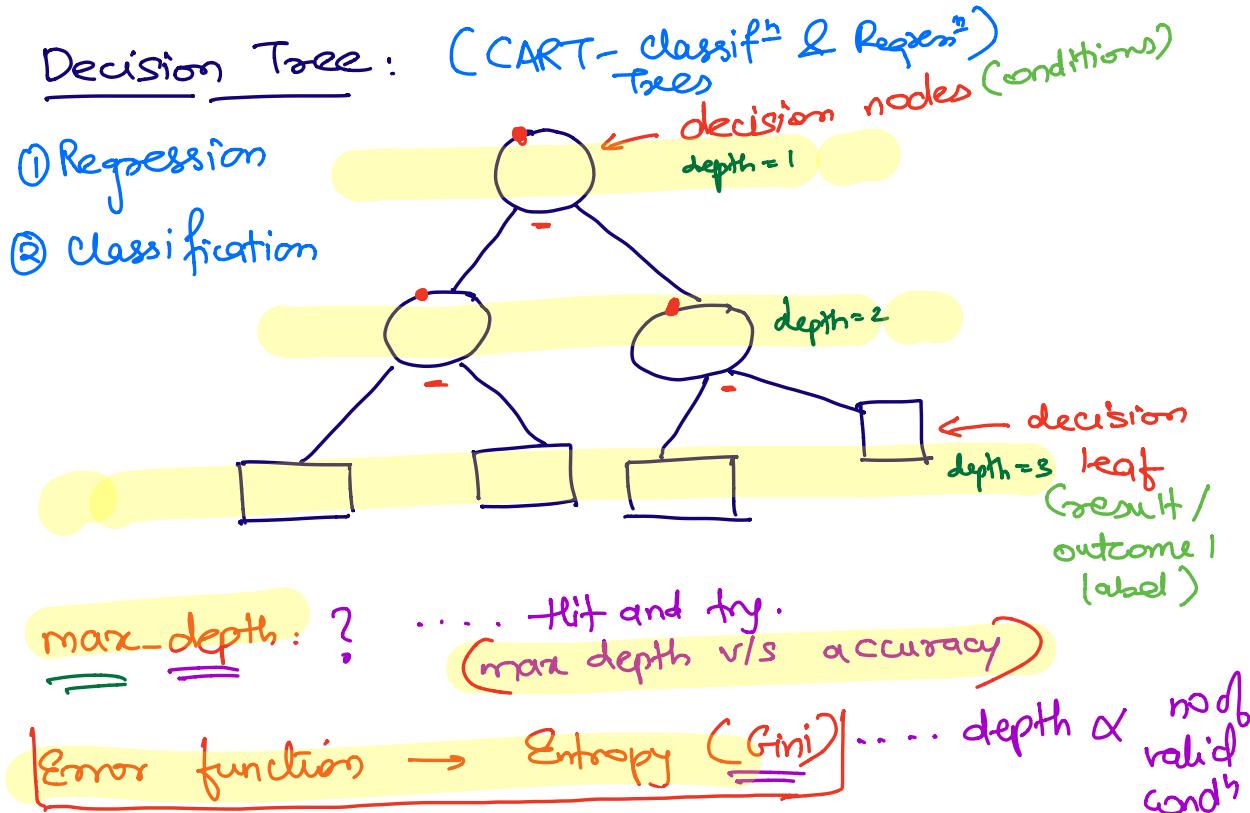


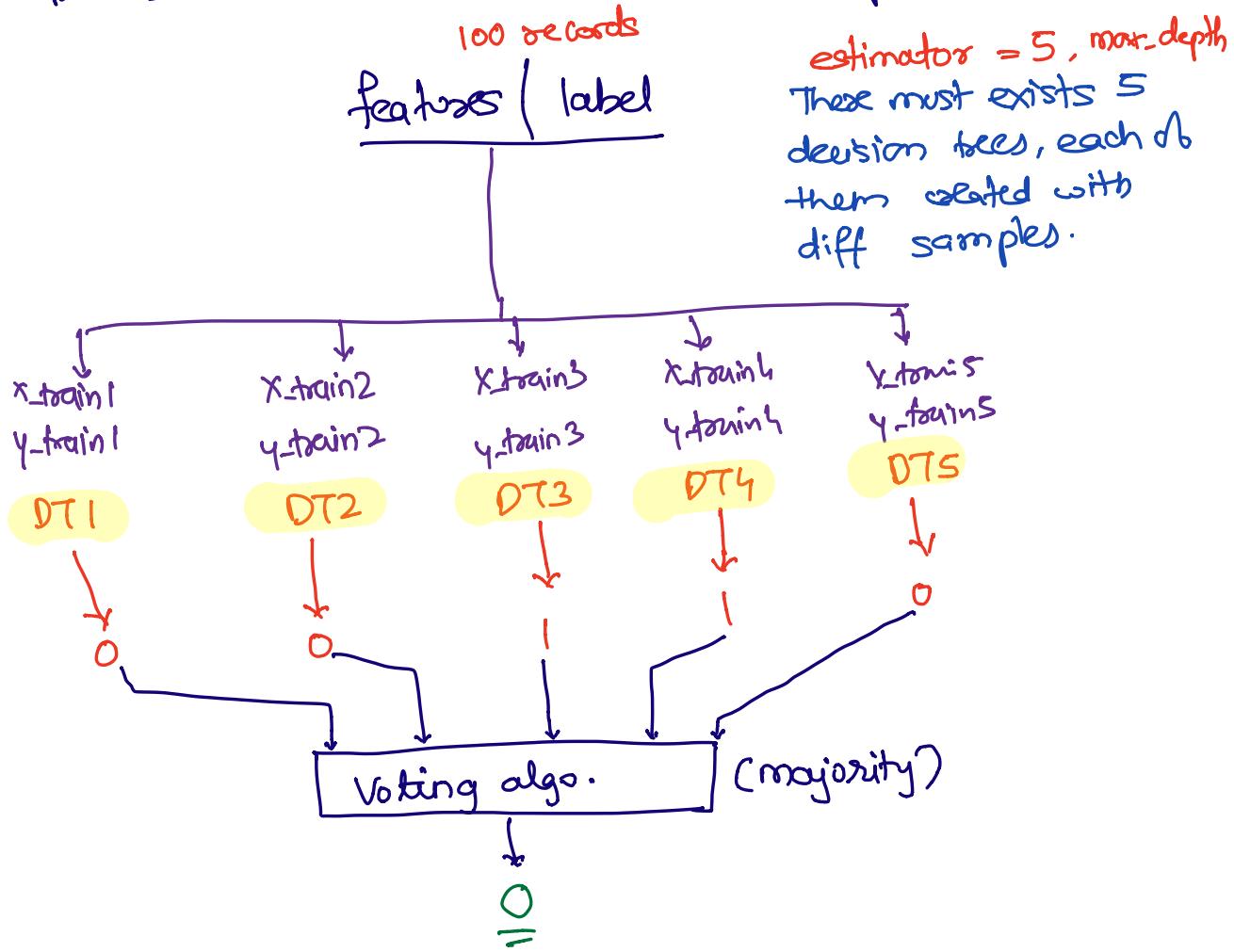
- Decision Tree and Randomforest
- Support Vector Machine
- Feature Selection and Sampling Techniques
- Unsupervised Learning (k-means algo).



Randomforest

- Ensemble Learning
 - Create multiple models with random sample and get the best result based on majority vote.

Random forest is all about creating multiple Decision Trees with different samples.



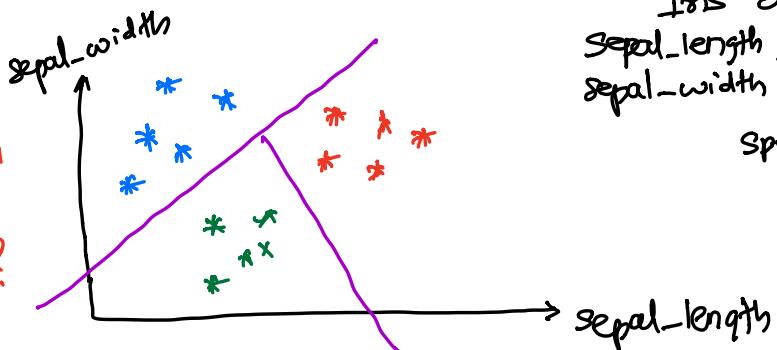
Support Vector Machine

- Works best when you have more than 2 feature columns. (multi-dimensional space)

- SVM
- Regression→ Classification}

All supervised methods.

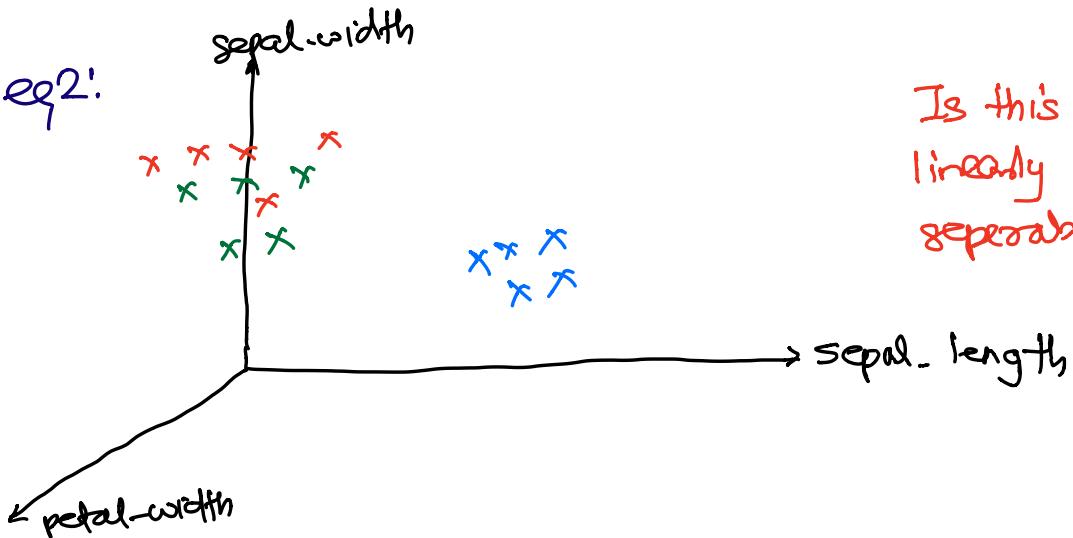
eg 1:
Is this data
linearly
separable?



Iris dataset
sepal-length, petal-length
sepal-width, petal-width

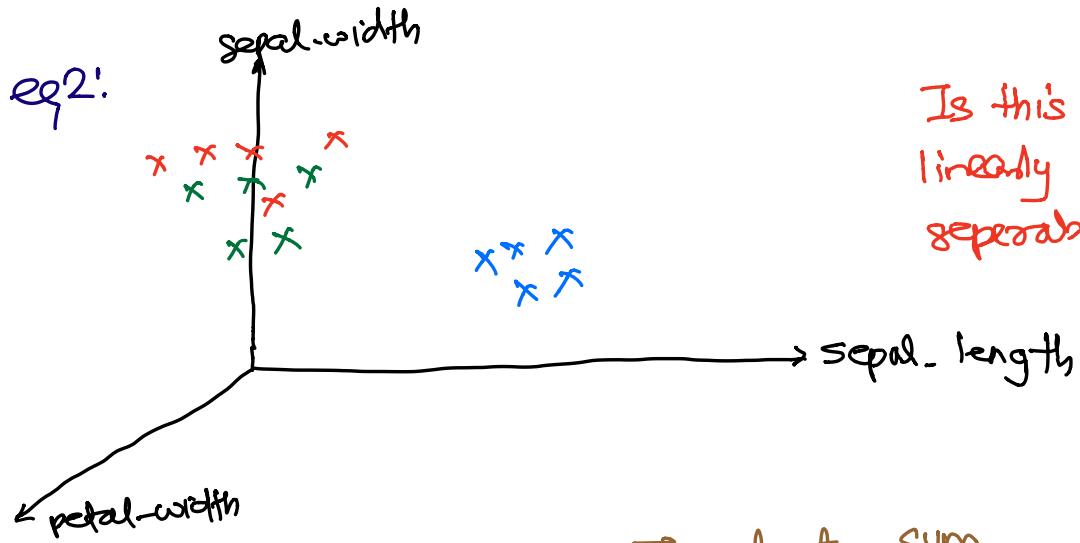
species
→ Setosa
→ Versicolor
→ Virginica

eg 2:



Is this data
linearly
separable?

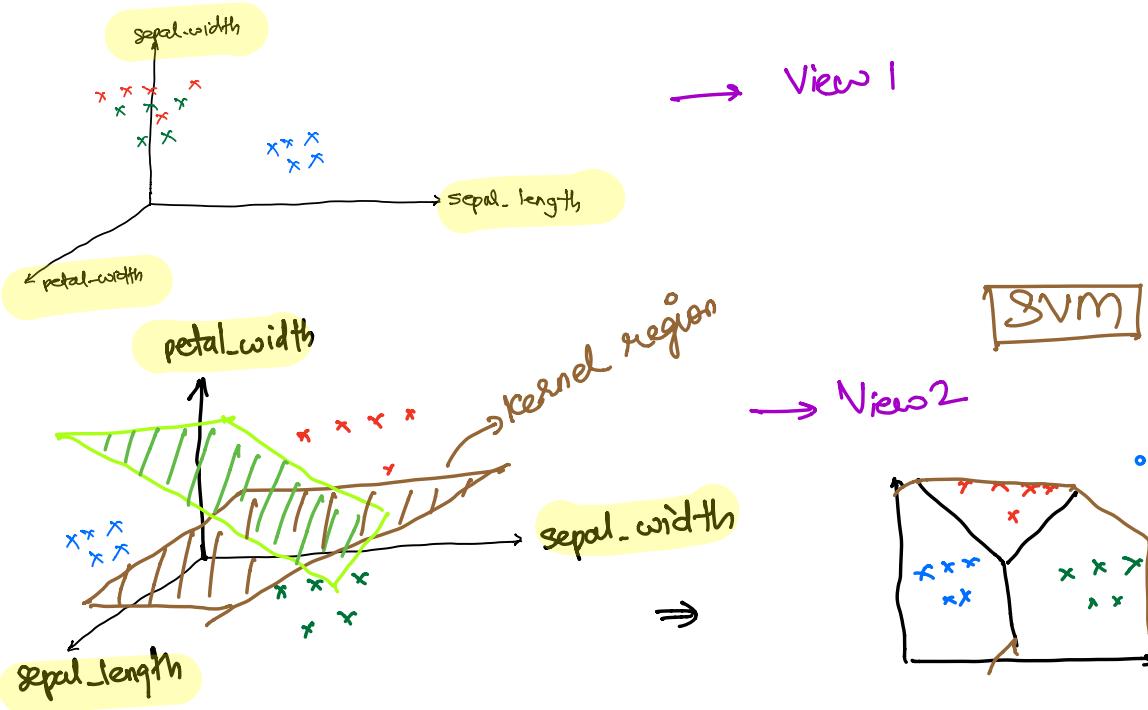
↓
SVM

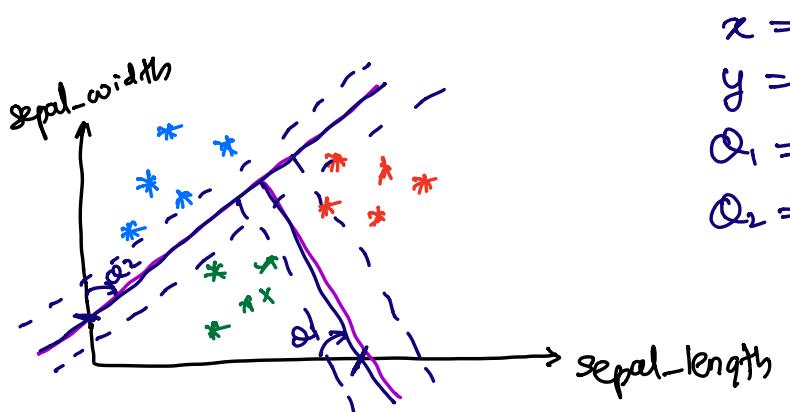


Is this data
linearly
separable?

Input to SVM

SVM tries to figure out can I make this data linearly separable by checking all possible dimensions?





$$x = ?$$

$$y = ?$$

$$\alpha_1 = ?$$

$$\alpha_2 = ?$$

How to figure out if my data can be linearly separated?

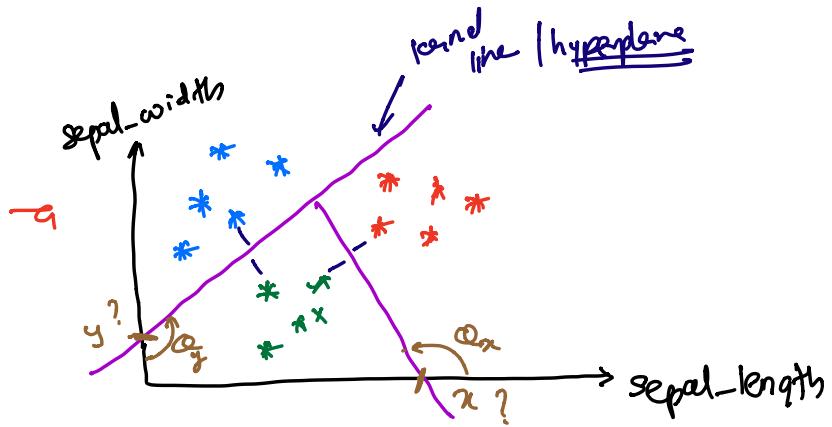
SVM (smiley face)

I will check all possible dimension to figure out if in any dimension, the data is linearly separable or not?

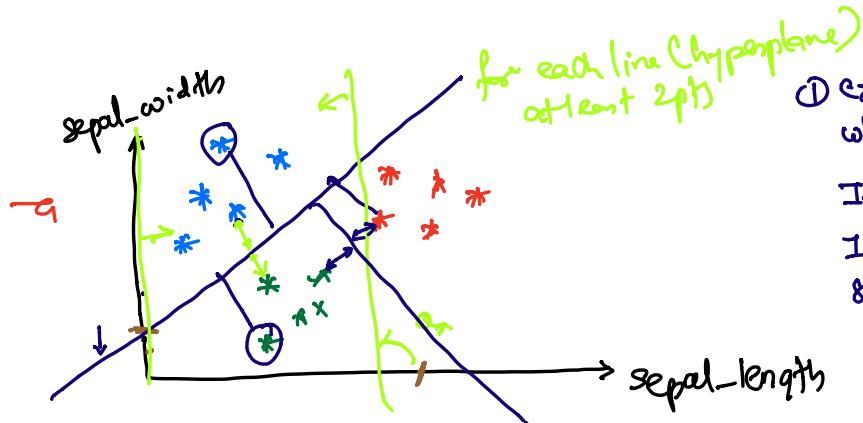
(smiley face) — data is separable
SVM will create the best model.

How to ensure my separator line (hyperplane) is placed in correct & generalized location?

$$x \Rightarrow ? \quad y \Rightarrow ? \quad \alpha_x \Rightarrow ? \quad \alpha_y \Rightarrow ?$$



Sum says
 Gamma parameters
 High → nearby pts from line
 Low → far pts from line



① Create one line where $\alpha = 0, y = 0$.
 Increment α till I receive atleast one separation

② change α_j such that the distance b/w two label pts is same

Feature Selection and Sampling Techniques

Questions we are trying to address:

- ① Is it mandatory to use all columns in the given dataset as feature or is there a way to figure out which feature column is the best for the given algo?
- ② Is the sample I am using for model creation, the best sample? KFold, Stratified Shuffle Split
- ③ What can be the maximum accuracy I can get from the given dataset? cross_val_score()

Linear Regression

Best features → $\xrightarrow{\text{Corr}()}$ OLS (Backward elimination.) ✓

what is
the best
accuracy I
can get? → → Cross validation for accuracy
score. (This technique is applicable for
all supervised learning algos)

Cross validation Technique for identifying best score for the given algorithm.

if you are using kfold,

Split depends on size.

Feature Selection

* Based on Statistics and not Domain.

Iris → sepal-length, sepal-width, petal-length, petal-width.

Best /
Worst

Statistics

Correlation matrix

Regressim

~~② OLS → Backward Elimination~~

R → UR, DT, RF, SVR

~~②~~ RFE → Recursive Feature Extraction Class → DT, RF

~~①~~ feature_importance } classification → DT RF

④ PCA → Principal Component Analysis } Reg. classf

✓ Principal Component Analysis

↓
Doesn't require any model at all. } Principal Component algo. → statistical algo