

# Support Vector Machine

S M T W T F S

DATE: \_\_\_\_\_

- Support Vector Machine are one of the more complex algorithms we will learn, but it all begins with a simple premise:
  - Does a hyperplane exist that can effectively separate classes?
  - To first understand this question, we first need to go through the history and development of support vector machine!

## Section Overview

- History
- Intuition and Theory for SVM
- SVM Classification Examples
- SVM Regression Example
- SVM Project Exercise and Solution

## 2. History of Support Vector Machines

- SVM will be one of the most recently developed machine learning algorithms we will learn about!
- Let's briefly review the recent history behind the development of support vector machines.

Quick Notes:-

- 1960s: Vladimir Vapnik obtained his PhD in statistic at the institute of control sciences in Moscow.
  - worked at the Institute from 1961-1990.
- 1960s: Alex Chervonenkis and Vladimir Vapnik begin the development of support Vector Machines.
- 1963: published "Generalized portrait Algorithm" for image analysis by computers.
  - SVM in image analysis -
- 1964: Aizerman, Braverman, and Rozonoer, publish "Theoretical foundation of the potential Function Method" in pattern recognition learning.
  - o introduces geometrical interpretation of the Kernels as inner products in a feature spaces.
- 1974 - Vapnik and Chervonenkis continue work publishing "Theory of pattern Recognition".

M T W T F S

DATE: \_\_\_\_\_

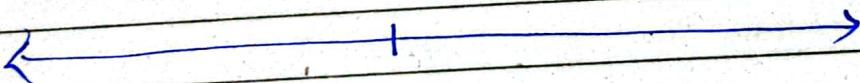
- 1960s - 1990s - Vapnik continues to work further develop SVM
  - 1992 - Bernhard Boser, Isabelle Guyon and Vladimir Vapnik suggested a way to create non-linear classifiers by applying the kernel trick to maximize margin hyperplane.
  - 1995: Corinna Cortes and Vladimir publish SVM incarnation utilizing a soft margin.
  - 1996 : Vladimir Vapnik, Harris Drucker, Christopher Burges, Linda Kaufman and Alexander Smola publish "Support Vector Regression Machines!", expanding SVM beyond classification tasks.
    - 30 years of work to get Support Vector machine we will be using today
- as we move
- We will see advance clearly from maximum margin classifiers, to support vector classifiers, and finally support vector Machines

### 3. SVM - Hyperplane and Margin

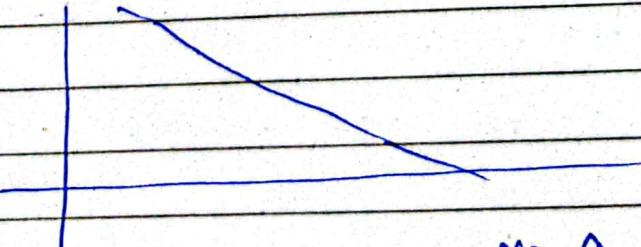
M I U V T F S

DATE: \_\_\_\_\_

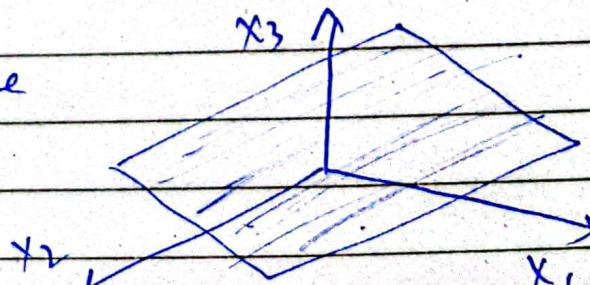
- We will slowly build up to SVMs:
  - Maximize Margin classifier
  - Support Vector classifier
  - Support Vector Machine
- Let's begin by understanding what a hyperplane is-
- In a N-dimensional space, a hyperplane is flat affine subspace of dimension  $N-1$
- For example:
  - 1-D Hyperplane is a single point
  - 2-D Hyperplane is a line
  - 3-D Hyperplane is a flat line.
- 1-D Hyperplane



- 2-D Hyperplane



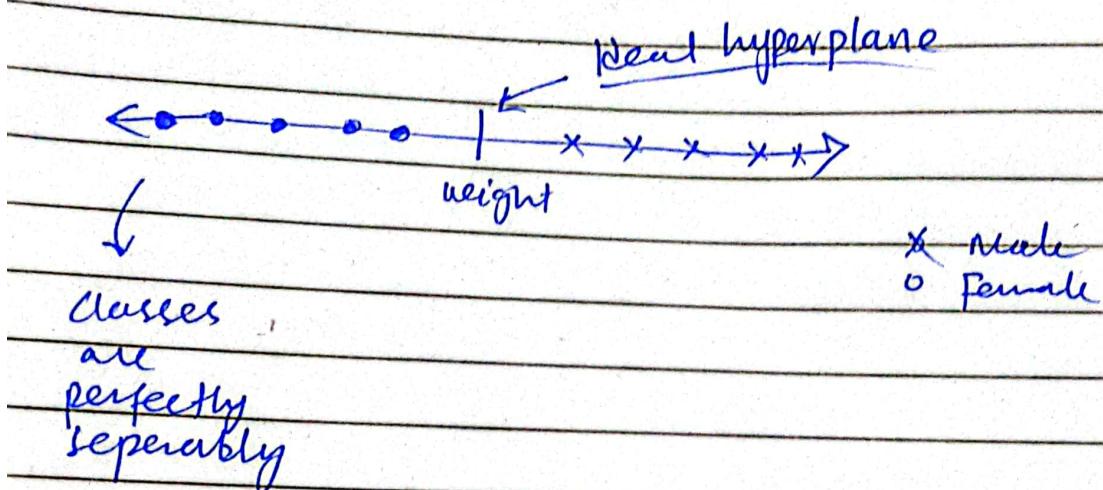
- 3-D Hyperplane



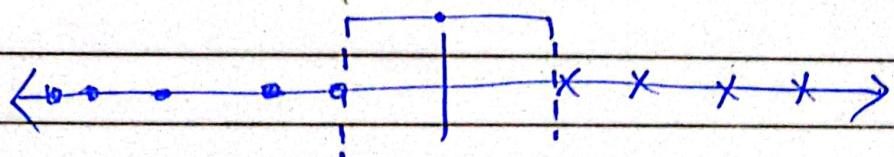
M T W T F S

DATE: \_\_\_\_\_

- The main idea behind SVM is that we can use hyperplane to create a separation between classes.
- Then new points will fall on one side of this separating hyperplane, which we can then use to assign a class.



- But for above case there could be many separation how to choose best fit line?
- We could choose separator that maximizes the margin between the classes.



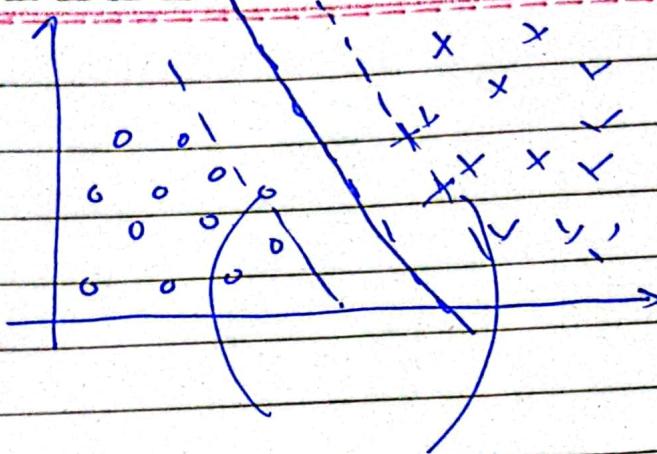
This is known as maximum marginal classifier.

- Same ideal is implemented for n-dimension -

- multiple possible margins  
choose the maximum margin

M I D W E E D F S

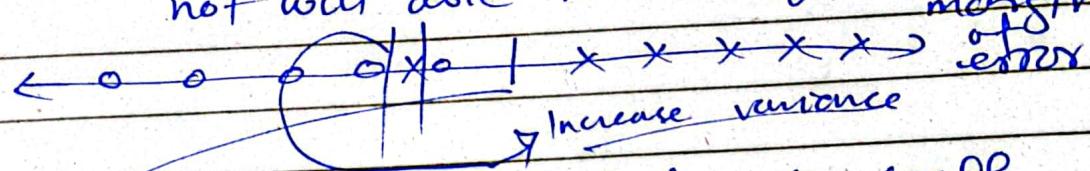
DATE: \_\_\_\_\_



Data points at the margin "support" separator.

- What happen if the classes are not perfectly separable?

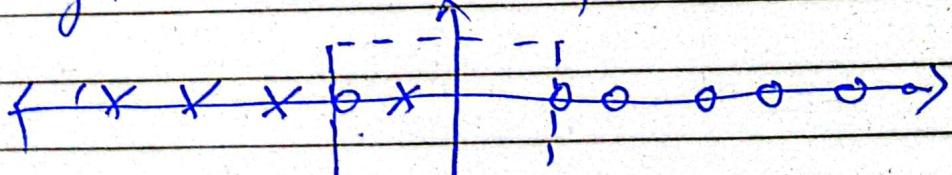
not will able to classify without margin



We will have a Bias-Variance tradeoff depending where we place this separator.

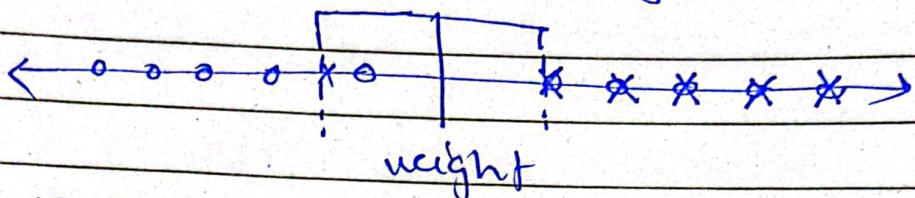
mean increase bias to lead to better long term result on future

- Here we allow more bias to lead to better long term result on future data.

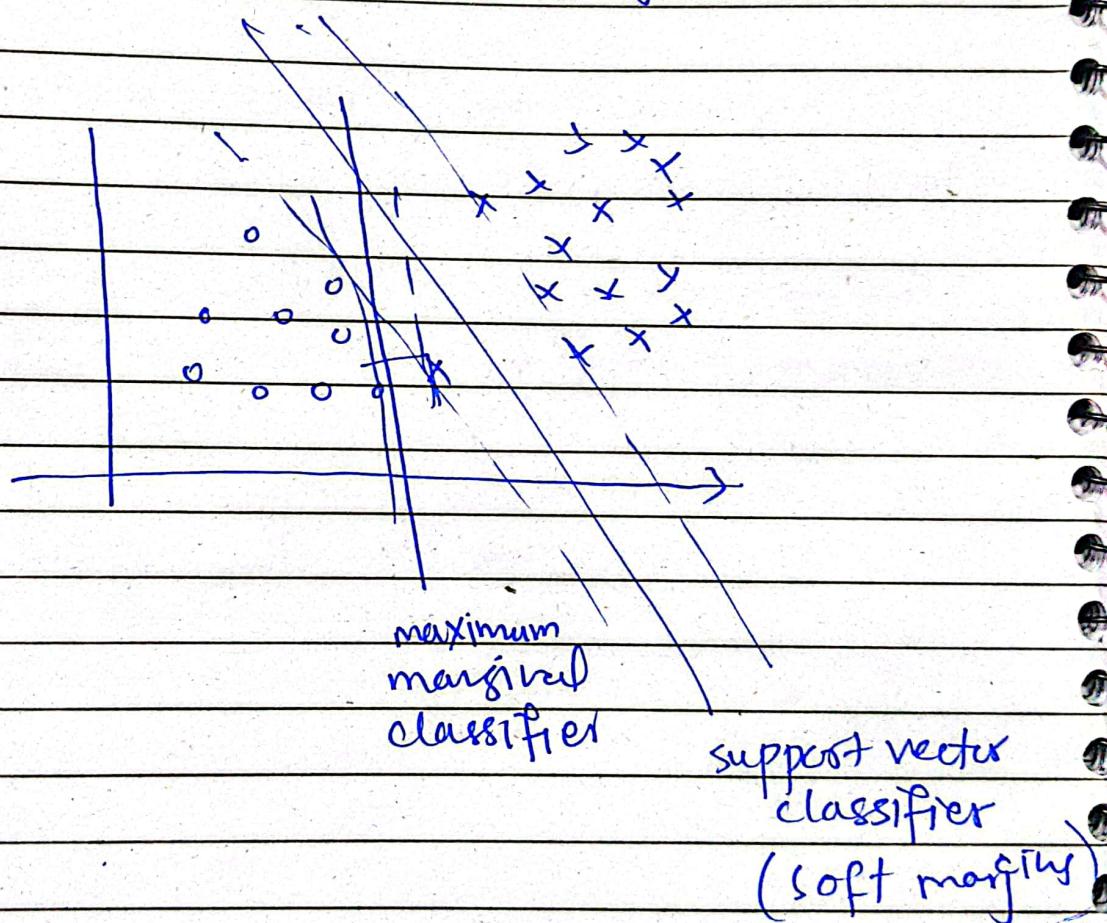


- Distance between threshold and the observation is a soft margin:

- Soft margin allows for misclassification inside the margin.
- There are many possible threshold splits if we allow for soft margins.



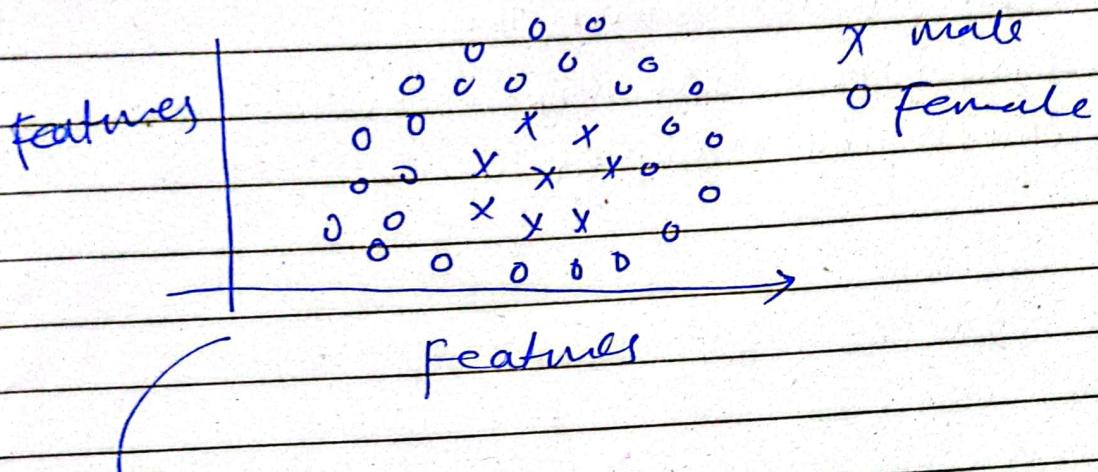
We can cross validation to determine the optimal size of the margins.



- We have only visualized cases where the classes are easily separated by the hyperplane in the original feature spaces.
- Allowing for some misclassifications still resulted in a reasonable results.

- What would happen in case where hyperplane performs poorly, even when allowing for missclassification?

~~X X X X O O O O X X X X~~



To resolve these cases, we move on from support vector classifiers, to support vector machines

- SVMs use kernels to project the data to a higher dimension, in order to use a hyperplane in this higher dimension to separate the data-

# 009 SVM - Kernels

M T W T F S

DATE: \_\_\_\_\_

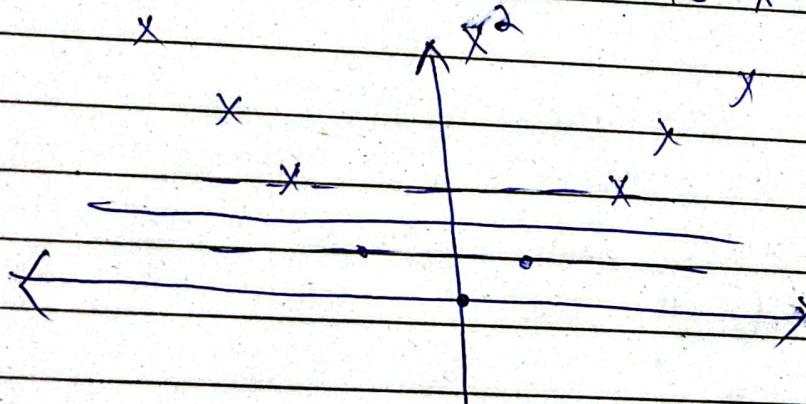
- Kernels allows us to move beyond a support vector classifier and use Support Vector Machines.
- There are a variety of kernels we can use to "project" the features to a higher dimension -
- Let's explore how this works through some visual examples

~~XX X 0 0 0 X X X~~

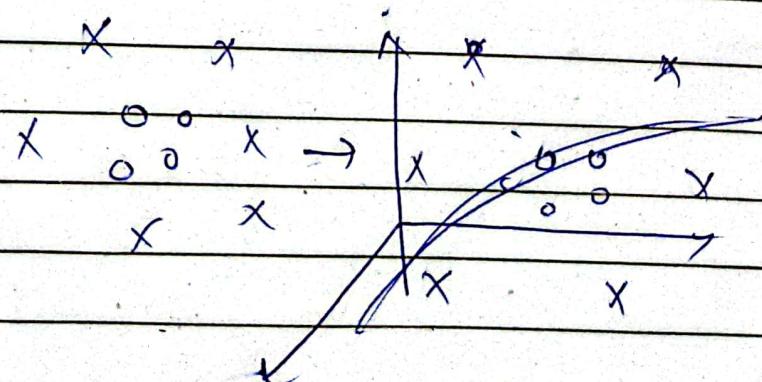
Using a polynomial kernel will expand onto  $x^2$  dimension

Default

No Default

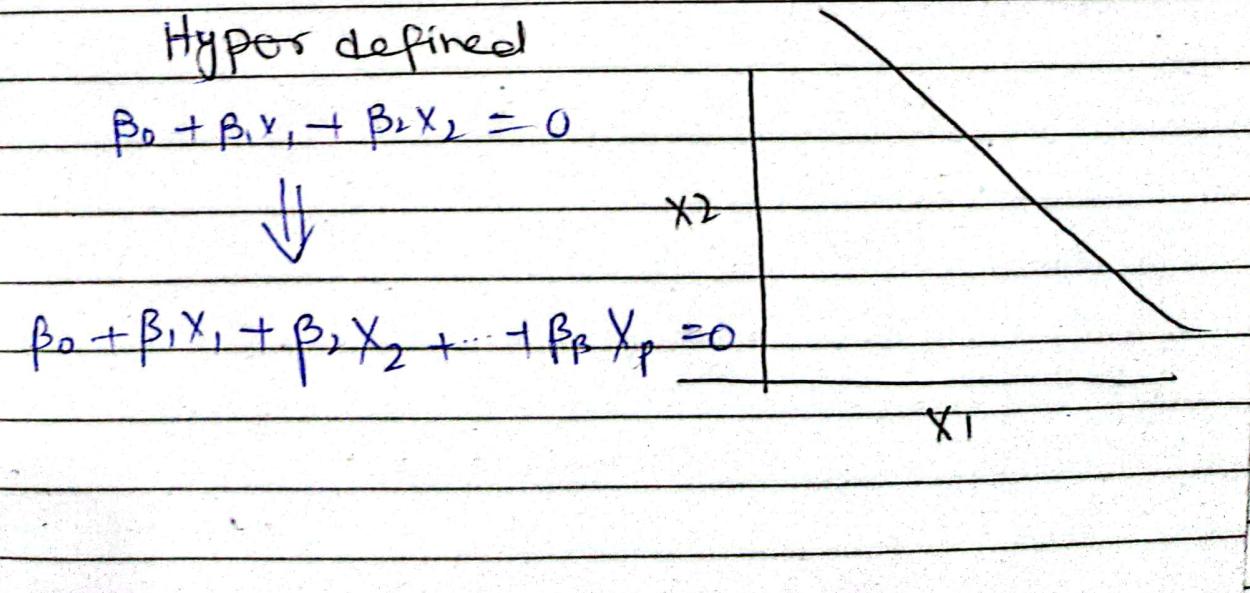


Imagine



- You may have heard of the use of kernels in SVM as the "kernel trick".
- We previously visualized transforming data points from one dimensional into higher dimension -
- Mathematically, the kernel trick actually avoids recomputing the points in a higher dimensional space!
- How does the kernel trick achieve this?
  - It takes advantages of dot products to achieve of the transposition of the data to achieve -
- In the next lecture we will go through the basic Mathematics idea behind the "kernel trick"!

## 5 SVM - Kernel Trick and Mathematics



00

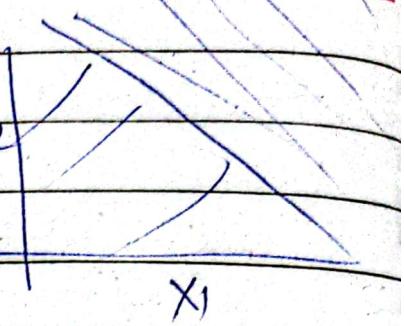
M T W T F S

DATE:

- Separating hyperplane

$$\beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_p x_p > 0$$

$$\beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_p x_p < 0$$



- Data points

$$x_1 = \begin{bmatrix} x_{11} \\ \vdots \\ x_{1p} \end{bmatrix}, \dots, x_n = \begin{bmatrix} x_{n1} \\ \vdots \\ x_{np} \end{bmatrix}$$

- Max Margin Classifier

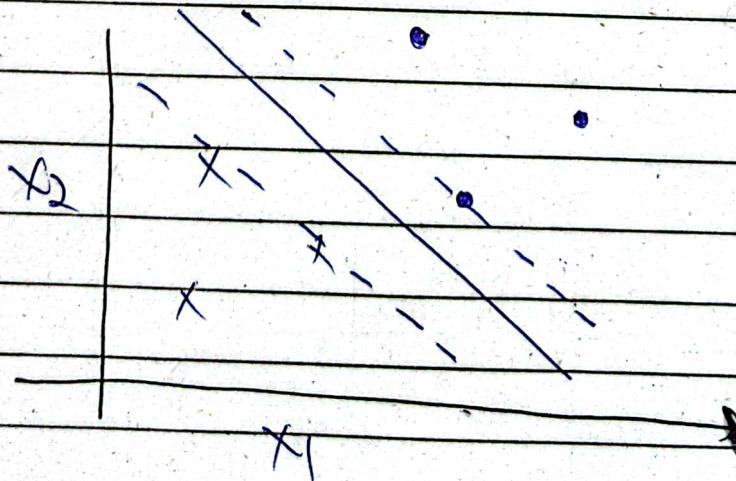
$$x_1 = \begin{bmatrix} x_{11} \\ \vdots \\ x_{1p} \end{bmatrix}, \dots, x_n = \begin{bmatrix} x_{n1} \\ \vdots \\ x_{np} \end{bmatrix}$$

maximize M

$$\beta_0, \beta_1, \beta_p, \dots, M$$

subject to  $\sum_{j=1}^p \beta_j^2 = 1$

$$y_i(\beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \dots + \beta_p x_{ip}) \geq M \quad \forall i = 1, \dots, n$$



• Support Vector Classifier

maximize . M

$$\beta_0, \beta_1, \beta_2, \dots, \beta_p, \epsilon_1, \dots, \epsilon_n, M$$

subject to  $\sum_{j=1}^p \beta_j^2 = 1$

$$\epsilon_i \geq 0, \sum_{i=1}^n \epsilon_i \leq C$$

$$y_i(\beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \dots + \beta_p x_{ip}) \geq M(1 - \epsilon_i)$$

• Note in Scikit-learn  
SVC!

C : float default = 1.0

Regularization parameter. The strength of the regularization is inversely proportional to C.  
Must be strictly positive. The penalty is squared L2 penalty.

• Support Vector Machine

$$x_1, x_2, x_3, \dots, x_p,$$



$$x_1, x_1^2, x_2, x_2^2, \dots, x_p, x_p^2$$

maximize M

$$\beta_0, \beta_1, \beta_2, \dots, \beta_p, \epsilon_1, \dots, \epsilon_n, M$$

Subject to  $y_i$

$$\left( \beta_0 + \sum_{j=1}^p \beta_j x_{ij} + \sum_{j=1}^p \beta_j x_{ij}^2 \right) \geq M(1 - \epsilon_i)$$

$$\sum_{i=1}^n \epsilon_i \leq C, \epsilon_i \geq 0, \sum_{j=1}^p \sum_{k=1}^2 \beta_{jk} = 1.$$

- How to deal with very large feature space?

- As a polynomial order grows larger, the number of computations necessary to solve for margin also grows!

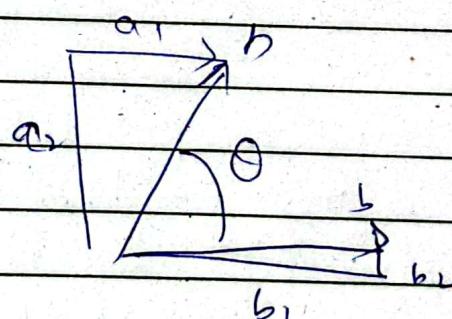
- The answer lies in the kernel trick which make use of the inner product of vectors, also known as the dot product.

- Dot product

$$\langle a, b \rangle = \sum_{i=1}^r a_i b_i$$

$$a \cdot b = a_1 b_1 + a_2 b_2$$

$$a \cdot b = |a| |b| \cos(\theta)$$



Notice that dot product can be thought of as a similarity between the vectors.

$$a \cdot b = |a| |b| \cos(\theta)$$

$$\cos(0) = 1$$

$$\cos(90) = 0$$

$$\cos(180) = -1$$

• Linear Support Vector Classified rewritten:

$$f(x) = \beta_0 + \sum_{i=1}^n \alpha_i \langle x, x_i \rangle$$

only non zero for support vector

calculate the inner product of all pairs of training observations

$$f(x) = \beta_0 + \sum_{i \in S} \alpha_i \langle x, x_i \rangle$$

$S$  collection of indices of these support points

• Kernel function :-

$$K(x_i, x_i') = \sum_{j=1}^p x_{ij} x_{i'j}$$

(A kernel is a function that quantifies the similarity of two observations.)

$$f(x) = \beta_0 + \sum_{i \in S} \alpha_i K(x, x_i)$$

$$\langle a, b \rangle = \sum_{i=1}^n a_i b_i$$

- polynomial kernel

$$K(x_i, x_i') = \left(1 + \sum_{j=1}^p x_{ij} x_{i'j}\right)^d. f(x) = \beta_0 + \sum_{i \in S} \alpha_i \langle x, x_i \rangle$$

$$\langle a, b \rangle = \sum_{i=1}^r a_i b_i$$

radial basis kernel

$$K(x_i, x_i') = \exp\left(-\gamma \sum_{j=1}^p (x_{ij} - x_{i'j})^2\right)$$

11

- The use of kernels as a replacement is known as the kernel trick
- Kernel allow us to avoid computations in the enlarged features spaces by only needing to perform computation for each distinct pair of training points.

# SVM - Classification with Scikit-learn

## part-one

NAME: \_\_\_\_\_

DATE: \_\_\_\_\_

# Code

```
df = df pd.read_csv("mouse-mal-study.csv")  
df.head()
```

```
sns.scatterplot(x="Med-1-ml", y="Med-2-ml",  
hue="Virus present", data=df)
```

```
x = np.linspace(0, 10, 100)
```

```
m = -1
```

```
b = 11
```

$$y = m \cdot x + b$$

```
plt.plot(x, y, 'black')
```

```
from sklearn.svm import SVC  
help(SVC)
```

```
y = df['Virus present']
```

```
x = df.drop(['Virus present'], axis=1)
```

```
model = SVC(kernel="linear", C=1000)
```

```
model.fit(x, y)
```

```
from SVM.mmargin_plot import plot_SVM_boundary
```

```
plot_SVM_boundary(model, x, y)
```

# SVM - Classification with scikit-learn

## part-2

M T W T F S

DATE: \_\_\_\_\_

model = SVC(kernel='rbf', C=1)

model.fit(X, y)

plot\_svm\_boundary(model, X, y)

Check for

C = 0,

gamma : {'scale', 'auto'} or float

for rbf, poly, sigmoid

if "gamma='scale'

$1/n\text{-features} \times \text{variance}$ ) as the  
value of gamma

If:  
auto

1

n-features

or we can give any float value.

model = SVC(kernel='sigmoid')

model.fit(X, y)

plot\_svm\_boundary(model, X, y)

model = SVC(kernel='poly', C=1, degree=1)

model.fit(X, y)

plot\_svm\_boundary(model, X, y)

from sklearn.model\_selection import GridSearchCV.

Svm = SVC()

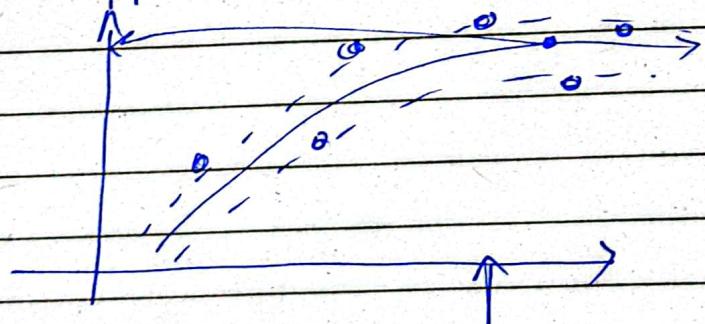
param\_grids = { 'C': [0.01, 0.1, 1, 10],  
'kernel': ['linear', 'rbf'] }

grid = GridSearchCV(Svm, param\_grids)  
grid.fit(X, y)

grid.best\_params

## 8 SVM with scikit-learn Regression Task

### o Support Vector Regression



### - Concrete Slump test :-

It measure the consistency ~~before~~ it ~~actually~~ sets of fresh concrete before it actually sets and its to ~~check~~ perform to check the workability of fresh freshly made concrete and therefore the ease

with which concrete flows.

also  $\text{b}$  can be used as an indicator of improperly mixed match - 



# code:

```
df = pd.read_csv("cement_slump.csv")
df.head()
plt.figure(figsize=(10,8))
sns.heatmap(df.corr(), annot=True)
```

`df.columns`

$$\underline{X = \alpha f}.$$

```
x = df("    ", axis=1)  
y = df("    ")
```

~~$X_{train}, t_{train}, X_{test}, y_{train}, t_{test} = train\_test\_split(X, y, test\_size=0.3)$~~

from sklearn.preprocessing

```
scaler = StandardScaler()
```

Scaled X train = 11 11

Scalcd-X-test = 1 1

```
from sklearn.svm import SVR, LinearSVR
```

MODULES

help(SVR)

SVR(\*, kernel, epsilon = 0.1)

gamma =

C = }

if zero it can  
be considered as  
overfit

base\_model = SVR( )

base\_model.fit(scaled\_X-train, y-train)

based-pred = base-model.predict(scaled-X-test)

mean\_absolute\_error(y-test, y-pred)

5.23690

mean = np.sqrt(mean\_squared\_error( ))

6.1542

param-grid = { 'C': [0.001, 0.01, 0.1, 0.5, 0.1],  
'kernel': ['linear', 'rbf', 'poly'],  
'gamma': ['scale', 'auto'],  
'degree': [2, 3, 4],  
'epsilon': [0, 0.01, 0.5, 1, 2] }

willy  
allow  
error

from sklearn.model\_selection import GridSearchCV

svr = SVR( )

grid = GridSearchCV(svr, param-grid)

grid.fit(scaled\_X-train, y-train)

M T W T F S

DATE: \_\_\_\_\_

grid.best\_params -

grid-preds = grid-predict(Scaled-X-test)

mean\_absolute\_error ( )

2.5128

np.sqrt( )

3.178