

Agenda :

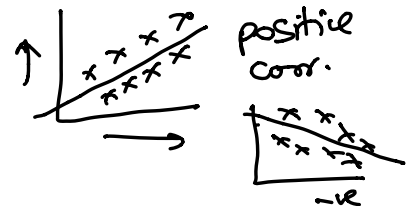
Support Vector machines

Linearity Behavior in a dataset



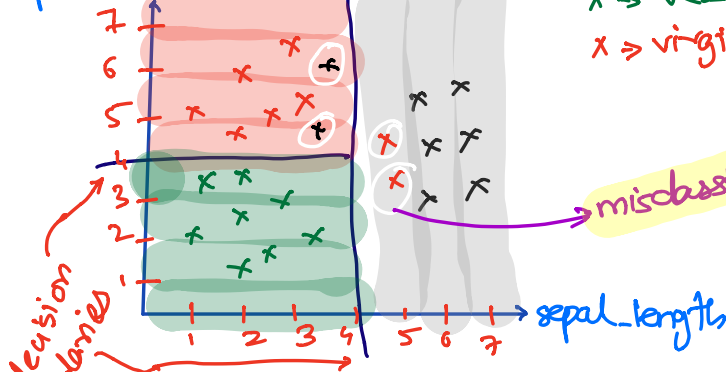
Is your dataset linearly separable or not

If there exists any kind of linear relationship in your dataset (Regression)  $\rightarrow \underline{\text{corr()}}$



classification Problem (iris.csv)

sepal\_width



x  $\rightarrow$  setosa (6)  
x  $\rightarrow$  versicolor (6)  
x  $\rightarrow$  virginica (6)

misclassified data.  
less than 7%

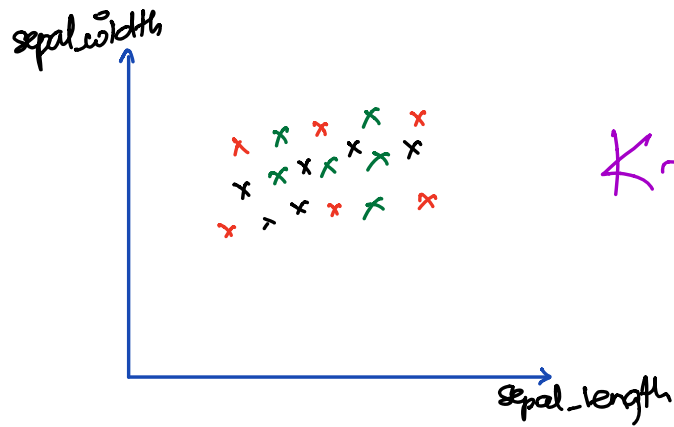
decision boundaries  
hyper plane  $\rightarrow$  SVM

Is this data linearly separable or not?

$x \leq 4$   
 $y \leq 4$   
 $\downarrow$   
versicolor

$x \leq 4$   
 $y > 4$   
 $\downarrow$   
virginica

$x > 4$   
 $\downarrow$   
setosa



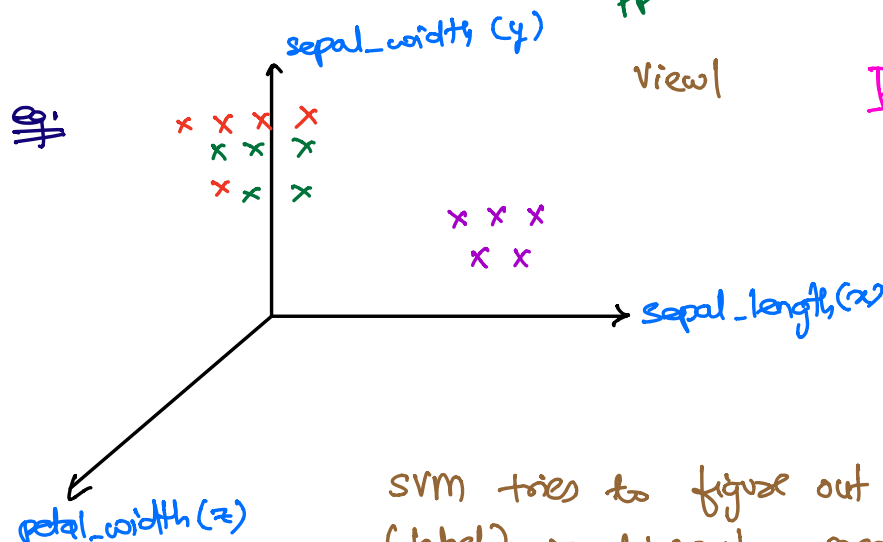
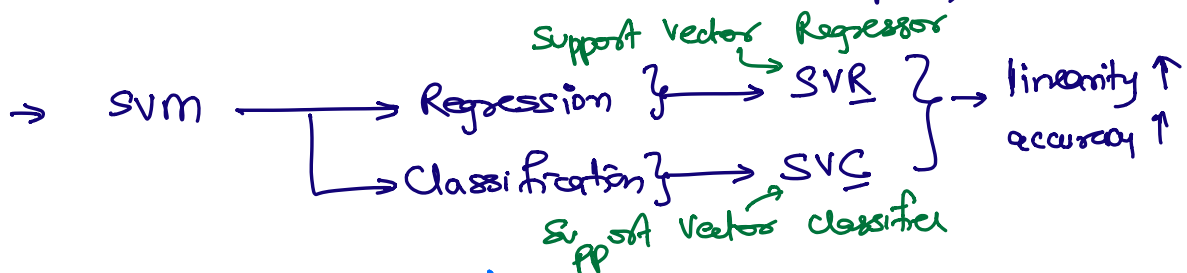
K-NN

Is this data linearly separable?

No-linear data expert to deal with (KNN)

What is Support Vector machine?

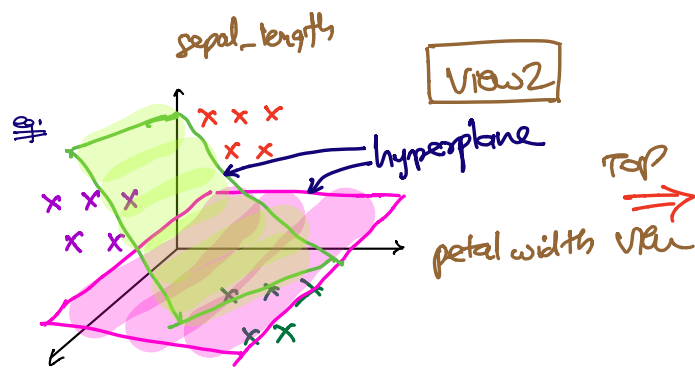
→ more than 2 features columns. (multi-dimensional space)



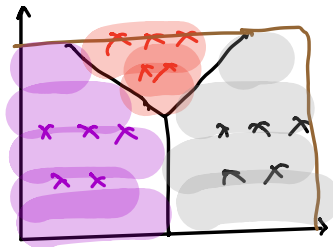
Is this data linearly separable or not?

↑  
answered by SVM

SVM tries to figure out whether the data (label) is linearly separable or not by checking all possible dimensions? (views)



Engineering Drawing  
Top view.



sepal width

hyperplane is used to separate regions

SVM tries to place the hyperplane in such that misclassification of data is very less.

Some of questions SVM tries to answer is,

① How to figure out whether my data is linearly separable or not?

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



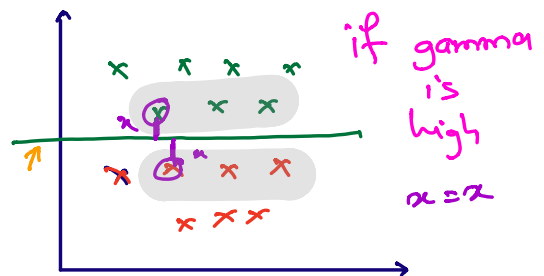
② How to ensure my hyperplane is placed in the correct & generalized location?

SVM → He has a parameter called gamma ( $\gamma$ )

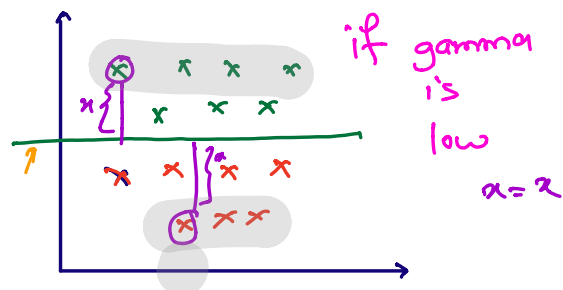
High (nearby pts)  
Low (far pts)

decision of hyperplane

placement is done using support vectors.



The points which is responsible to place the hyperplane at correct location is called **SUPPORT VECTOR**.



## Sampling Techniques.

↑  
how the data is picked and maintained in training & testing set.

Sampling Techniques helps us to understand whether the data has all possible combinations of Population.

Questions we are trying to address:

→ accuracy } score.  
→ generalization

① Is it mandatory to use all feature columns in the given dataset or can we remove features such that population & prediction both are not affected?

SelectByModel, RFE, ANOVA, OLS (regression)

② Is the sample I am using for model creation, the best sample?

→ kFold validation  
→ StratifiedShuffleSplit.

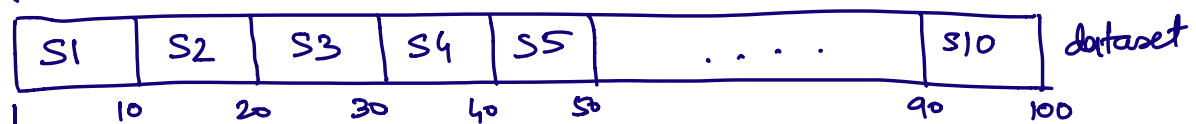
③ What can be the minimum accuracy I can get from the given best sample?

→ cross validation score.

### Cross validation

CV = 10 → dataset (100 records)      no of splits & no of iterations

(feature, label)



10 iterations:

1<sup>st</sup> iteration      S1 → testing      } → model → { score - ? score1  
                                 remaining → training

2<sup>nd</sup> iteration      S2 → testing      } → model → { score - ? score2  
                                 remaining → training

⋮

10<sup>th</sup> iteration      S10 → testing      } → model → { score - ? score10  
                                 remaining → training

Result:

[score1, score2, ..., score10]