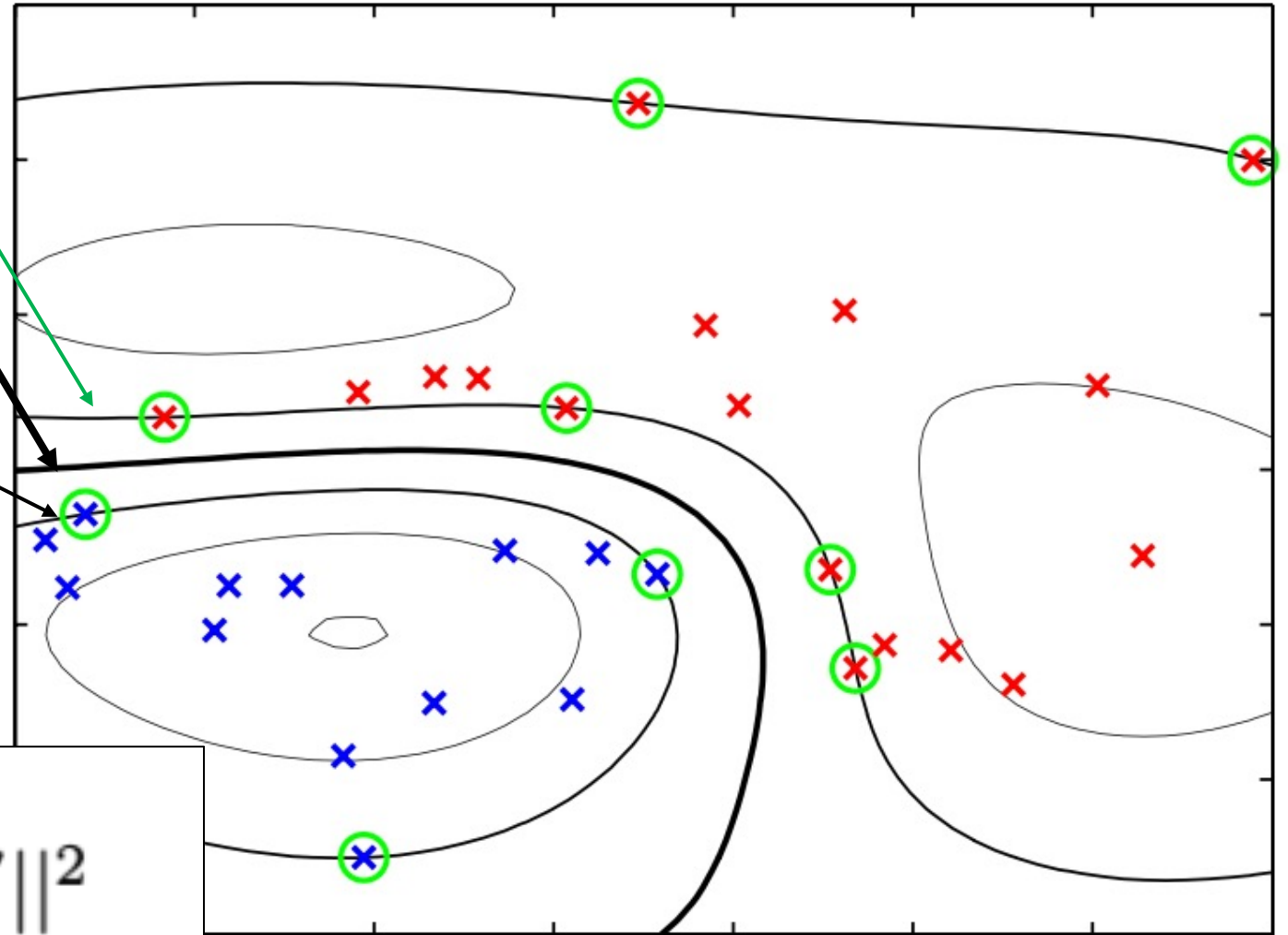$$\sum_{n=1}^{N} \lambda_n^* \cdot t_n = 0$$

- By using kernel trick, maximum margin classifier can be applied efficiently to feature spaces whose dimensionality exceeds the number of data points. (M>>N)
- By using kernel trick, we can make training data linearly separable.

# The Outcomes of Gaussian Kernel SVM
(1) **Decision hyperplane**
(2) Support Vectors (Samples)
(3) Margins



Gaussian Kernel

$$\kappa(x, x') = \exp \frac{-\|x - x'\|^2}{2\sigma^2}$$

From the dual function, we compute $\lambda * (n = 1, .. N)$.
Then we can build a maximum margin classifier.
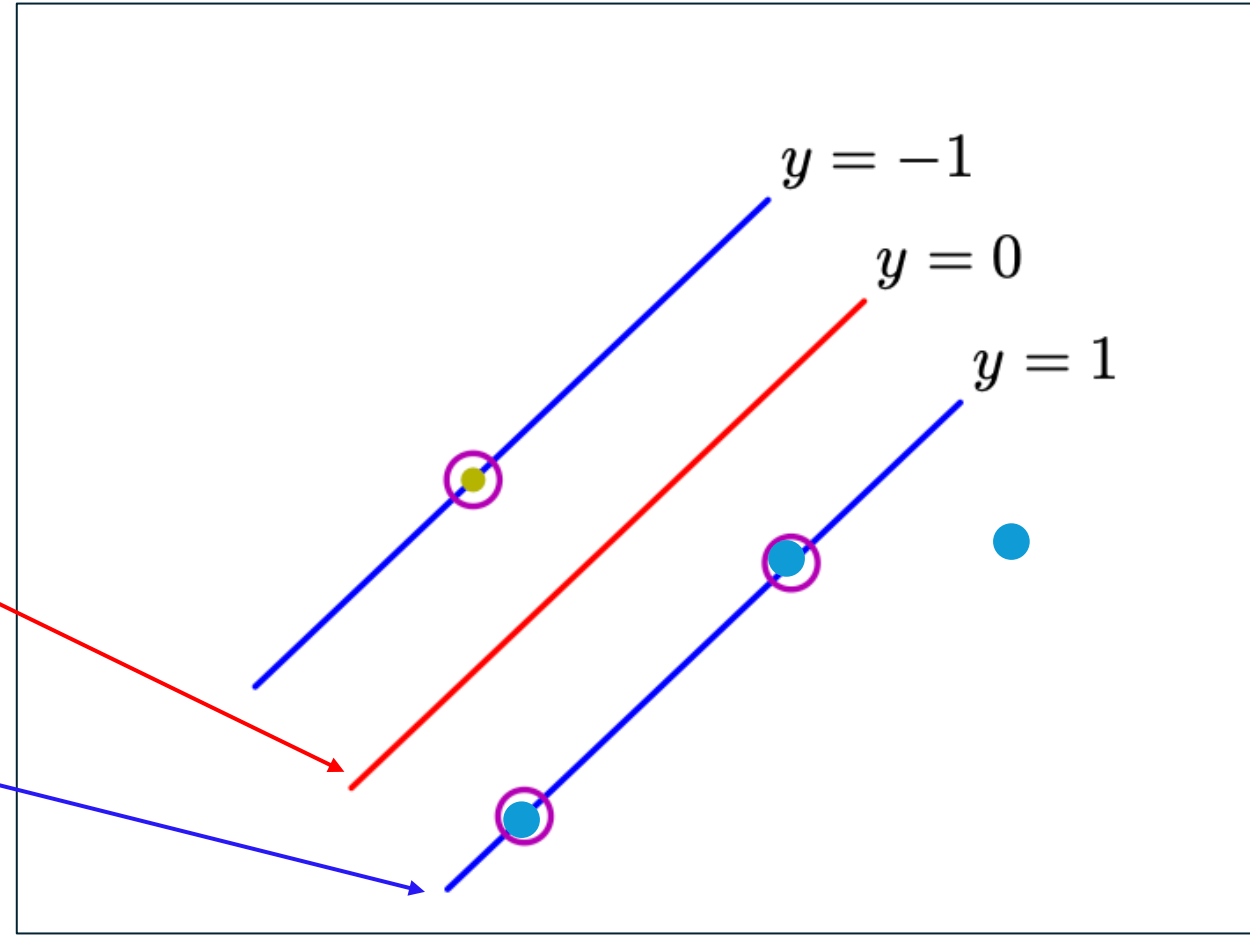A subset of the training data points are used to build the classifier.

$$y(x) = \sum_{n=1}^{N} \lambda *_n t_n \kappa(x_n, x) + b$$

The second and third KKT $=\begin{cases} \lambda *_n = 0 & \text{if} \quad w *^t x_n - 1 > 0 \\ \\ \lambda *_n > 0 & \text{if} \quad w *^t x_n - 1 = 0 \quad \text{Support Vectors } x_n! \end{cases}$

The Outcomes of Linear Kernel SVM
(1) Decision hyperplane.
(2) Support Vectors (Samples)
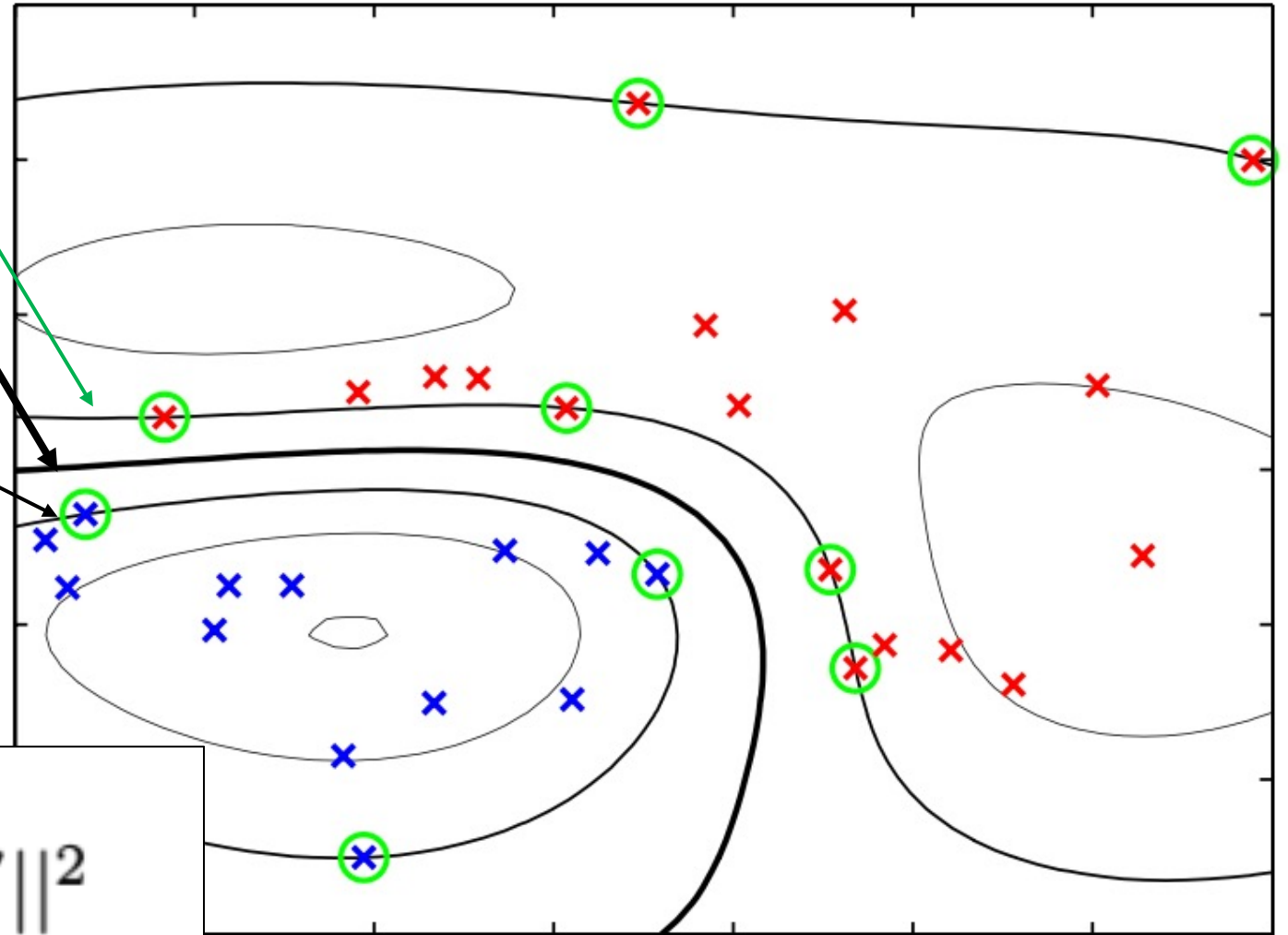(3) Margins.

$y = -1$

$y = 0$

$y = 1$

# The Outcomes of Gaussian Kernel SVM
(1) **Decision hyperplane**
(2) Support Vectors (Samples)
(3) Margins



Gaussian Kernel

$$\kappa(x, x') = \exp \frac{-\|x - x'\|^2}{2\sigma^2}$$

# About Gaussian Kernel

Gaussian Kernel

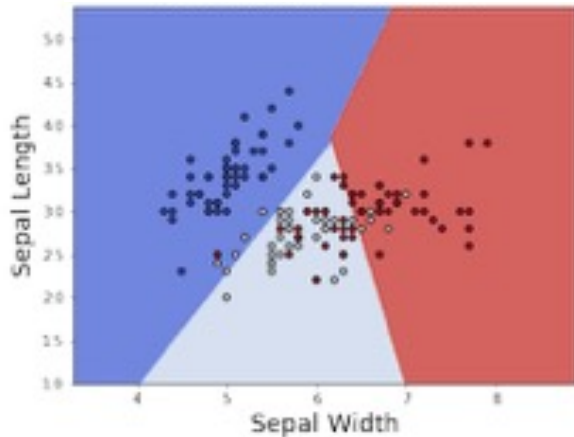$$\kappa(x, x') = \exp \frac{-||x - x'||^2}{2\sigma^2}$$

- Gaussian kernel embeds the infinite dimensional feature space. However, SVM with Gaussian is not sensitive to # data points because it the maximum margin boundary can be defined by two $+/-$ data samples. (performance/complexity)

- The model complexity depends on $\sigma$.

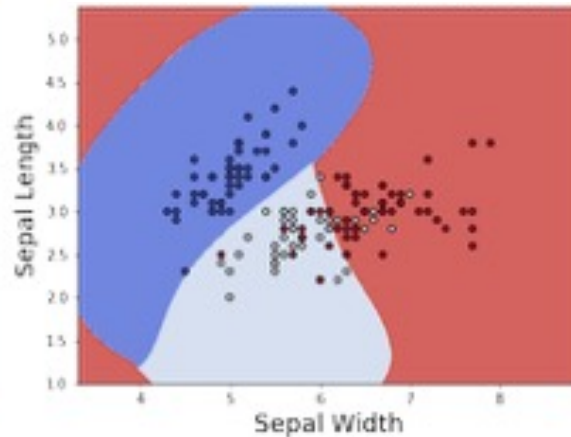The Effect of Gamma on the number of Support vectors & decision Boundary

$$\gamma = \frac{1}{\sigma^2}$$

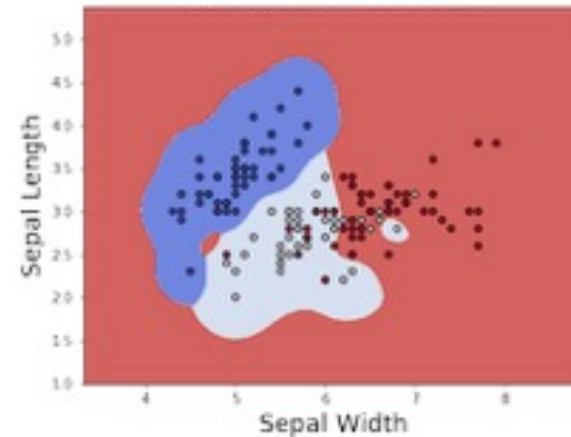From https://www.kaggle.com/code/gorkemgunay/understanding-parameters-of-svm
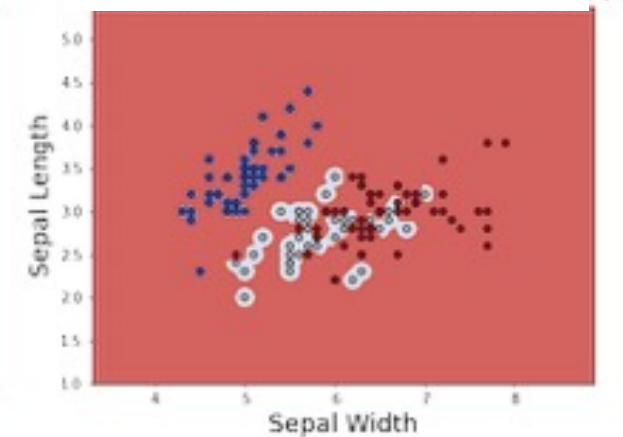
Gamma : 0.01    *Gamma* :1    Gamma :10    Gamma :500



- Small gamma: some representative samples become support vectors.
- Large gamma: individual samples become support vectors

What happens for SVM if data samples are not linearly separable?

For each data points,
the constraint to define the margin is
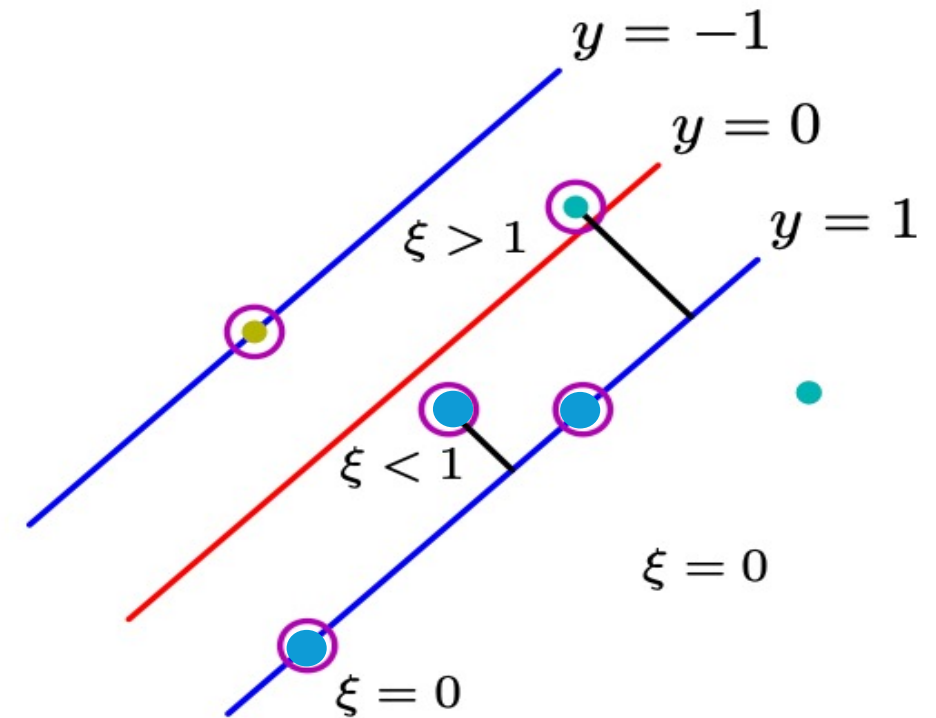relaxed by the slack variable $\xi_n \geq 0$

- hard margin

$$t_n(w^t x_n + b) \geq 1 \quad \forall n$$

- soft margin

$$t_n(w^t x_n + b) \geq \boxed{1 - \xi n} \quad \text{and} \quad \xi n \geq 0 \quad \forall n$$

+ margin can be reduced; even it can be a negative value.

$y = -1$

$y = 0$

$y = 1$

$\xi > 1$

$\xi < 1$

$\xi = 0$

$\xi = 0$

- Hard Margin Case (Exact-Separable)

$$w*, b* = \arg\min_{w,b} \frac{1}{2}||w||^2$$

$$\text{subject to} \quad t_n(w^t x_n + b) \geq 1 \quad \forall n$$

- Soft Margin Case (Non-Separable)

$$w*, b* = \arg\min_{w,b} C \cdot \sum_{n=1}^{N} \xi_n + \frac{1}{2}||w||^2$$

$$\text{subject to} \quad t_n(w^t x_n + b) \geq 1 - \xi_n \quad \forall n$$

$$\text{subject to} \quad \xi_n \geq 0 \quad \forall n$$

# Soft Margin SVM Primal and Dual Problem

- Primal

$$w*, b* = \arg\min_{w,b} C \cdot \sum_{n=1}^{N} \xi_n + \frac{1}{2}||w||^2$$

$$\text{subject to} \quad t_n(w^t x_n + b) \geq 1 - \xi_n \quad \forall n$$

$$\text{subject to} \quad \xi_n \geq 0 \quad \forall n$$

- Dual

$$\lambda*_{n=1}^{N} = \arg\max_{\lambda*} \sum_{n=1}^{N} \lambda *_n - \frac{1}{2} \sum_{n=1}^{N} \sum_{m=1}^{N} \lambda_n^* \cdot t_n \lambda_m^* \cdot t_m \cdot \kappa(x_n, x_m)$$

$$\text{subject to} \quad 0 \leq \lambda_n^* \leq C$$

$$\sum_{n=1}^{N} \lambda_n^* \cdot t_n = 0$$

# KKT conditions

$$\nabla_w L(w, b, \lambda^N_{n=1}, \xi^N_{n=1}, \mu^M_{n=1}) = \vec{w} - \sum_{n=1}^{N} (\lambda_n \cdot t_n \cdot \vec{x_n})$$

$$\nabla_b L(w, b, \lambda^N_{n=1}, \xi^N_{n=1}, \mu^M_{n=1}) = \sum_{n=1}^{N} \lambda_n \cdot t_n$$

$$\nabla_{\xi_n} L(w, b, \lambda^N_{n=1}, \xi^N_{n=1}, \mu^M_{n=1}) = C - \lambda_n - \mu_n$$

$$\lambda_n \cdot \{t_n(w^t x_n) + b - 1 + \xi_n\} = 0$$

$$\mu_n \xi n = 0$$

# Soft Margin SVM Primal and Dual Problem

- Primal

$$w*, b* = \underset{w,b}{\arg\min} \; C \cdot \sum_{n=1}^{N} \xi_n + \frac{1}{2}\|w\|^2$$

$$\text{subject to} \quad t_n(w^t x_n + b) \geq 1 - \xi_n \quad \forall n$$

$$\text{subject to} \quad \xi_n \geq 0 \quad \forall n$$

- Dual

$$\lambda*_{n=1}^{N} = \underset{\lambda*}{\arg\max} \sum_{n=1}^{N} \lambda *_n - \frac{1}{2}\sum_{n=1}^{N}\sum_{m=1}^{N} \lambda_n^* \cdot t_n \lambda_m^* \cdot t_m \cdot \kappa(x_n, x_m)$$

$$\text{subject to} \quad 0 \leq \lambda_n^* \leq C$$

$$\sum_{n=1}^{N} \lambda_n^* \cdot t_n = 0$$

- $\lambda*_n = 0$

- Support vectors: $0 < \lambda*_n < C$

- Support vectors: $\lambda*_n = C$

# All data samples must satisfy the KKT condition.

- $\lambda*_n = 0$

- Support vectors: $0 < \lambda*_n < C$

- Support vectors: $\lambda*_n = C$

$\longrightarrow$

"correctly labeled with a room to spare"

$$y_n(w^t x_n + b) \geq 0$$
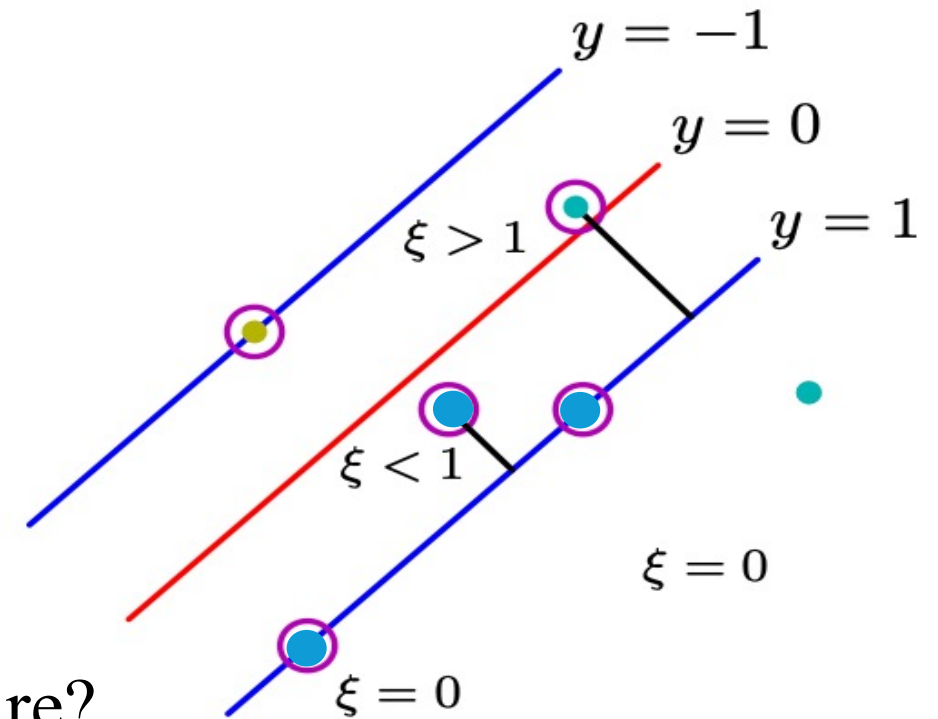
$$y_n(w^t x_n + b) = 1 \quad \text{"unbound"}$$

$$y_n(w^t x_n + b) \leq 0$$

"incorrectly labeled or line within the margin"

- $\lambda*_n = 0$
- Support vectors: $0 < \lambda*_n < C$
- Support vectors: $\lambda*_n = C$

$y = -1$

$y = 0$
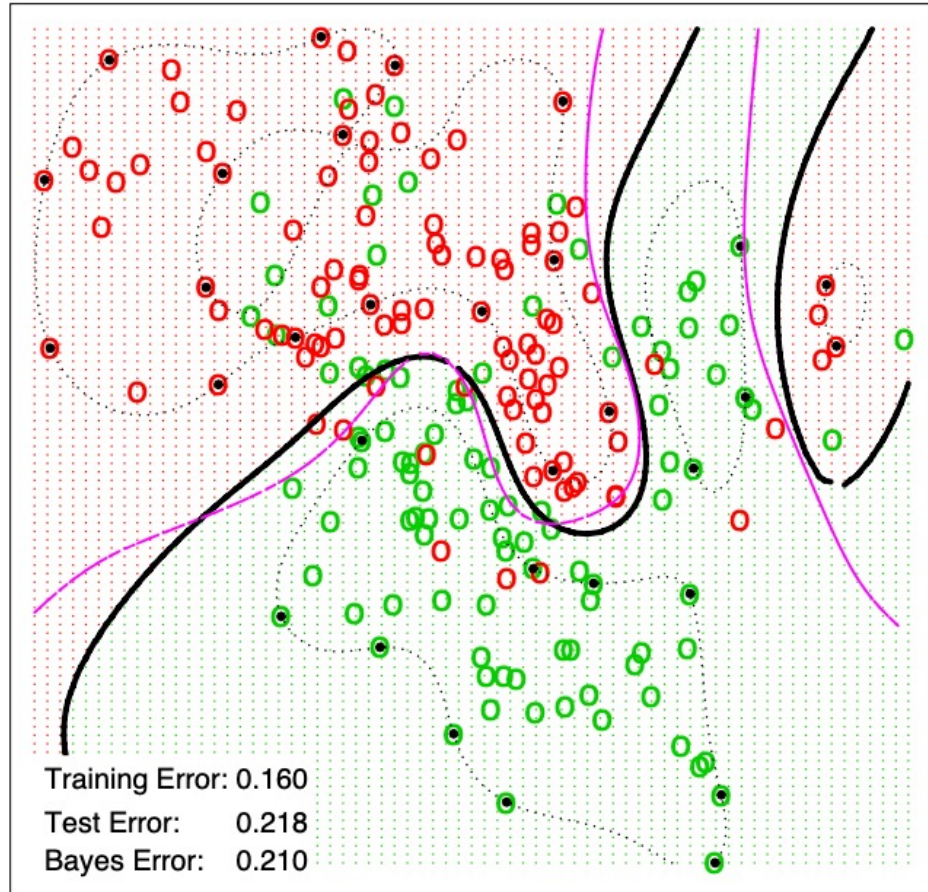
$\xi > 1$

$y = 1$

$\xi < 1$

$\xi = 0$

$\xi = 0$

Q: How many support vectors are in this figure?

Q: Soft margin SVM can handle non-separable data.
Can soft margin algorithm function as regularization for SVM?

Radial Kernel: $C = 2, \gamma = 1$

Training Error: 0.160
Test Error:     0.218
Bayes Error:    0.210

Radial Kernel: $C = 10,000, \gamma = 1$

Training Error: 0.065
Test Error:     0.307
Bayes Error:    0.210

For small C, Large Slack
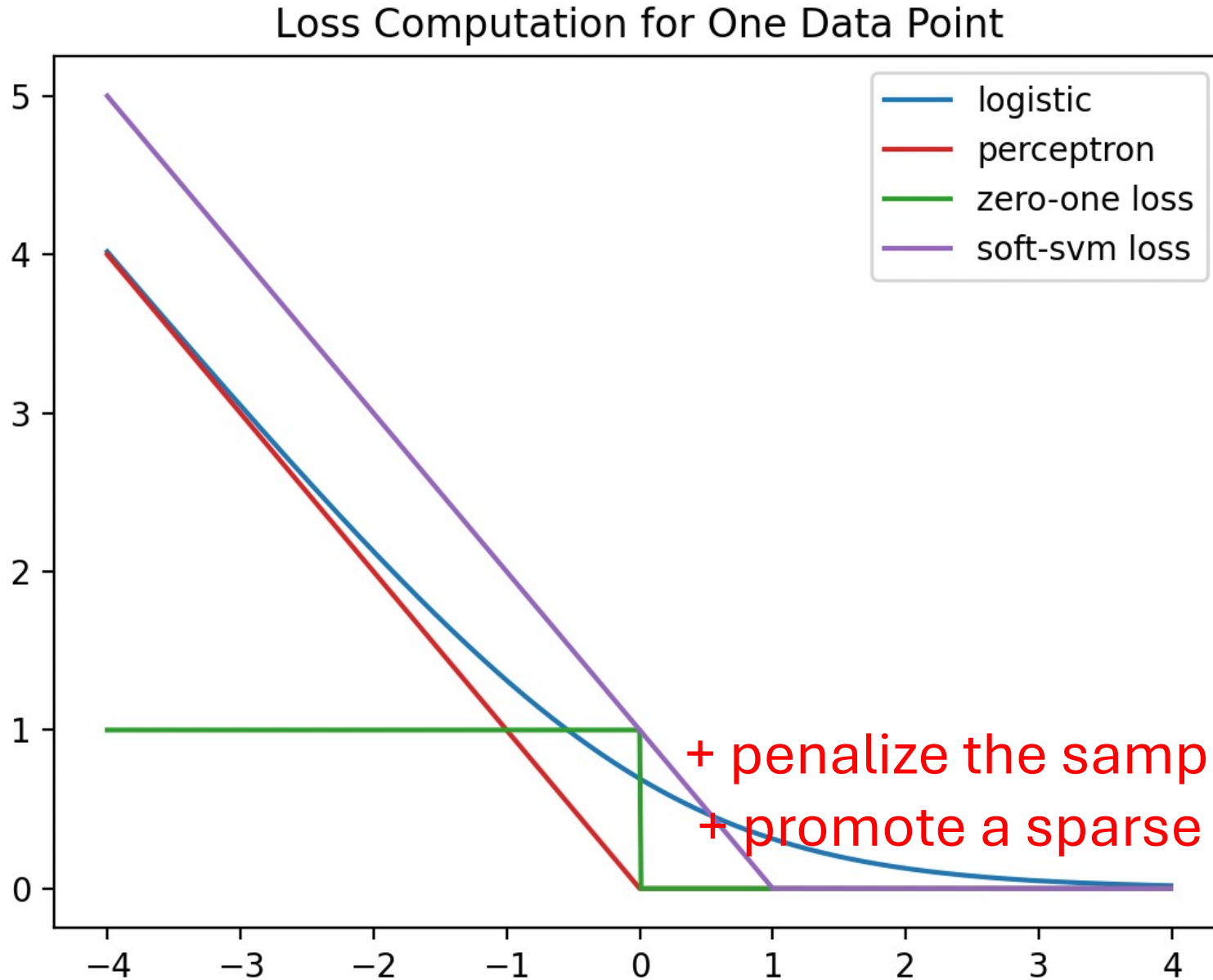
For large C, Small Slack

- Small C: Smooth boundaries for Soft-Margin so better generalization

We need to conduct cross validation to select proper $\sigma^2$ and C.

# The Loss for One Data Point of soft Margin SVM

$$\text{Soft Margin SVM Loss } (x, y) = \begin{cases} \xi_n = 0 & \text{for} \quad yw^t x \geq 1 \\ \\ \xi_n = 1 - yw^t x & \text{for} \quad yw^t x < 1 \end{cases}$$

# Loss Comparison for One Data Point  (Objective Function)



Loss Computation for One Data Point

+ penalize the samples within the margi[n]
+ promote a sparse solution

# Solving Minimal Optimization (SMO)

- Suppose given $\mathcal{D} : \{(x, y): x \in \mathbb{R}^M \text{ and } t \in \{-1, 1\}\}$ and $\kappa(x, x')$.
  A soft SVM dual function is defined below.
  How could we find the optimal $\lambda_1, \lambda_2 \ldots, \lambda_n$?

$$\lambda *_{n=1}^N = \underset{\lambda*}{\arg\max} \sum_{n=1}^N \lambda *_n - \frac{1}{2} \sum_{n=1}^N \sum_{m=1}^N \lambda_n^* \cdot t_n \lambda_m^* \cdot t_m \cdot \kappa(x_n, x_m)$$
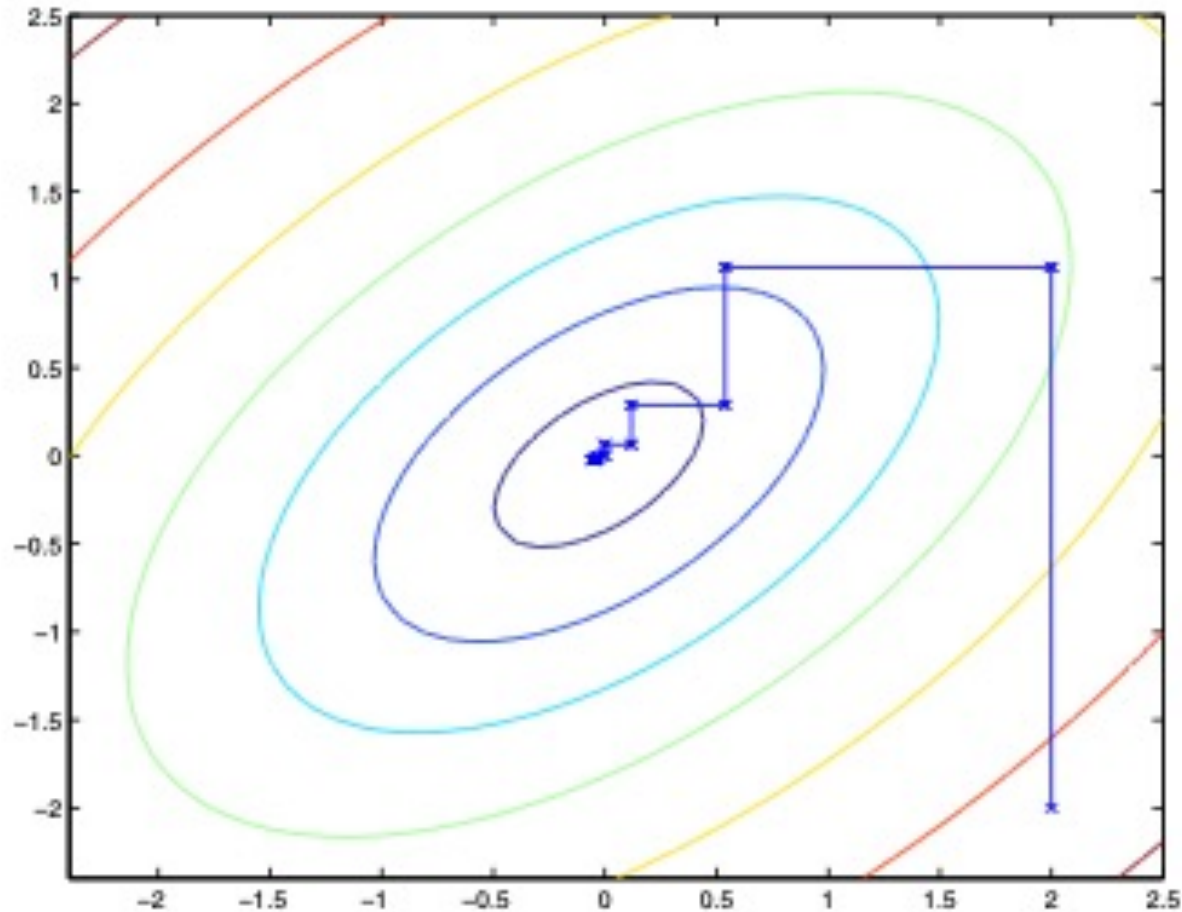
$$\text{subject to} \quad 0 \le \lambda_n^* \le C$$

$$\sum_{n=1}^N \lambda_n^* \cdot t_n = 0$$

Q: Can we solve this problem by using gradient descent?

# Coordinate Ascent Algorithm: One Coordinate at One Time.

At each step, the algorithm finds a minimum along the axis, we selected.

- Suppose given $\mathcal{D} : \{(x, y): x \in \mathbb{R}^M \text{ and } t \in \{-1, 1\}\}$ and $\kappa(x, x')$
  we defined a soft margin SVM dual function below.
  How could we find the optimal $\lambda_1, \lambda_2 \ldots, \lambda_n$?
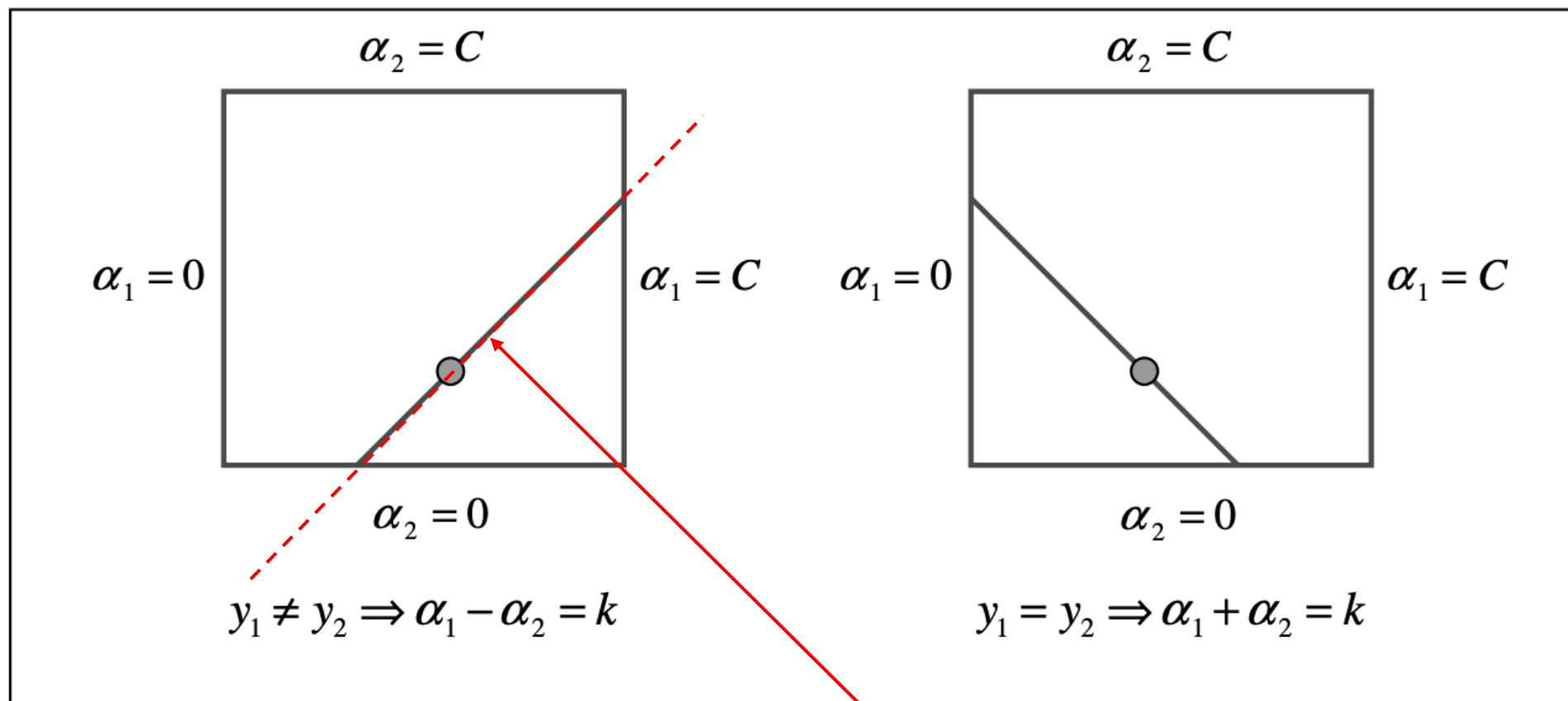
$$\lambda *_{n=1}^{N} = \arg\max_{\lambda *} \sum_{n=1}^{N} \lambda *_n - \frac{1}{2} \sum_{n=1}^{N} \sum_{m=1}^{N} \lambda_n^* \cdot t_n \lambda_m^* \cdot t_m \cdot \kappa(x_n, x_m)$$

$$\text{subject to} \quad 0 \leq \lambda_n^* \leq C$$

$$\sum_{n=1}^{N} \lambda_n^* \cdot t_n = 0$$

Let apply the concept of coordinate ascent to solve dual problem,
but we have <u>the two</u> constraints.

In stead of updating one coordinate,
we will update a pair of two coordinates: $(\alpha_1, \alpha_2)$ but others will be fixed.



$\alpha_2 = C$

$\alpha_1 = 0$          $\alpha_1 = C$

$\alpha_2 = 0$

$y_1 \neq y_2 \Rightarrow \alpha_1 - \alpha_2 = k$

$\alpha_2 = C$

$\alpha_1 = 0$          $\alpha_1 = C$

$\alpha_2 = 0$

$y_1 = y_2 \Rightarrow \alpha_1 + \alpha_2 = k$

$$y_1\alpha_1 + y_2\alpha_2 = -\sum_{i=3}^{N} y_i\alpha_i$$

$$y_1\alpha_1 + y_2\alpha_2 = k$$

$L = \max(0, \alpha_2 - \alpha_1), \qquad H = \min(C, C + \alpha_2 - \alpha_1)$

- we are going to find a minimum along the line.

# Finding a minimum $\alpha_{2(new)}$ along the line : $y_1\alpha_1 - y_2\alpha_2 = k$

$$\alpha_{2new} = \alpha_{2old} + \frac{y_2(E_1 - E_2)}{\eta}$$

This is derived by $\frac{dJ'(\alpha_2)}{d\alpha_2} = 0$
The new alpha point is the minimum solution along the line.

$$\eta = \kappa(x_1, x_1) + \kappa(x_2, x_2) - 2\kappa(x_1, x_2)$$

$$E_1 = \left(\sum_{n=1}^{N} \lambda_n^* t_n \kappa(x_n, x_1) + b\right) - y_1$$

$$E_2 = \left(\sum_{n=1}^{N} \lambda_n^* t_n \kappa(x_n, x_2) + b\right) - y_2$$

For derivations:

https://dsmilab.github.io/Yuh-Jye-Lee/assets/file/teaching/2017_machine_learning/SMO_algorithm.pdf
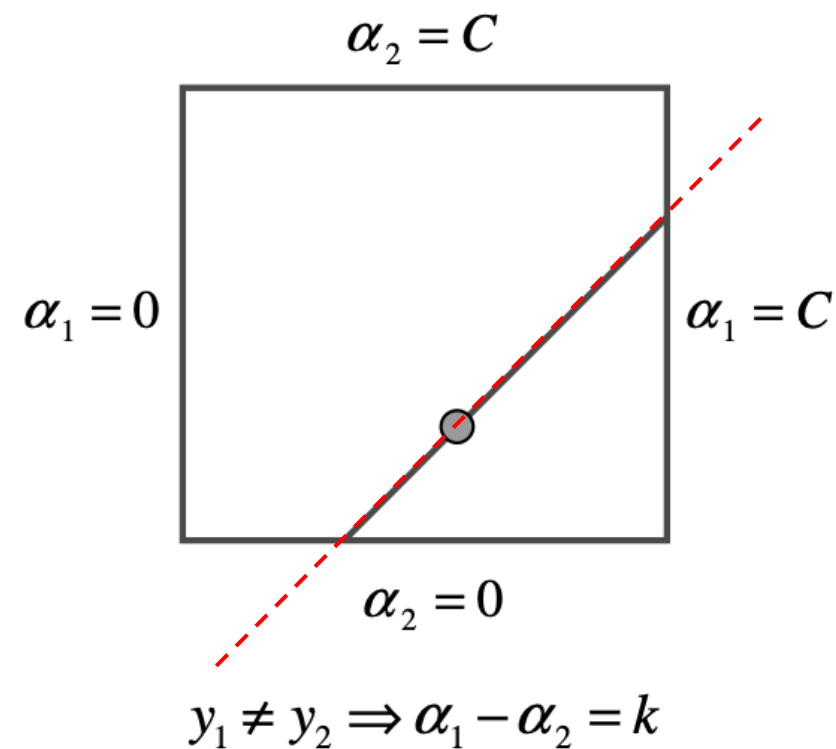
Clipping $\alpha_2$ :
the minimum solution  did not consider the bounds.
Hence, we need clipping



$$\alpha_2 = C$$

$$\alpha_1 = 0 \qquad \alpha_1 = C$$

$$\alpha_2 = 0$$

$$y_1 \neq y_2 \Rightarrow \alpha_1 - \alpha_2 = k$$

$$\alpha_2^{\text{new,clipped}} = \begin{cases} H & \text{if} & \alpha_2^{\text{new}} \geq H; \\ \alpha_2^{\text{new}} & \text{if} & L < \alpha_2^{\text{new}} < H; \\ L & \text{if} & \alpha_2^{\text{new}} \leq L. \end{cases}$$

Once we find $\alpha_2$, then we can update $\alpha_1$.
How?

$$\alpha_2 = C$$

$$\alpha_1 = 0 \qquad \alpha_1 = C$$

$$\alpha_2 = 0$$

$$y_1 \neq y_2 \Rightarrow \alpha_1 - \alpha_2 = k$$

# SMO Algorithm

Repeat until KKT conditions are satisfied for all N training samples within certain tolerance (preset usually $10^{-3} \sim 10^{-2}$)

1. Pick two alphas ($\alpha_1 \ and \ \alpha_2$).
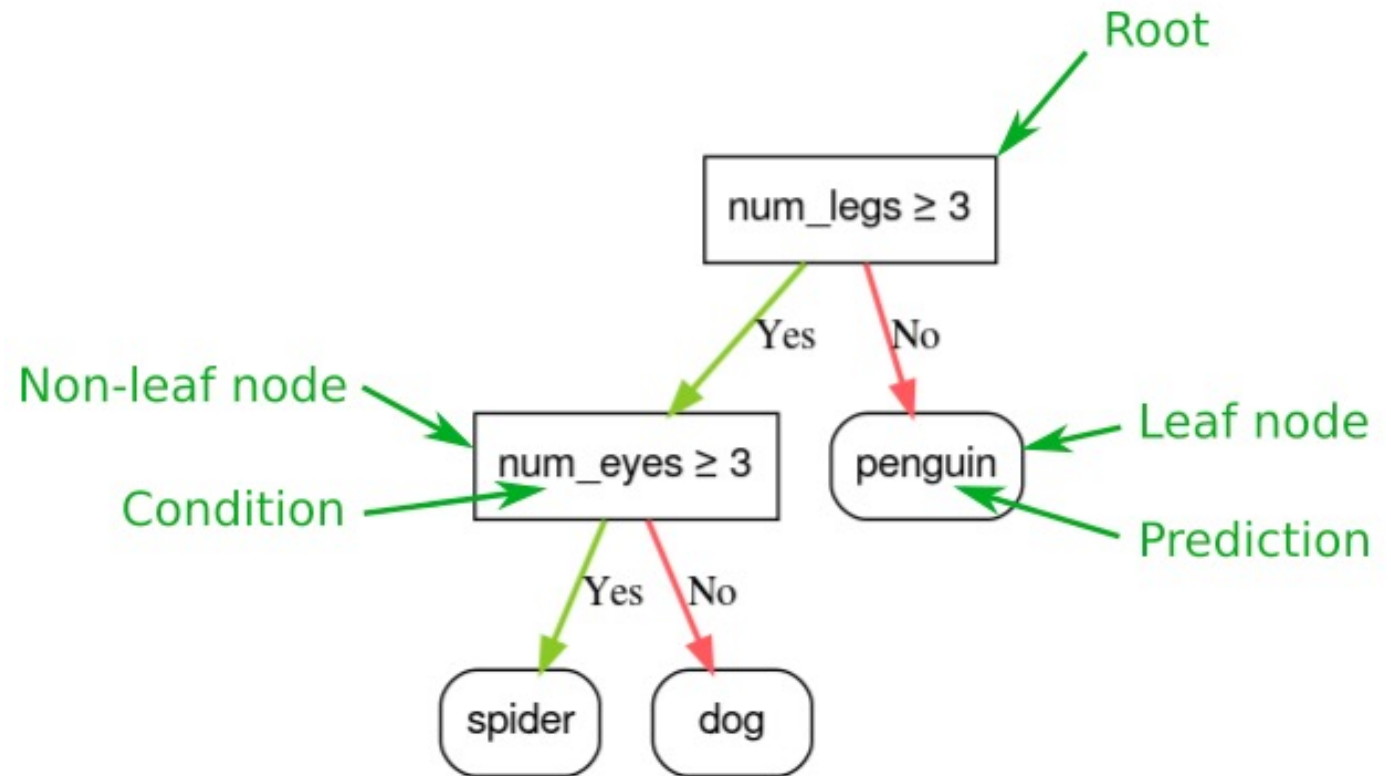
2. Define the range L and H for $\alpha_2$.

3. Compute
$$\alpha_{2new} = \alpha_{2old} + \frac{y_2(E_1 - E_2)}{\eta}$$

4. Clipping by L and H.

5. Update $\alpha_1$.
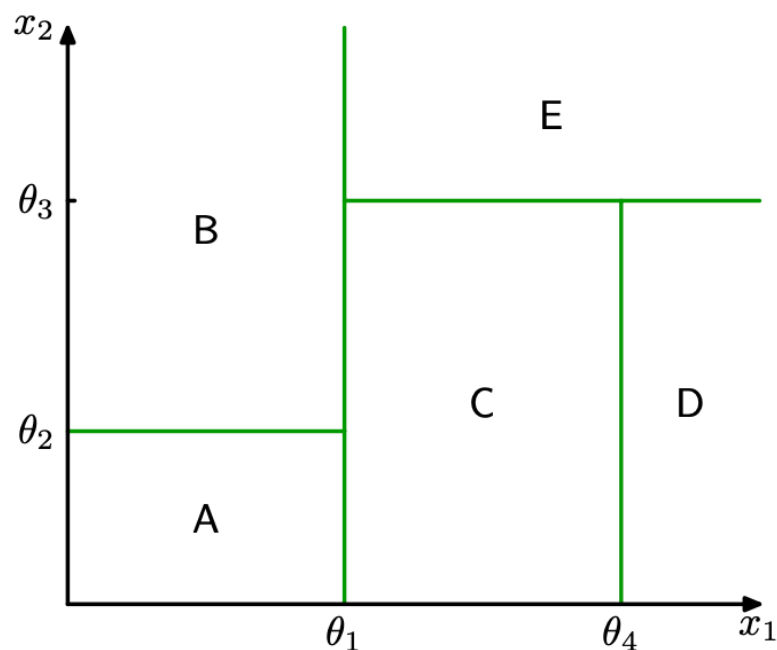
# Decision Tree

# Decision Tree

## :partitioning input space into the regions whose edges are aligned with axes.

Suppose there exist a logical flow in our mind
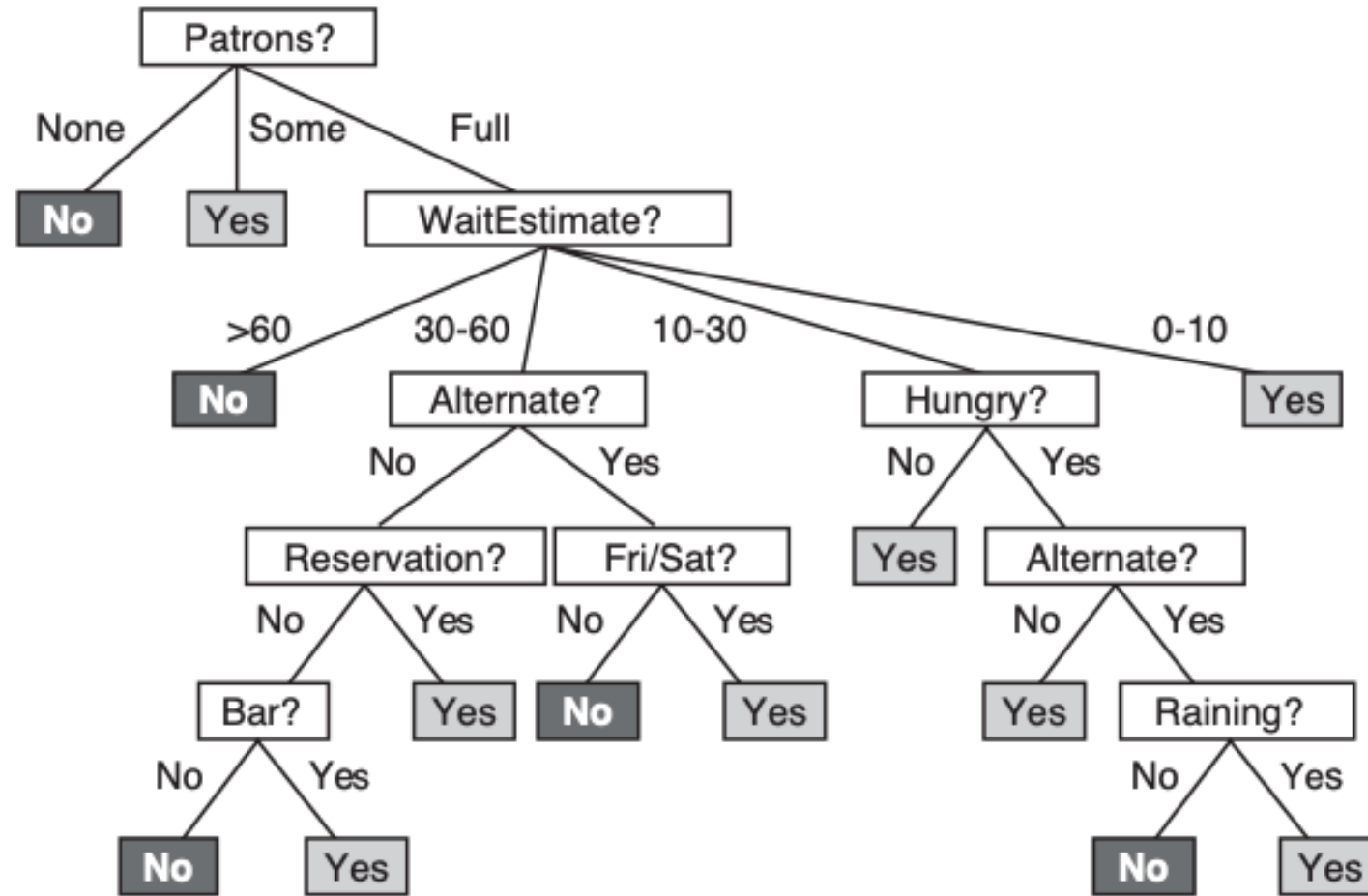as deciding to wait a restaurant. We want to estimate the tree structure.



**Figure 18.2** A decision tree for deciding whether to wait for a table.

# Learning a decision tree.
## Suppose we have categorical data representation.

| Example | Input Attributes | | | | | | | | | | Goal |
| :---: | :---: | :---: | :---: | :---: | :---: | :---: | :---: | :---: | :---: | :---: | :---: |
| | Alt | Bar | Fri | Hun | Pat | Price | Rain | Res | Type | Est | WillWait |
| $x_1$ | Yes | No | No | Yes | Some | $$$ | No | Yes | French | 0–10 | $y_1 = Yes$ |
| $x_2$ | Yes | No | No | Yes | Full | $ | No | No | Thai | 30–60 | $y_2 = No$ |
| $x_3$ | No | Yes | No | No | Some | $ | No | No | Burger | 0–10 | $y_3 = Yes$ |
| $x_4$ | Yes | No | Yes | Yes | Full | $ | Yes | No | Thai | 10–30 | $y_4 = Yes$ |
| $x_5$ | Yes | No | Yes | No | Full | $$$ | No | Yes | French | >60 | $y_5 = No$ |
| $x_6$ | No | Yes | No | Yes | Some | $$ | Yes | Yes | Italian | 0–10 | $y_6 = Yes$ |
| $x_7$ | No | Yes | No | No | None | $ | Yes | No | Burger | 0–10 | $y_7 = No$ |
| $x_8$ | No | No | No | Yes | Some | $$ | Yes | Yes | Thai | 0–10 | $y_8 = Yes$ |
| $x_9$ | No | Yes | Yes | No | Full | $ | Yes | No | Burger | >60 | $y_9 = No$ |
| $x_{10}$ | Yes | Yes | Yes | Yes | Full | $$$ | No | Yes | Italian | 10–30 | $y_{10} = No$ |
| $x_{11}$ | No | No | No | No | None | $ | No | No | Thai | 0–10 | $y_{11} = No$ |
| $x_{12}$ | Yes | Yes | Yes | Yes | Full | $ | No | No | Burger | 30–60 | $y_{12} = Yes$ |

**Figure 18.3**  Examples for the restaurant domain.
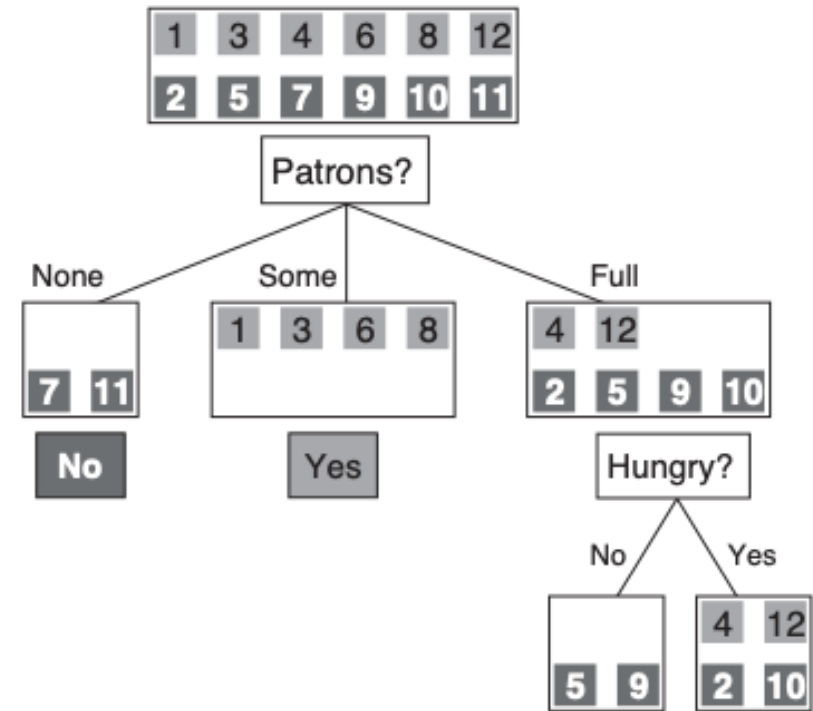
# A Simple Tree Algorithm



1. Choose the best feature to split

2. For each value that feature takes,
   sort training example to leaf nodes.

3. Stop if the leaf contain all training examples with same labels (nothing to divide )

4. Assign each leaf majority vote of labels of training examples.

5. Repeat 1,2,3,4 and create subtrees recursively until training error falls below some threshold.

Two Major Problems in Learning a  Decision Tree

1. <u>Choose the best feature to split (based on mutual information)</u>

2. Decide threshold to make the data points categorical