

Linearity

Supervised Learning

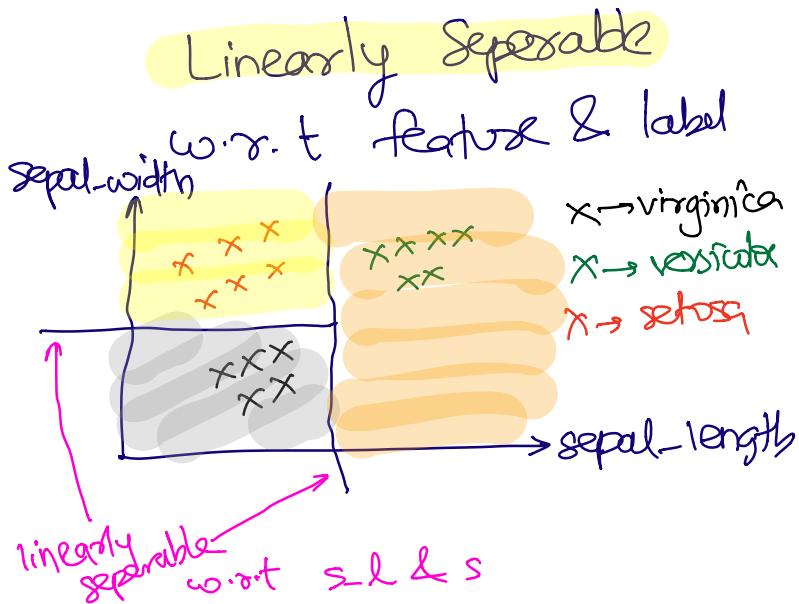
Regression

classification

is your data linearly separable or not.

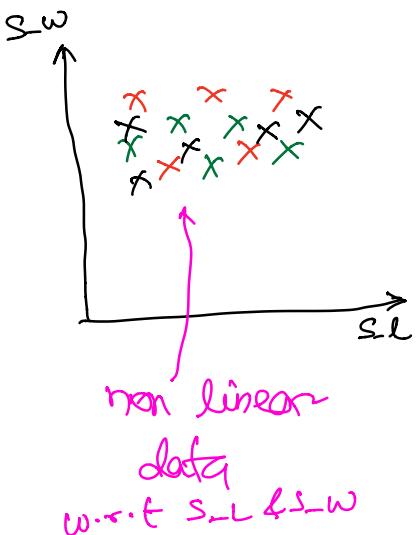
If there exists any kind of linear relationship corr()

Classification Problem



SVM

Non-linear data



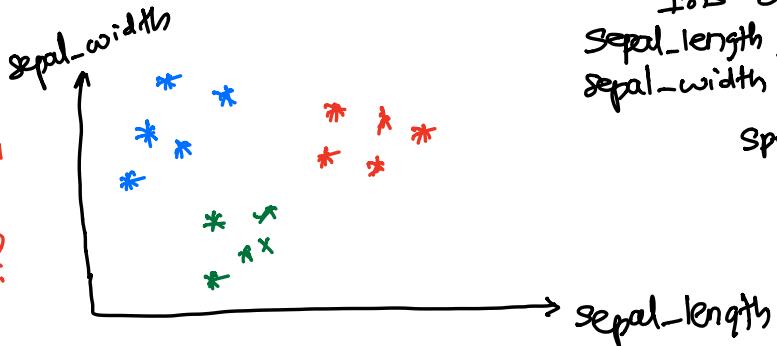
K-NN

Support Vector Machine

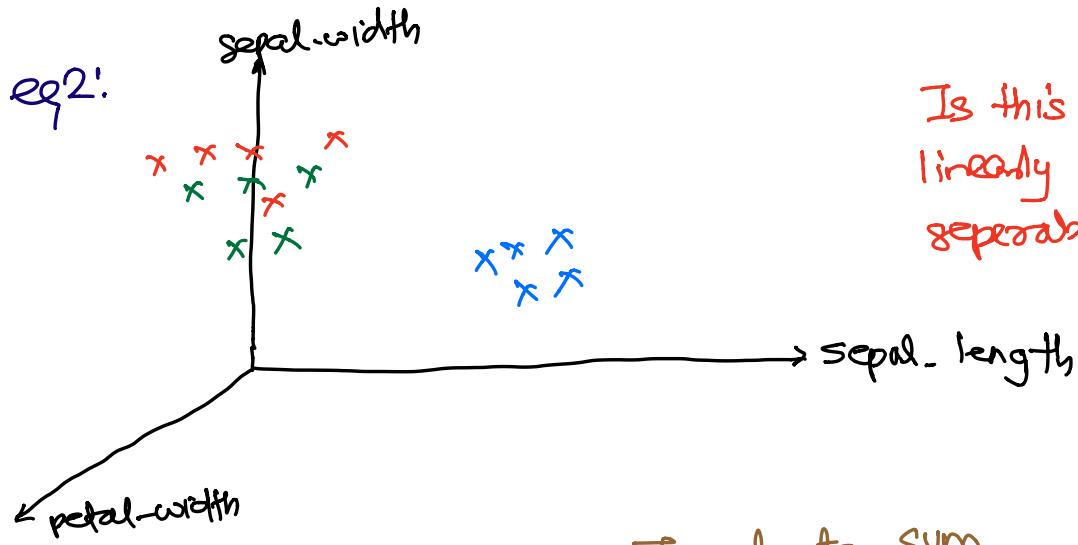
- Works best when you have more than 2 feature columns. (multi-dimensional space)

- SVM → Regression
→ Classification } All supervised methods.

eg 1:
Is this data linearly separable?



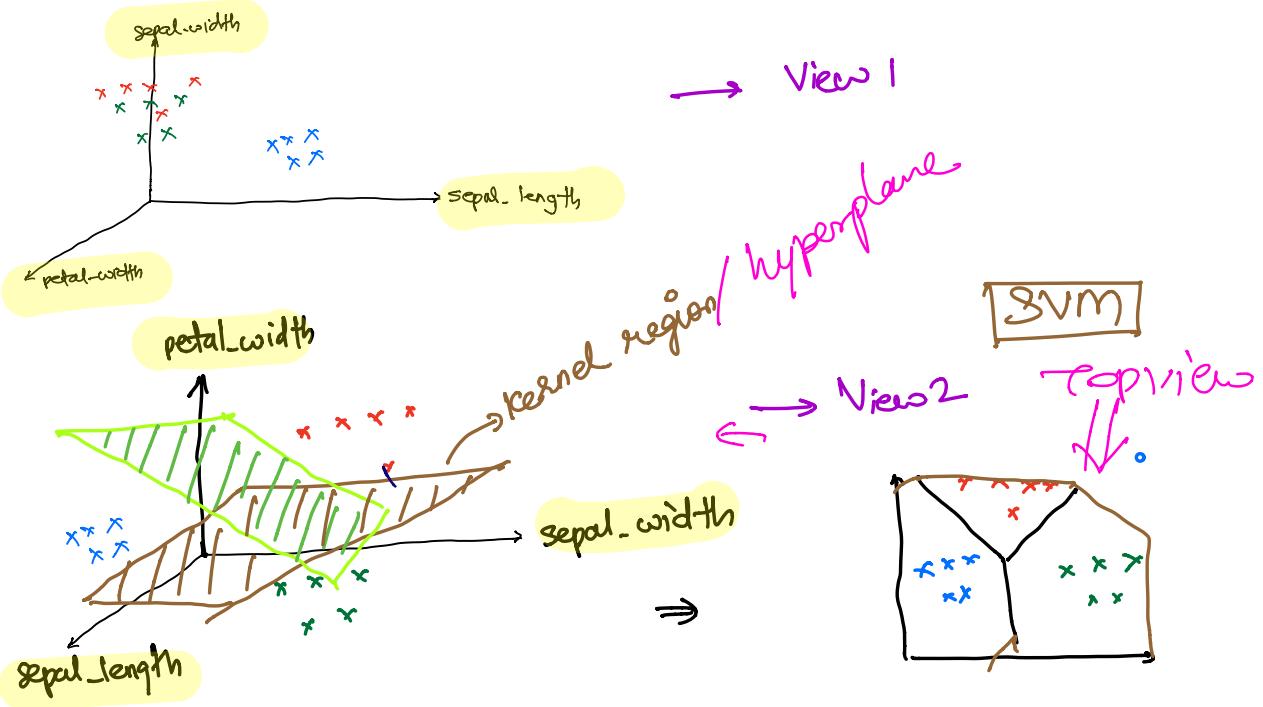
Iris dataset
sepal-length, petal-length
sepal-width, petal-width
species
→ Setosa ✕
→ Versicolor ✕
→ Virginica ✕

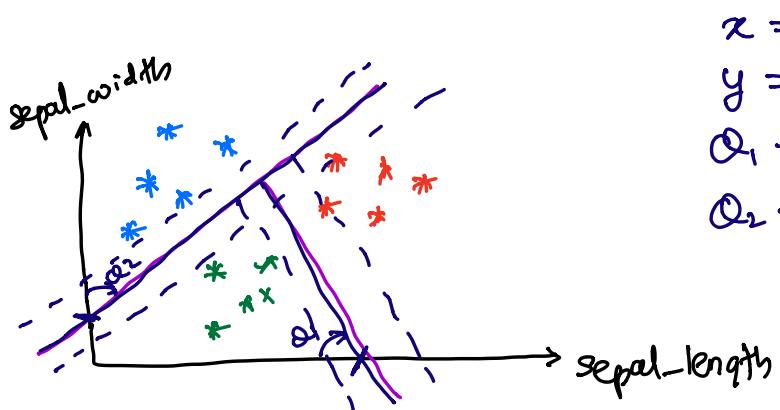


Is this data
linearly
separable?

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?

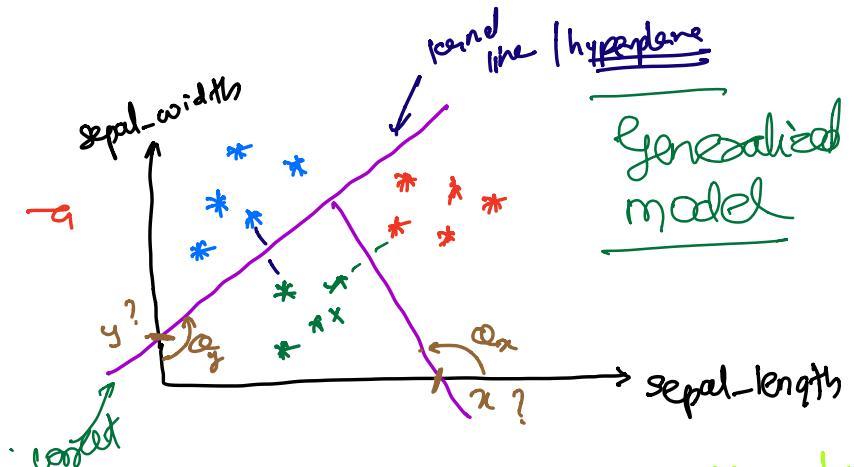
SVM

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

data is separable
SVM will create the best model.

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

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

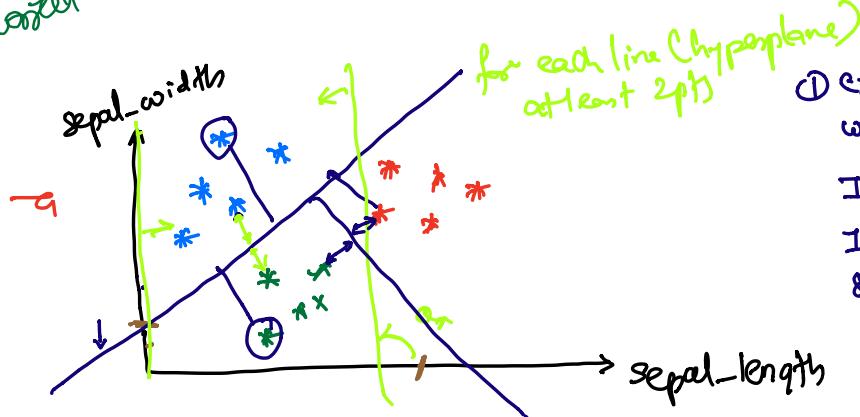


SVM says

Gamma parameters

High → nearby pts from line

Low → far pts from line



- ① Create one line
where $x=0, y=0$.
Increment x till
I receive atleast one
separation

leamma is high !

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 Classif. Regression

~~①~~ feature_importance } classification → DT RF

④ PCA → Principal Component Analysis } Reg. class

✓ Principal Component Analysis

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

K-fold Cross Validation Technique

dataset \rightarrow 100 records

