

## (2) Polynomial Regression

- It is the form of Linear Regression in which the relationship b/w the independent variable ( $x$ ) & dependent variable ( $y$ ) is modeled as an  $n^{\text{th}}$  degree polynomial.
- Polynomial regression fits a nonlinear relationship b/w the value of  $n$  & the corresponding conditional mean of  $y$ , denoted by  $E(y/x)$

### Why Poly-Regr?

- For curvilinear hypothesis relationships.
- Inspection of Residuals (y axis), if we try to fit a linear model to curved data, a scatterplot of residuals (y axis) on the predictor (x axis) will have patches of many free residuals in the middle. Hence, it is not appropriate.
- All Indp. variables being Indpt. is not satisfied in Poly Regs.

### Uses:-

- to know,
- Growth Rate of tissues
- Progression of disease epidemics
- Distribution of Carbon Isotopes in lake sediments.

$$y = \beta_0 + \beta_1 x + \beta_2 x^2 + \dots + \beta_m x^m$$

$$\Rightarrow y = \beta_0 + \sum_{i=1}^m \beta_i x_i + (\epsilon_p) \rightarrow \text{Polynomial function.}$$

### Advantage (using degree of polynomial inc.)

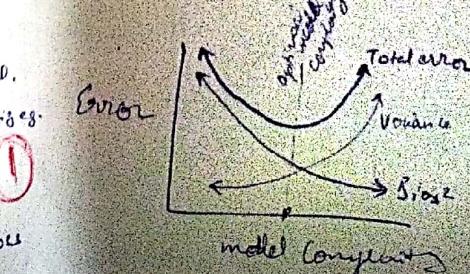
- Broad range of function can be fit under it
- Polynomial basically fits wide range of curvature.
- Provides best approximation of Relationship b/w dependent and independent variable.

- High Bias indicates underfit  
To over come under fitting, we need to decrease complexity of the model i.e.  $\text{RMSE} \propto R^2 \uparrow$

- degree is  $\uparrow$  then error  $\uparrow$

### Disadvantage

- Too sensitive to outliers
- presence of 1 or more outliers can affect the results.
- few availability of tools for detection of outliers in nonlinear Regs than that for linear Regs.
- Chances of overfitting
- High variance means high degree of polynomial.  
 $\Rightarrow$  training error is low for small no. of training.  
is training + error  $\downarrow$
- For high variance & no. of samples in the training set size will decide the gap b/w Cross Validation & training error



## Polynomial Regression # 5 minute Egg.

- Here linear model is used on Non-linear data set.
- Relationship b/w dependent & independent variable is non linear in nature.
- Simple linear regg. eqn we know,

$$y = \beta_0 + \beta_1 x_1 \quad \leftarrow \text{linear regg.}$$

$$\rightarrow y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_3 + \beta_4 x_4 \quad \leftarrow \text{(it is 1 degree polynomial)}$$

$$\Rightarrow y = \beta_0 + \sum_{i=1}^m \beta_i x_i \quad \text{--- (i)}$$

- Check Independent's power, we'll get to know degree of polynomial

- 0 degree polynomial:  $y = \text{constant}$

1 degree polynomial:  $y = mx + c$

2 degree polynomial:  $y = ax^2 + bx + c$

$$y = \beta_0 + \beta_1 x + \beta_2 x^2 + \beta_3 x^3 + \dots + \beta_n x^n$$

$$\Rightarrow y = \beta_0 + \sum_{i=1}^m \beta_i x_i^i + F_p$$

Polynomial function, it tries to add variables

we know, for the below,  
we can't apply  
linear models  
 $\therefore$  if we want to apply  
linear model in

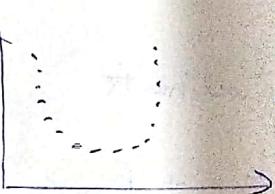
this  $\rightarrow$  i.e. (i) we normally  
had this as a linear/multiple regg.

Standard eqn: i.e.  $[x_1, x_2]$ , now

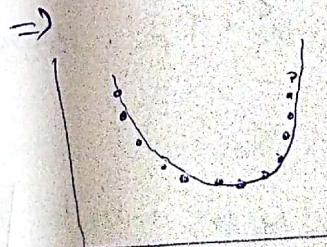
for (ii) we see this  $F_p$ ,  $F_p$  here will  
inc. the power of variable i.e.  
 $[x_1, x_2, x_1^2, x_2^2]$ , this  $F_p$  is the

polynomial function, jisne add or  
power ko inc karne ka raam liya hai.

$F_p$   $\rightarrow$  Kitne Variable chahiye kyun ke variable ka degree  
Kitni chahiye tak linear model ko hum  
use kar ske for non linear data set.



isme kuch 2 degree ki zarurat hai :/  
2 degree polynomial  $\rightarrow$  hota hai; e  
parabolic function  
k form hota hai



## ~~QUESTION~~ SVR

(16)

### SVM:-

- Stands for Support Vector machine, is a classifier.
- Classifiers perform classification, predicting discrete categorical labels.

### SVR

- Stands for Support Vector Regression, is a Regressor.
- Regressor performs regression, predicting continuous ordered variable.

SVR & SVM use very similar algorithms, but predict different type of variables.

In Simple regres'ion we try to minimize the error rate. (pred - actual)

In SVR — fit the error within a certain threshold. (boundary)

### Important terms in SVR :-

Kernel: function used to map a lower dimensional data into a higher dimensional data.

Hyperplane: In SVM, it is the separation line b/w the data classes.

In SVR, it is the line that will help us predict the continuous target value.

### Boundary line:-

In SVM there are two lines other than hyperplane which creates a margin.

The support vector can be on the boundary line or outside it.

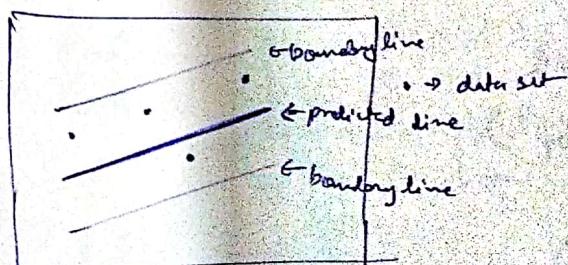
The boundary line separates the two classes.

In SVR it is same concept.

### Support Vectors :-

These are the data points which are closest to the Boundary.

The dist. of the pts is minimum or least.



## Random Forest

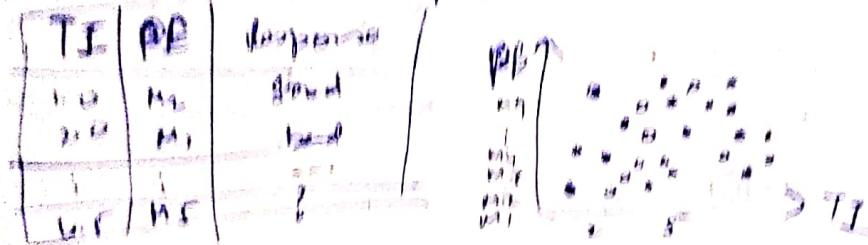
- It is an ensemble classifier that consists of many decision trees which make predictions and then take average to output final prediction.
- It combines Performance (Accuracy) like the generalization of classifier.

## Decision tree

- There are individual learners that are combined, then commonly used
  - One type of decision tree is called Classification & Regression Tree
- CART involves finding space best split of independent variables.  
→葉子 should be pure w.r.t. response variable.  
→ Example model is fit for each response variable rate for classification, concentration of response.

### Decision tree dimensions greatly, recursive partitioning

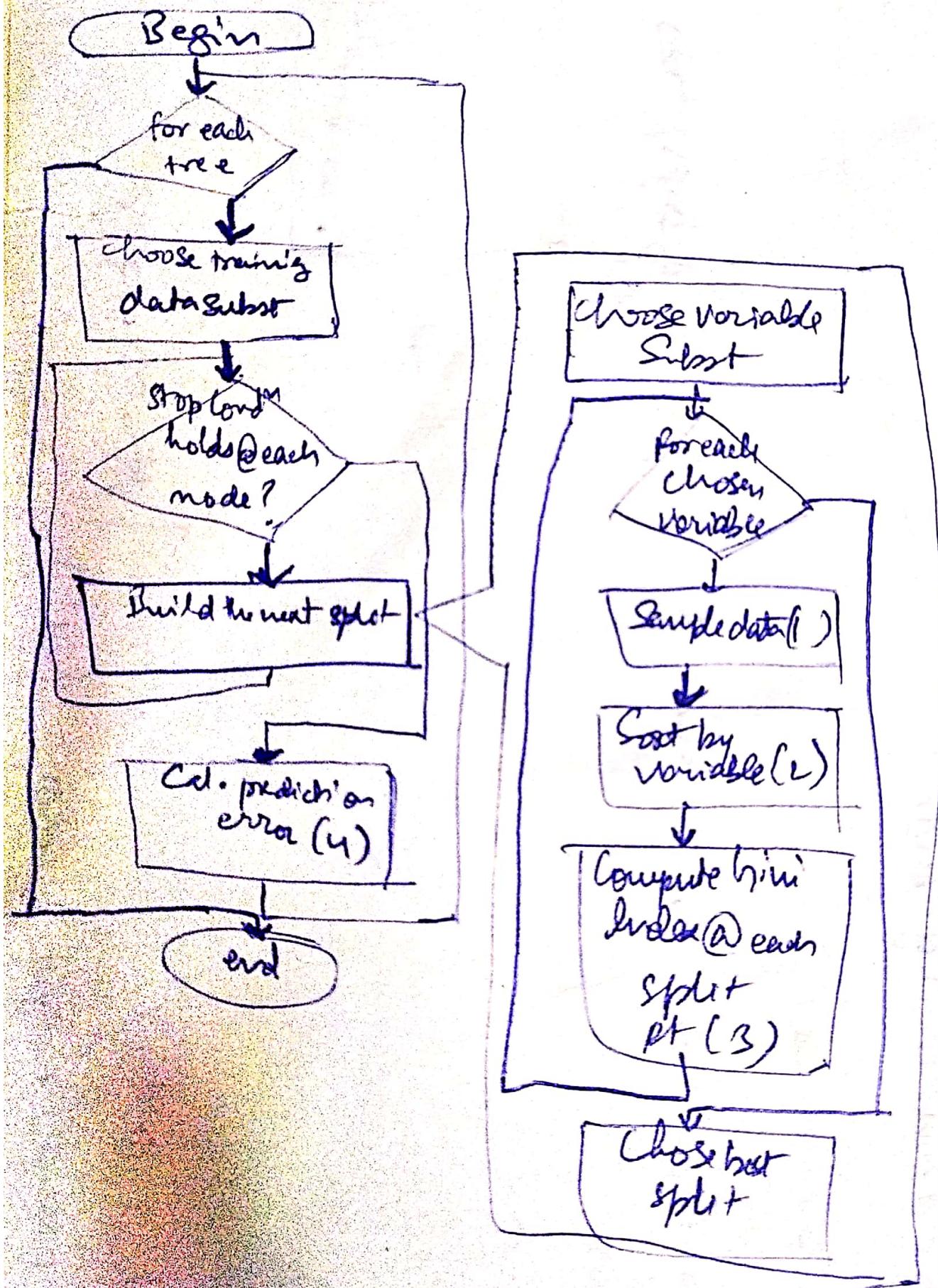
- Simple dataset with four predictors



- Grows, recursive partitioning along TI & PB



# Algorithm Flowchart



→ Points are proportional to purity measure (e.g. 29.0% (negative), 61.0% (positive))

→ N is fixed

Build trees until the error no longer dec

→ M is selected

Try to get around defaults, half of them & twiddle them & pick the best

## Random Forest

(1)

- ① Most Accurate learning algorithms available ; produces highly accurate classifier.

- ② runs efficiently on large databases.
- ③ Can handle thousands of input variable without variable selection.
- ④ Gives estimation of which variables are important in the ~~the~~ <sup>selection</sup>.
- ⑤ Generates an internal unbiased estimate of the generalization error as the forest building progresses.
- ⑥ It has effective method for estimating missing data & maintains accuracy when a large proportion of data are ~~decided~~ missing.
- ⑦ Has methods for balancing errors in class popn unbalanced data sets.
- ⑧ Generated forest can be saved for future use on other data.
- ⑨ It offers an experimental method for detecting variable interactions.

### Disadvantages

- ① It has been observed to overfit for some datasets with ~~no~~ noisy classification regression tasks.
- ② For data including categorical variables with diff. no. of levels, random forests are biased in favor of those attributes with more levels.

### Random Forest estimating the test error

- While growing forest, estimate test error from training samples.
- For each tree grown (33% to 36%) samples are not selected in Bootstrap called out of Bootstrap (OOB) samples.
- Using oob samples as input to the corresponding tree, predictions are made as if they were novel test samples.
- Through book keeping majority vote, average reg., is computed for all oob samples from all trees.

estimating the importance of each predictor.

- Denoted by  $\hat{e}$  the oob estimate of the loss when using original training set,  $D$ .
- For each predictor  $x_p$  where  $p \in \{1, \dots, k\}$ 
  - Randomly permute  $p^{\text{th}}$  predictor to generate a new set of samples  $D' = \{(y'_1, x'_1), \dots, (y'_N, x'_N)\}$
  - Compute oob estimate  $\hat{e}'$  of prediction error with the new samples.
- A measure of importance of predictor  $x_p$  is  $\hat{e}' - \hat{e}$ , the increase in error due to random perturbation of  $p^{\text{th}}$  predictor.