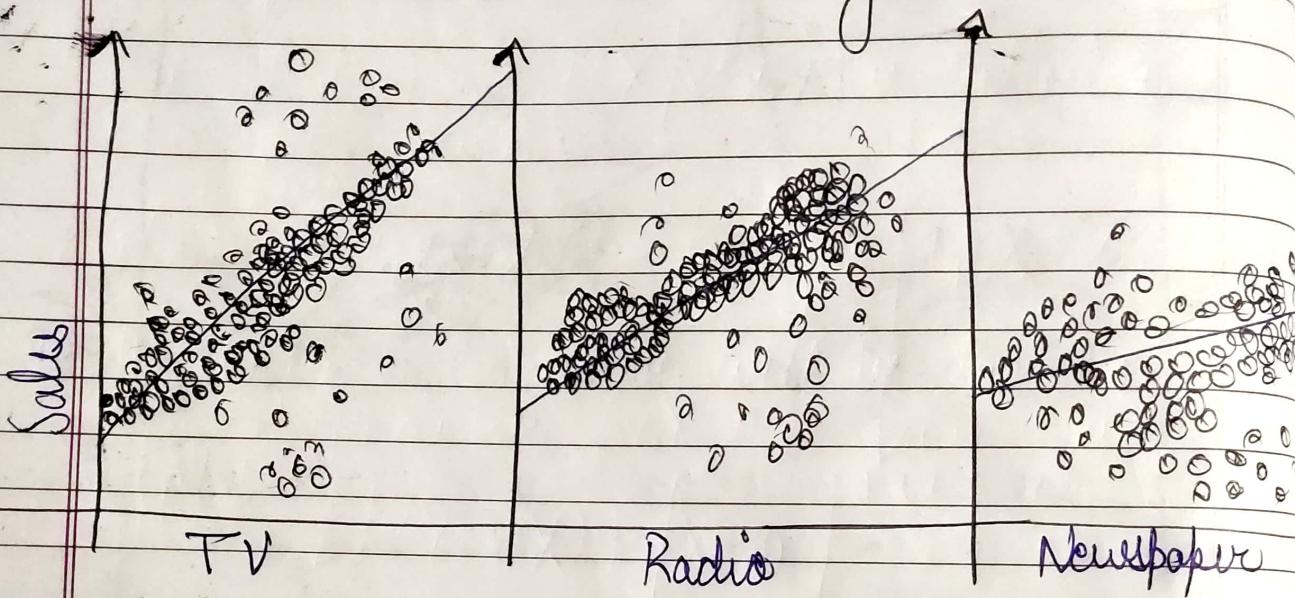


$$h \Rightarrow \lambda$$

## What is Statistical learning



Showers above is the Sales vs TV, Radio & Newspaper with a blue linear regression line fit separately to each

But can we predict Sales using all of these three?

The answer is Yes. The statistical learning helps you do so.

We here like to know how these work together with each other i.e.

$$\text{Sales} \approx f(\text{TV}, \text{Radio}, \text{Newspaper})$$

Here, Sales is a response / Target or Y.

And Input Variable / Features  $\{x_1, x_2, x_3\}$

where,  $\begin{pmatrix} X_1 \\ X_2 \\ X_3 \end{pmatrix}$  represents TV.  
 $X_2$  represents Radio.  
 $X_3$  represents Newspaper.

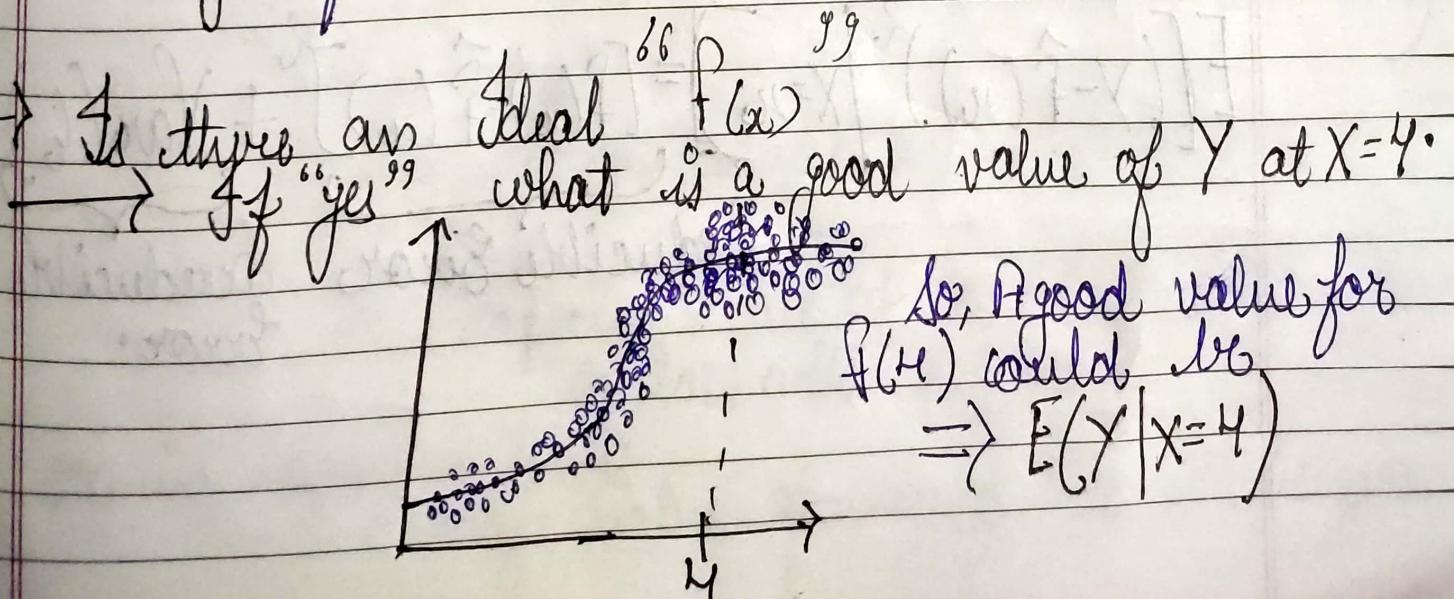
So input Vector  $X = \begin{pmatrix} X_1 \\ X_2 \\ X_3 \end{pmatrix}$

Now, we can write  

$$Y = f(x) + \epsilon$$
  
 where  $f$  is an irreducible error / measurement error

→ So, with a good "f"

- we can do good predictions
- we can understand which input feature are important.
- we can also understand how each component  $x_j$  of  $X$  impacts  $Y$ .



$f(x) = E(Y|X=x)$  what does this specify?

This specifies  $E(Y|X=x)$  means expected value / weighted average of value  $\underline{Y}$  at  $\underline{X=x}$ .

\* This ideal  $f(x) = E(Y|X=x)$  is called a regression function.

$$\rightarrow f(x) = f(x_1, x_2, x_3) = E(Y|X_1=x_1, X_2=x_2, X_3=x_3)$$

$\rightarrow$  The  $f(x)$  presented above will be the ideal/best optimal predictors for  $Y$  that minimizes

$$MSE = E[(Y - f(x))^2 | X=x] \quad ***$$

$\rightarrow \epsilon = Y - f(x)$  is the irreducible error - i.e. even if we still make errors in predictions. Because even at  $X=x$  there is a cloud/distribution of  $Y$  values.

$\rightarrow$  We represent our estimate by  $\hat{f}(x)$  where ideal estimate is  $f(x)$

$$E[(Y - \hat{f}(x))^2 | X=x] = [f(x) - \hat{f}(x)]^2 + \text{Var}(\epsilon)$$

$\underbrace{\text{Reducible Error}}$        $\underbrace{\text{Irreducible Error}}$

\* \* \* This topic is irrelevant for Non stats Background \* \* \*

Proof:  $E[(Y - \hat{f}(x))^2 | X=x] = (\hat{f}(x) - \hat{f}(x))^2 + \text{Var}(\epsilon)$

add & subtract  $f(x)$  to the LHS.

$$\begin{aligned} & E[(Y + f(x) - f(x) - \hat{f}(x))^2 | X=x] \\ &= E[(Y - E[Y|X])^2] + (E[Y|x] - \hat{f}(x))^2 \\ &\quad + 2E[(Y - E[Y|X])(E[Y|x] - \hat{f}(x))] \end{aligned}$$

We will let the first 2 terms stay, & deal with 3rd one.

In 3<sup>rd</sup> term we can apply property of conditional expectation that states

$$\text{So, } E(Y) = E[E(Y|X=x)]$$

$$\begin{aligned} & E[(Y - E[Y|X])(E[Y|x] - \hat{f}(x))] \\ &= E[(E[Y - E[Y|X]]|X)] (E[Y|x] - \hat{f}(x)) \end{aligned}$$

This will be reduced to zero

as, Expected Value of Expected Value of loss will be zero.

or in layman terms average value of average loss will be zero

so the <sup>final</sup>  
error term,

$$E[(Y - E[X|x])^2] + E[(E[X|x] - \hat{f}(x))^2]$$

Reducible error  
Irreducible error

This will the Variance  
across of  $X$  at a given  $x$   
at all  $x$ 's.

$$\therefore E(\epsilon) = \text{Var}(\epsilon)$$

because mean is 0

$$\& \text{Var}(\epsilon) = E(\epsilon^2) + (E(\epsilon))^2$$

This is 0

$$= E((f(x) - \hat{f}(x))^2)$$
$$= (f(x) - \hat{f}(x))^2$$

because this is linear

\* hint take 2 linear equations  
& list

$$\therefore E[(Y - \hat{f}(x))^2 | X=x]$$

$$= \text{Var} \epsilon + (f(x) - \hat{f}(x))^2$$

XXX

XX X

Typically we've a few values of  $y$  at  $x=4$

so, we can't find exact predictor of  $Y$  but we can predict  $y$  with good accuracy

- Typically we've a few values at  $x=4$
- So we can't compute  $E(Y|X=x)$

So, what can we do.

$$\hat{P}(x) = \text{Ave}(Y | x \in N(x))$$

where  $N(x)$  is neighborhood of  $x$ .

Nearest Neighbor can be pretty good when

- $p \leq 4$  (less predictors)
- $n \gg p$  (large  $n$ ), so that we've a high density neighborhood.
- it might not work good when  $p$  is large.

Reason: Curse of Dimensionality. Nearest Neighbor tend to be far in high dimensions.

→ We need to get a reasonable fraction of  $N$  values of  $y_i$  to bring the variance down. ex: 10%.

→ If 10% neighborhood in high dimensions will no longer be local, hence contradicting  $E(Y|X=x)$ .

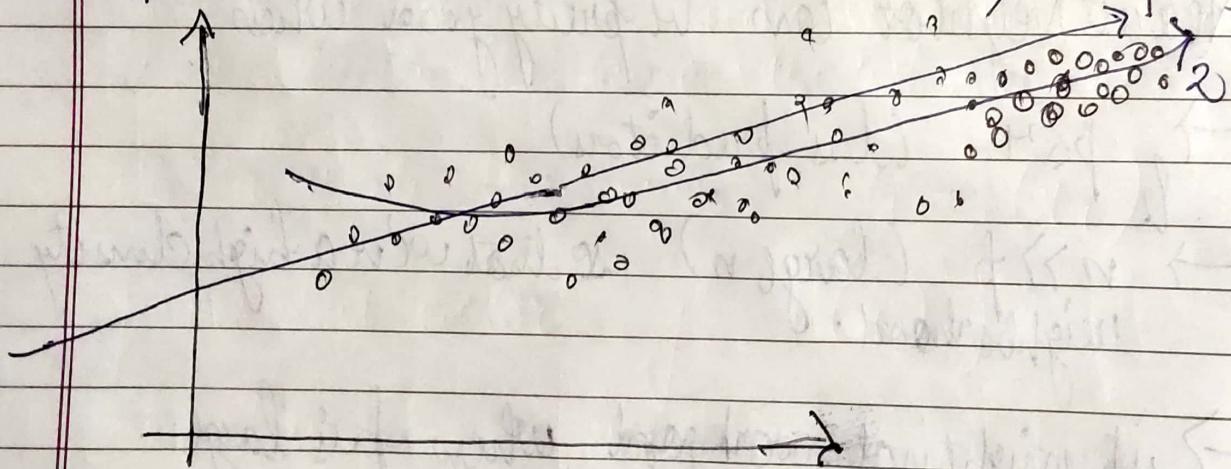
So, what should we rather do:

### Parametric & Structured Models:

The linear model is a good example of parametric models.

$$f(x) = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_p x_p$$

- A linear model is specified in terms of "fit" parameters.
- We estimate the parameters by fitting the model to the data (training).
- Although it is almost never correct, a linear model serves as a good & interpretable approximation to the true  $f(x) \approx E(Y|X=x)$



1. Represents Linear Model

$$f(x) = \beta_0 + \beta_1 x$$

2. Represents Quadratic Model

$$f(x) = \beta_0 + \beta_1 x_1 + \beta_2 x_2^2$$

\* These both are linear models but with different parametric forms  
(It fits slightly better)

Parametric Methods involve a two-step method approach:

1. First, is to take an assumption of function form or shape of  $f$ .  
For example, a linear function.
2. Second step involves training the model & estimating the values of  $\beta_s$ .

Non Parametric Methods:

These do not make explicit assumptions about the functional forms of  $f$ .

Instead they seek to estimate ' $f$ ' & get as close as possible to it.

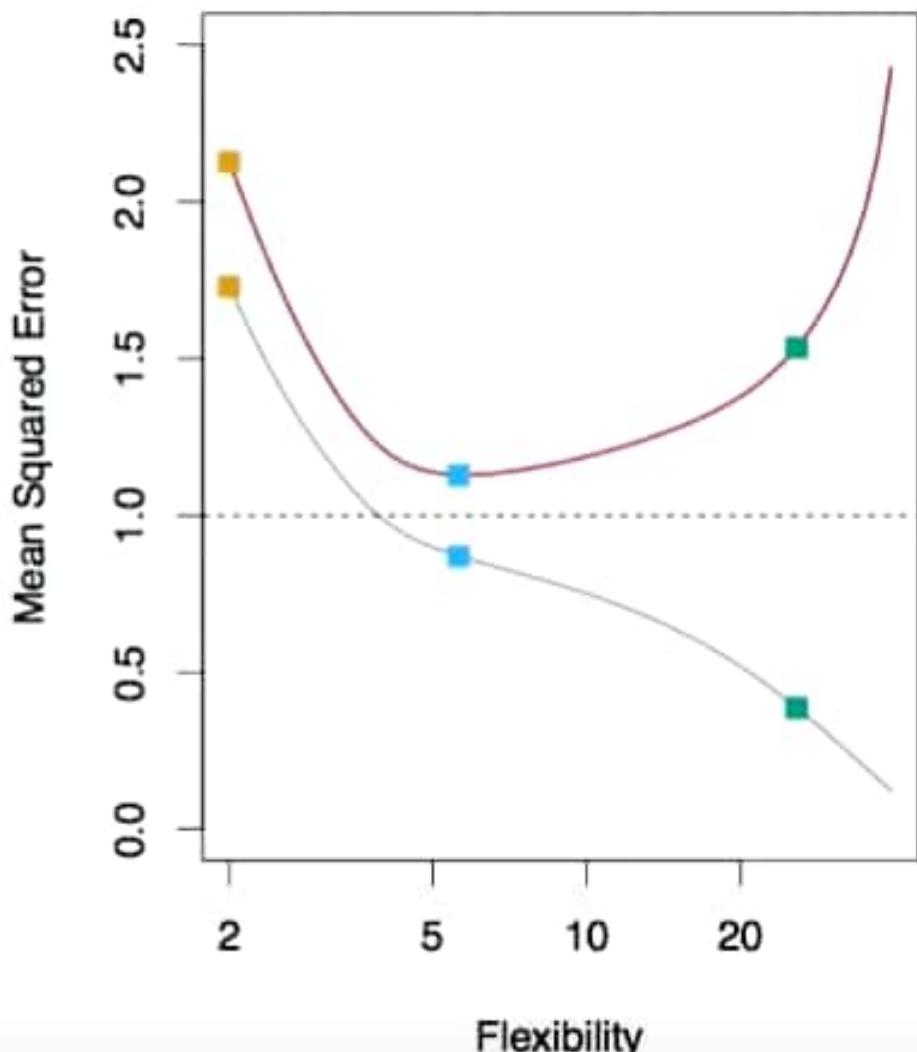
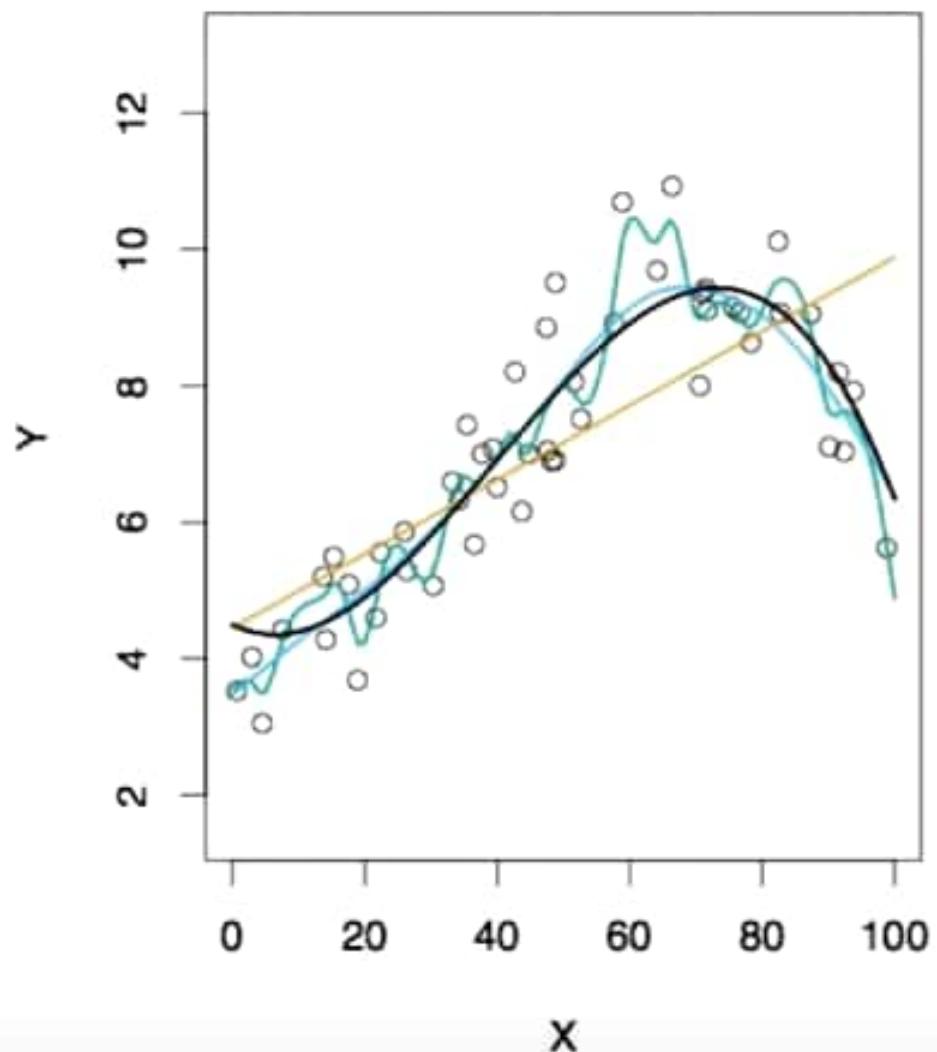
The ~~&~~ the only disadvantage is it ends up overfitting the data.

### The Tradeoff Between Prediction Accuracy & Flexibility

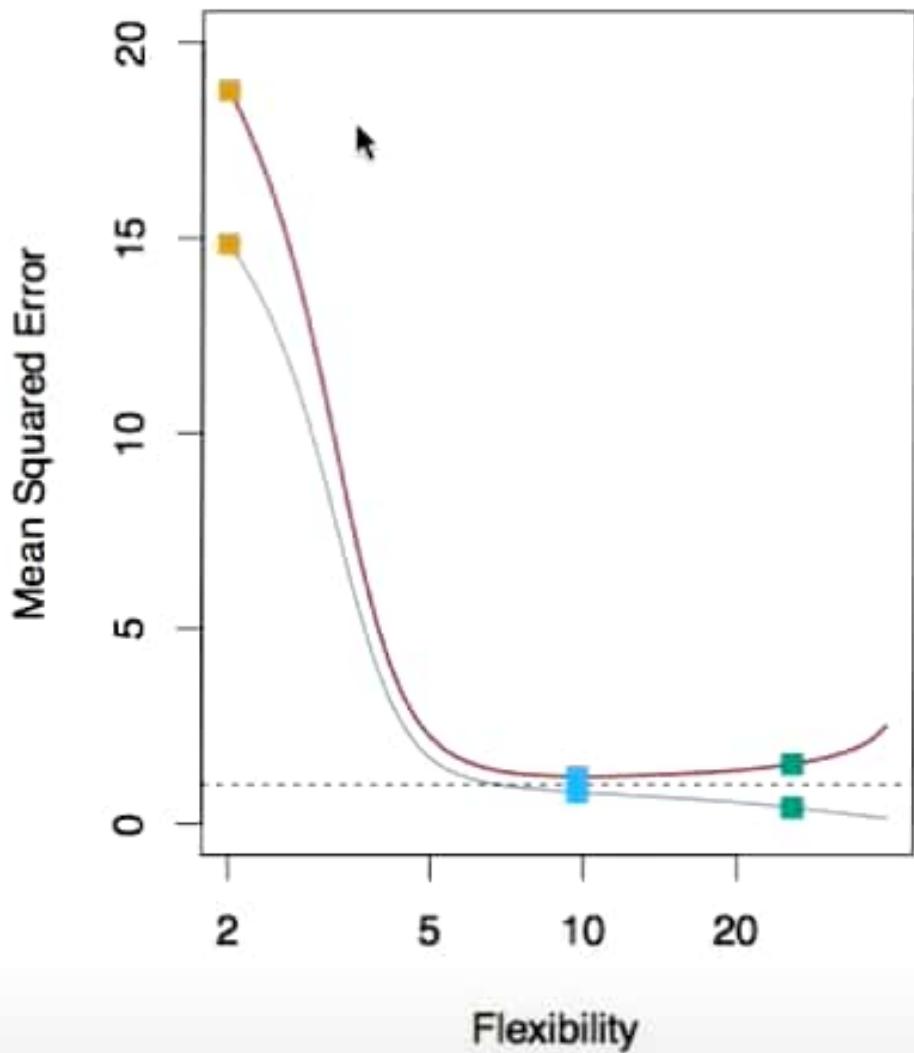
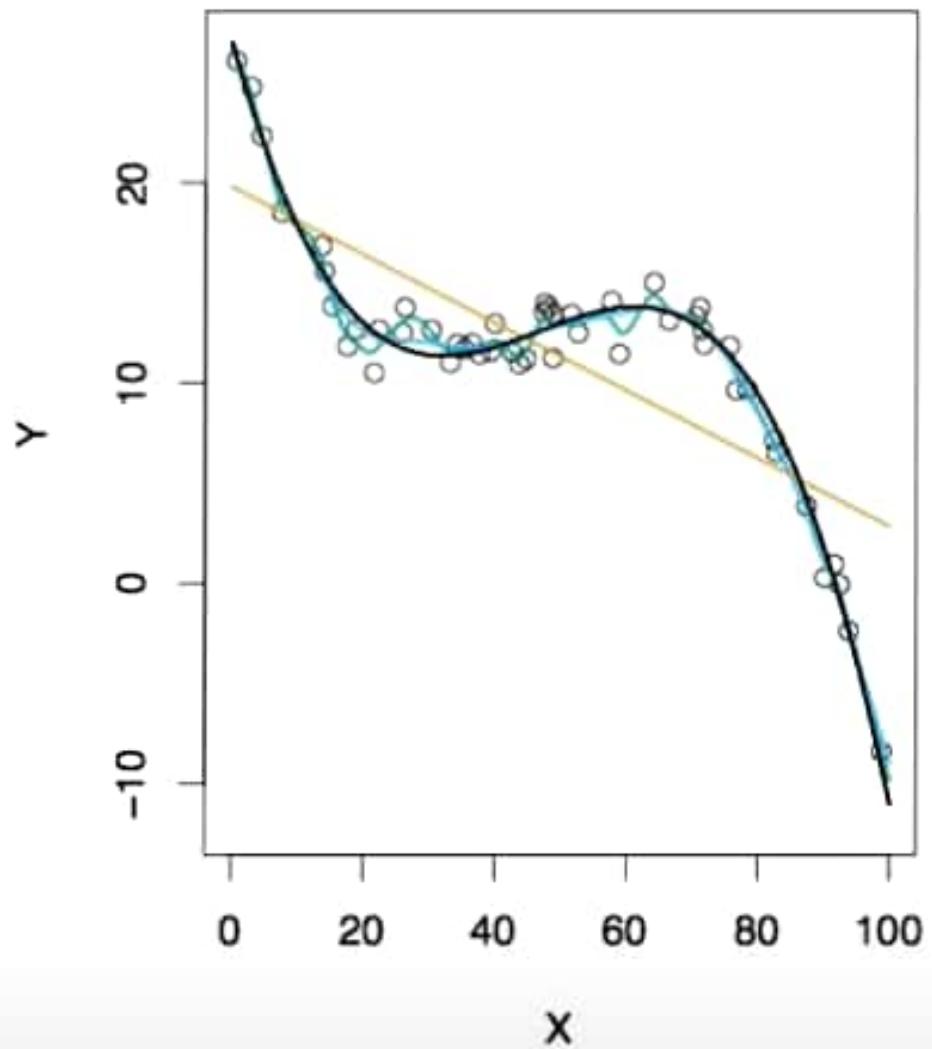
There are multiple methods we can use for predictions. Some are less flexible or more restrictive.

For ex: Linear Regression is very inflexible approach as it can only create linear functions.

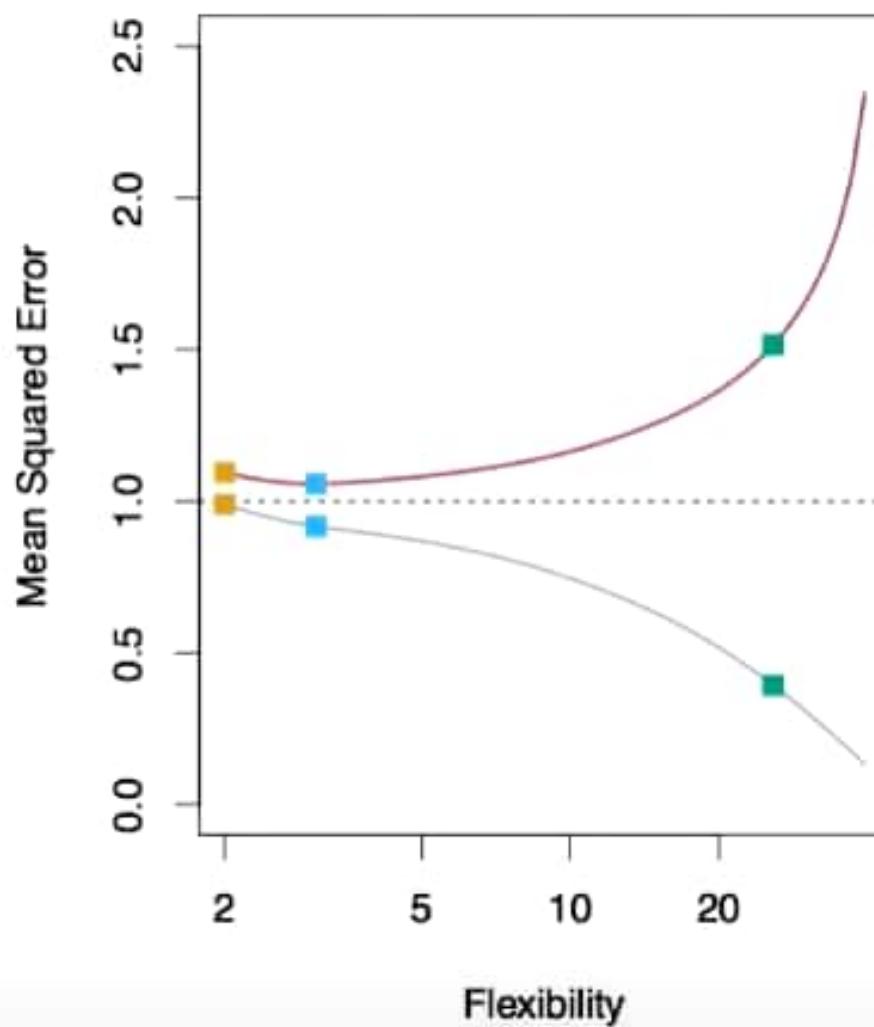
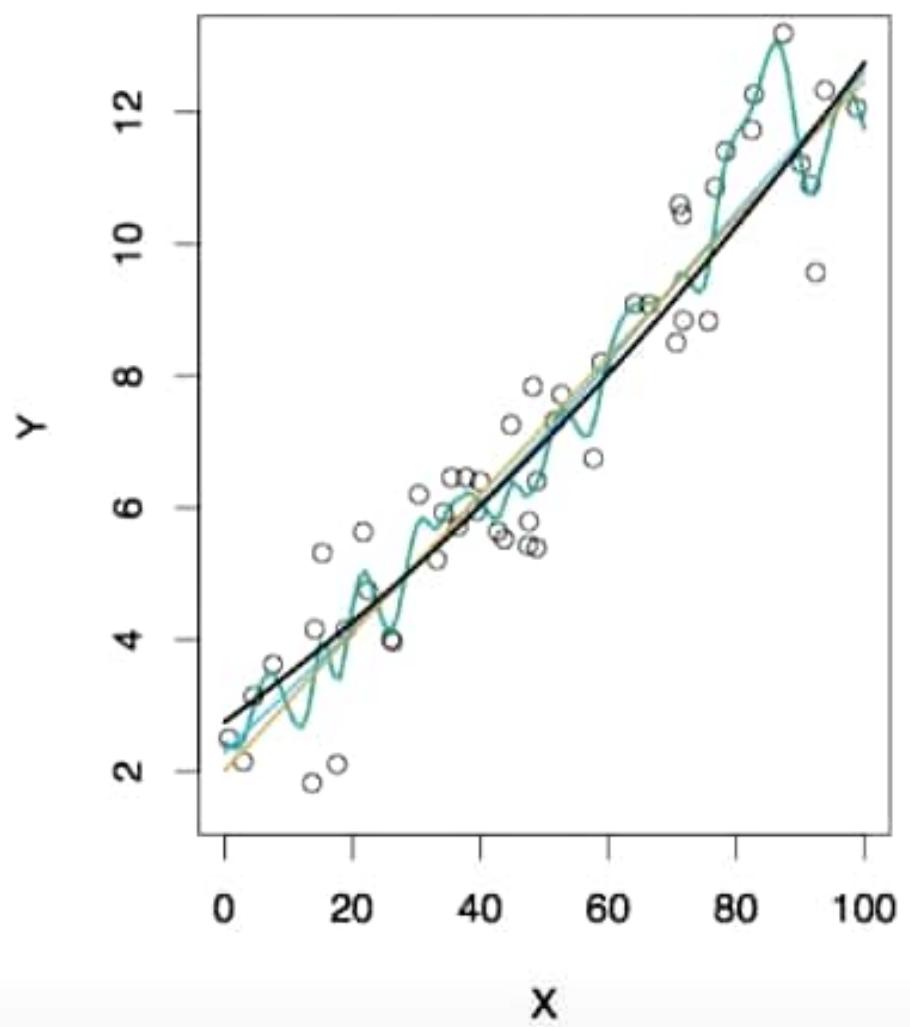
Other methods such as SVM's or thin plate splines are more flexible.



Black curve is truth. Red curve on right is  $MSE_{Te}$ , grey curve is  $MSE_{Tr}$ . Orange, blue and green curves/squares correspond to fits of different flexibility.



Here the truth is wiggly and the noise is low, so the more flexible fits do the best.



Here the truth is smoother, so the smoother fit and linear model do really well.

So, this question arises, why would we ever choose restrictive models?

There are several reasons:

- They are more interpretable
- They prone to overfit

② Added a photo above.

## Assessing Model Accuracy:

Let us suppose we fit a model to a linear data & that is  $f(x)$  & we wish to see how it performs.

To do so, we calculate Mean Squared Error  
It is exactly the same as it sounds:

We calculate the error:  $y_i - \hat{f}(x_i)$

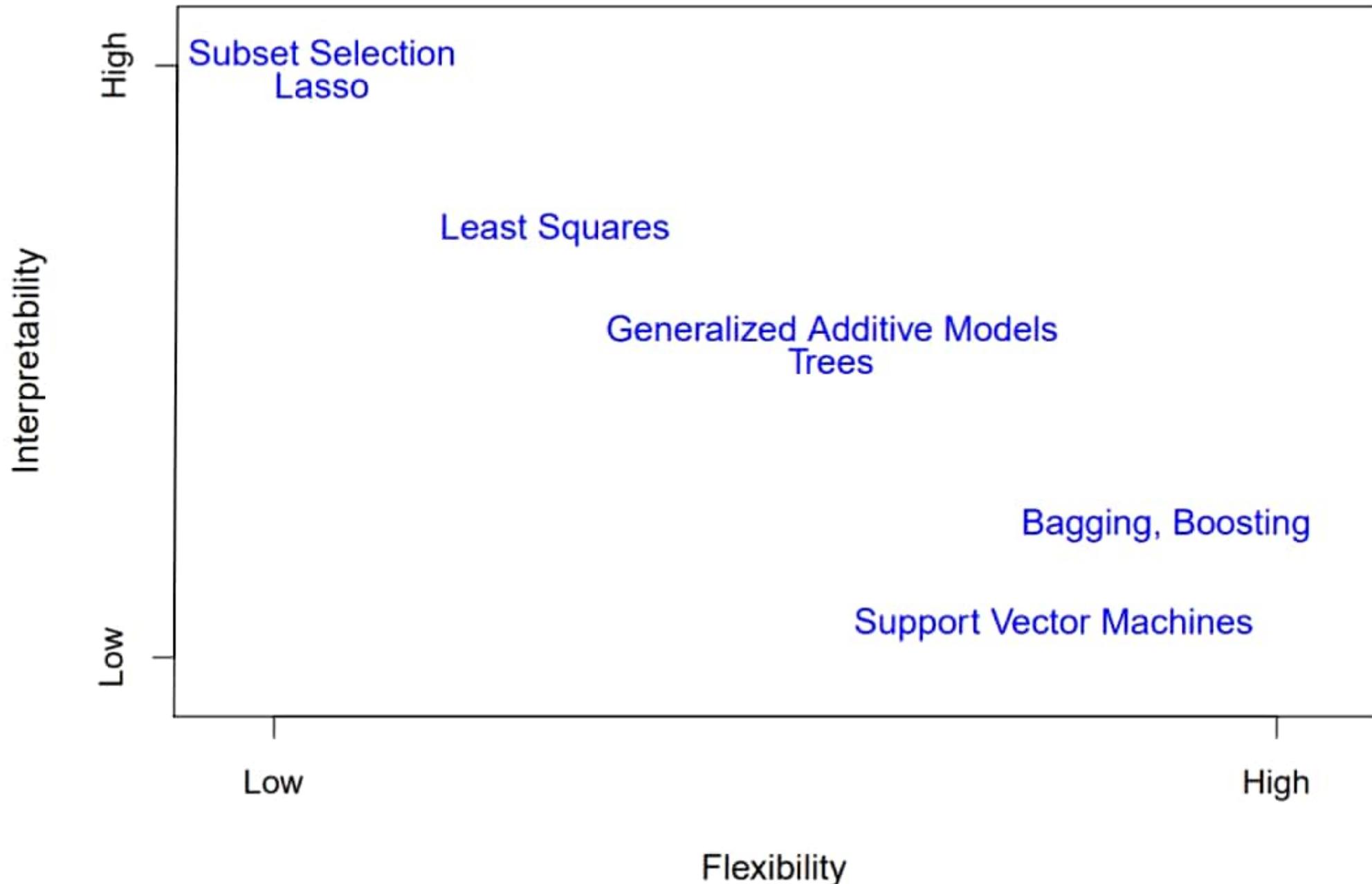
Square the error:  $(y_i - \hat{f}(x_i))^2$

Take its mean over all the training data.

$$\boxed{\frac{\sum_{i=1}^{N_T} (y_i - \hat{f}(x_i))^2}{N_T}}$$

But this may be biased toward overfit models  
So, we ~~can't~~ calculate:

$$MSE_T = \frac{\sum_{i=1}^{N_T} (y_i - \hat{f}(x_i))^2}{N_T}$$



where,  $N_{TE}$  &  $N_{Tr}$  is the number of samples of Training & Testing Data respectively.

### Bias-Variance Tradeoff:

Suppose, we're to fit a model  $\hat{f}(x)$  to some training data  $T_r$ , & let  $(x_0, y_0)$  are test samples drawn from population & true model is

$$Y = f(x) + \epsilon$$

$$E(y_0 - \hat{f}(x_0))^2 = \text{Var}(\hat{f}(x_0)) + [\text{Bias}(\hat{f}(x_0))]^2 + \text{Var}(\epsilon)$$

$$\text{where } \text{Bias}(\hat{f}(x_0)) = E[\hat{f}(x_0)] - f(x_0)$$

So, as flexibility of  $\hat{f}$  increases Bias increase (Overfitting) & variance increases (as the new testing data will have increased variance i.e. chances of misclassification due to overfitting). In high flexible data a little change in  $x_0$  can change  $\hat{f}$  very much hence variance of  $\hat{f}$  increase.

### Classification Problem:

Here,  $Y$  is Qualitative/Categorical not a continuous variable say  $C = \{\text{spam, ham}\}$  or  $C = \{0, 1, \dots, 9\}$  digits.  
 Our goals are:

- Building a classifier to predict classes for unlabelled data.
- Assess uncertainty in each classification.
- Understand the role of all the predictors.

The concepts like <sup>bias variance Tradeoff</sup> & <sup>Flexibility Interpretability Tradeoff</sup> are transferable here.

Here, the model is assessed using:

$$\frac{1}{n} \sum_{i=1}^n I(y_i \neq \hat{y}_i)$$

which returns 1 when  $y_i \neq \hat{y}_i$   
& 0 when  $y_i = \hat{y}_i$

so, it gives the proportion of Training Error.

### Naive Bayes Classifier:

$$P_{k,i}(x) = P_k(Y=k | X=x)$$

where  $k = 1, \dots, K$

This is called a conditional probability at a given  $x$ .

The Bayes Classifier at  $x$  is:

$$c(x) = j \text{ if } P_j(x) = \max(P_1(x), P_2(x), \dots, P_K(x))$$

where, Bayes Error Rate is given by  $P_{K(x)}$

$$1 - E(\max P_k(Y=j | X))$$

## Nearest Neighbor:

The nearest neighbor averaging can be done as before.  
 This also breaks down as dimensions grow.

(60)

→ Model is accessed using

$$= E(I[y_i \neq \hat{y}_i])$$

$$= \frac{1}{N} \sum I[y_i \neq \hat{y}_i]$$

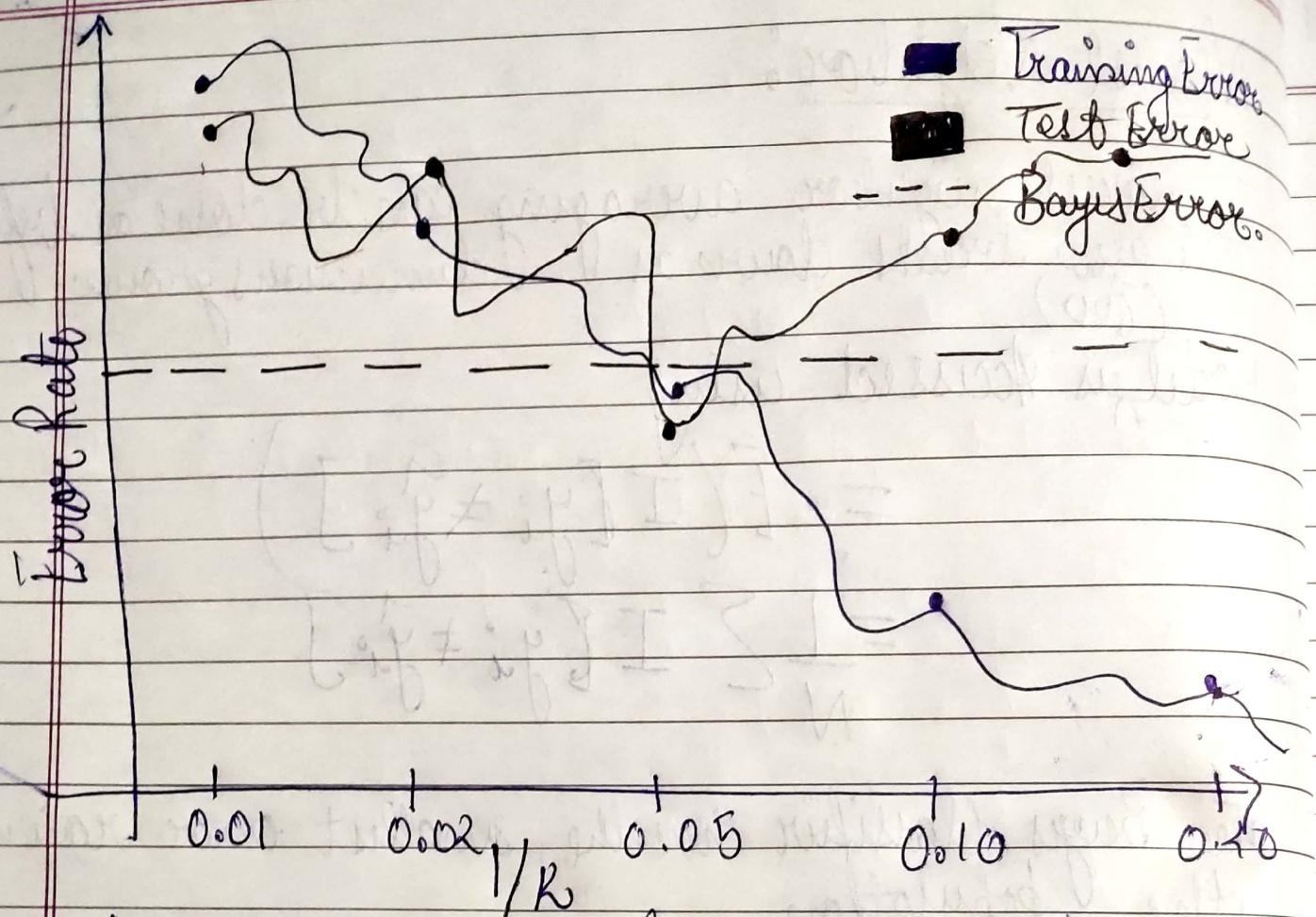
→ The Bayes classifier has the smallest error rate in the population.

→ SVM builds structured model for  $\ell(x)$ .

→ also models like logistic regression, boosting & bagging.

In K-nearest neighbors we calculate the "K" smallest distant point from the data say  $(x_0)$  selects  $(x_{10}, x_{20}, x_{30}, x_{40}, \dots, x_{(n-K)})$  as averages & takes the output with high probabilities.

I/P	O/P
$\{0, 0, 0, 1, 1\}$	0
$\{0, 0, 0, 0, 1\}$	0
$\{1, 1, 1, 0, 0\}$	1
$\{1, 1, 0, 0, 0\}$	0
$\{1, 1, 1, 0, 1\}$	1



We can see that test error rates increase after some time again & obviously nothing does better job than naive bayes theoretically.