

# General overview of the game we play in machine learning.

- Nature generates data according to some function  $f: x \rightarrow y$  we don't know
- In machine learning we want a good approximation of  $f$  on out of sample data.
- We start with set of possible functions (Hypothesis space  $H$ ) and try to pick a function  $\hat{f}$  best approximating true  $f$  (which we don't know)
  - Complex  $H$  means good chance of approximating  $f$  but hard to find approximation?
  - Simple  $H$  means good generalization on out of sample data.

Bias variance analysis effectively asks following questions.

- How well  $H$  can approximate  $f$  overall.
- Given the sample how good is our approximation.

Let's use square error analysis on linear regression Let's say  $f$  is true function and  $\hat{f}$  is our choice of target function based on sample dataset  $D$  in hypothesis space  $H$  Then

$$E_{out}(\hat{f}_D) = E_x(\hat{f}_D(x) - f(x))^2$$

This depends on a particular sample  $D$ . If we take another same number of sample this will change. To do error analysis we need to get rid of a particular sample  $D$  choice.

Let's take expectation with respect to distribution over  $D$

$$\begin{aligned} E_D(E_{out}(\hat{f}_D)) &= E_D E_x(\hat{f}_D(x) - f(x))^2 \\ &= E_x E_D(\hat{f}_D(x) - f(x))^2 \end{aligned}$$

Define average hypothesis

$$\bar{f} = E_D(\hat{f}_D(x))$$

$$\begin{aligned} E_D(\hat{f}_D(x) - f(x))^2 &= E_D[\hat{f}_D(x) - \bar{f}(x) + \bar{f}(x) - f(x)]^2 \\ &= E_D[(\hat{f}_D(x) - \bar{f}(x))^2 + (\bar{f}(x) - f(x))^2 + 2(\hat{f}_D(x) - \bar{f}(x))(\bar{f}(x) - f(x))] \\ &= E_D[(\hat{f}_D(x) - \bar{f}(x))^2] + [\bar{f}(x) - f(x)]^2 \end{aligned}$$

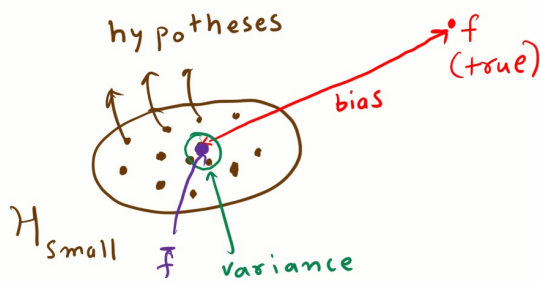
$$\underbrace{E_D[(\hat{f}_D(x) - \bar{f}(x))^2]}_{\text{variance}(x)} + \underbrace{[\bar{f}(x) - f(x)]^2}_{\text{bias}(x)}$$

bias = a measure of best we can do in our hypothesis set

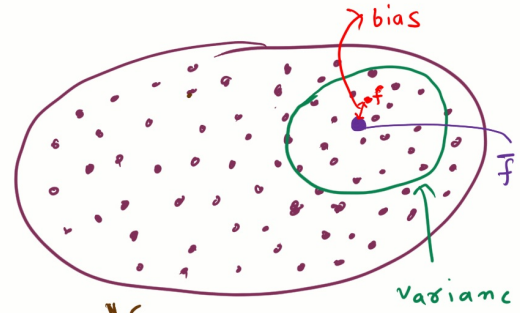
variance = how far our hypothesis based on  $D$  is away from average

$$\begin{aligned} E_x E_D(\hat{f}_D(x) - f(x))^2 &= E_x[E_D[(\hat{f}_D(x) - \bar{f}(x))^2] + [\bar{f}(x) - f(x)]^2] \\ &= E_x[\text{bias}^2(x) + \text{variance}(x)] \\ &= \text{bias}^2 + \text{variance} \end{aligned}$$

# BIAS VARIANCE TRADEOFF

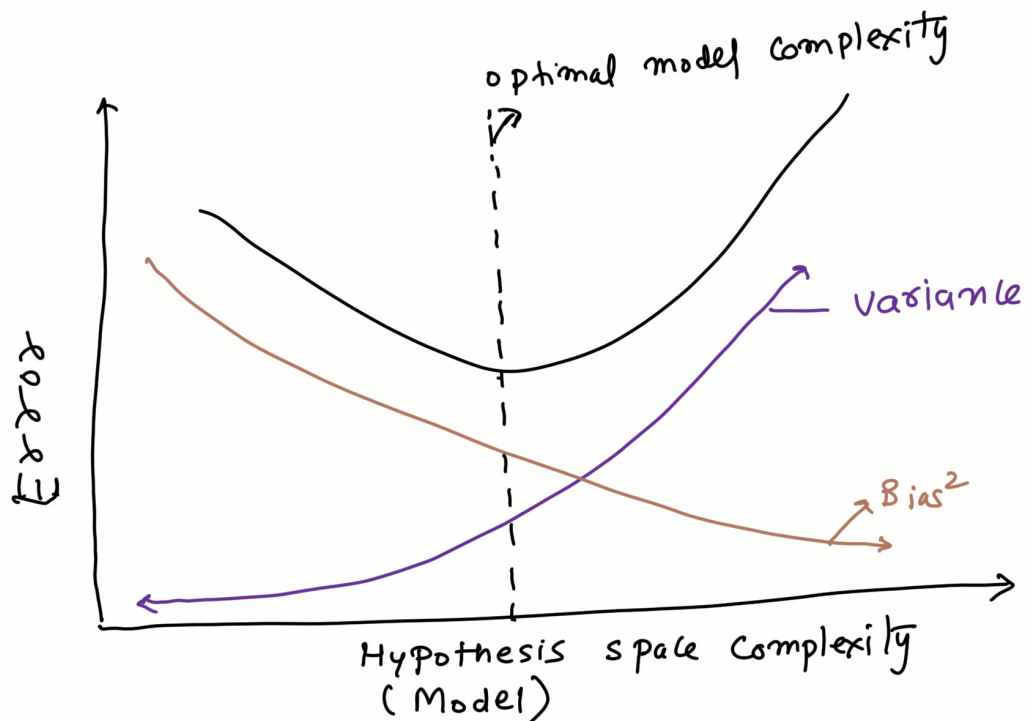


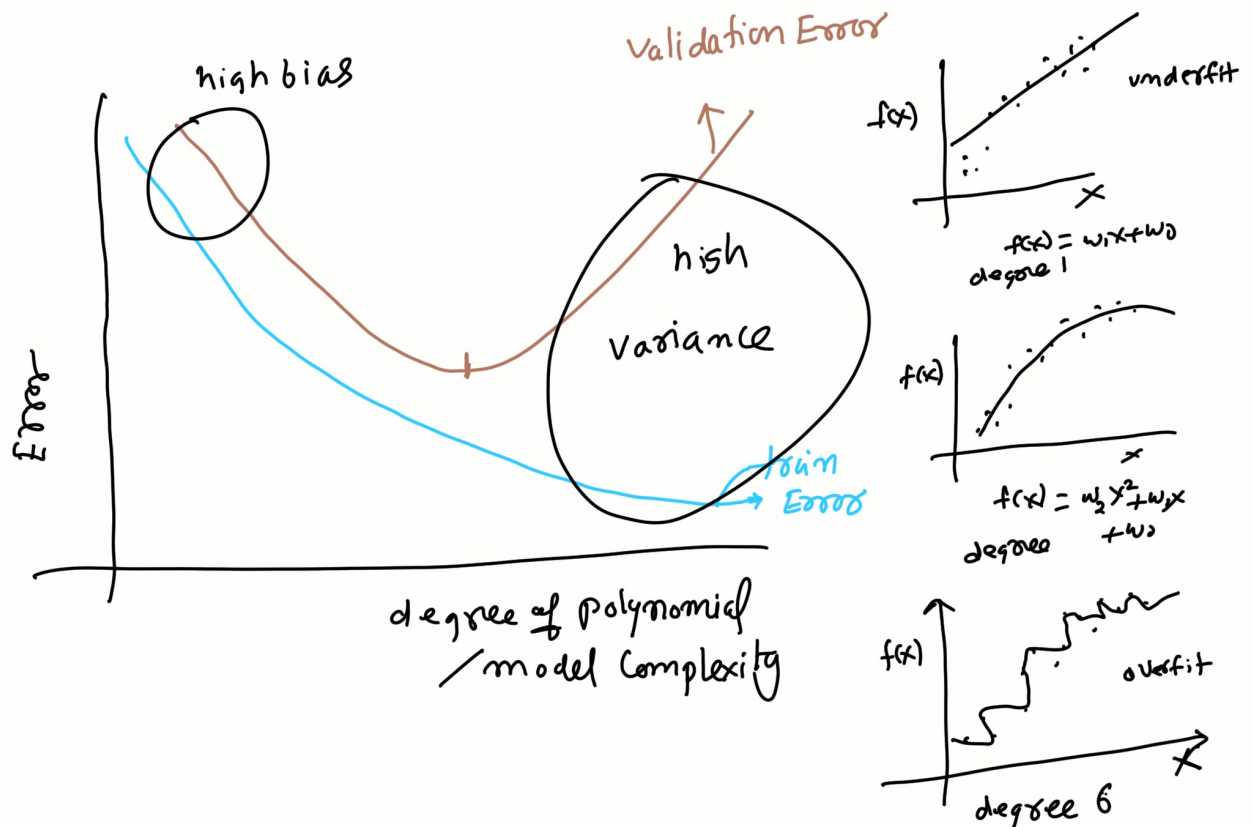
- High bias
- Low variance



- Low bias
- High variance

## Bias Variance Trade off





## Evaluation of Learning algorithms

We need a way to measure how well learning algorithm performs. Till now we have seen following measures

- Accuracy
- Precision, Recall, ROC curve

## let's make it more formal

- Given an Hypothesis space  $H$ , using training data  $D$  we select a hypothesis  $\hat{f} \in H$
- Given a feature vector  $x$  we predict using  $h$  as  $\hat{y} = \hat{f}(x)$ . Let say true value is  $y$ .

## For Regression

- MAE (Mean absolute error)  $\frac{\sum_i^N |\hat{f}(x_i) - y_i|}{N}$
- MSE (Mean square error)  $\frac{\sum_i^N \|\hat{f}(x_i) - y_i\|^2}{N}$  ## For Classification
  - Miss classification error =  $\frac{\sum_i^N 1(\hat{f}(x_i) \neq y_i)}{N}$
  - Build Confusion matrix, Remember TP, TN, FP, FN for binary classification. You can build confusion matrix for K class classification too.
  - ROC

## Confusion Matrix for a particular value of threshold $\tau$

	$y = 1$	$y = 0$	
$\hat{y}=1$	TP	FP	$\hat{N}_+ = TP + FN$
$\hat{y}=0$	FN	TN	$\hat{N}_- = TN + FP$
-----	:-----:	-----:	:-----
	$N_+ = TP + FN$	$N_- = TN + FP$	$N = TP + FN + TN + FP$

$$\text{precision} = \frac{TP}{TP + FP}$$

$$\text{recall} = \text{True positive rate (TPR)} = \frac{TP}{N_+} \approx p(\hat{y} = 1 | y = 1)$$

$$\text{false positive rate (FPR)} (\text{type I error rate}) = \frac{FP}{N_-} \approx p(\hat{y} = 1 | y = 0)$$

## Sample error vs True Error(In classification setting)

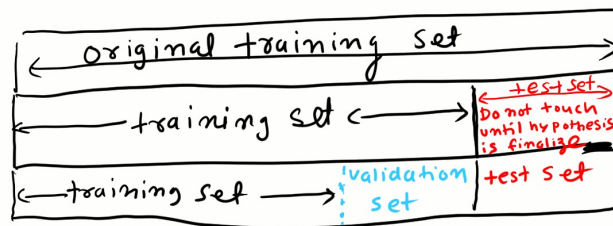
**Sample error** : The sample error of a hypothesis  $\hat{f}(x)$  with respect to true function(target)  $f$  on a sample  $D = \{x_i, y_i = f(x_i)\}_{i=1}^N$  is miss classification error  $\frac{\sum_{i=1}^N \mathbb{1}(\hat{f}(x_i) \neq f(x_i))}{N}$

**True Error** The true error of a hypothesis  $\hat{f}(x)$  with respect to true function(target)  $f$  and distribution  $P$ , is the probability of misclassifying a random sample from distribution  $P$ ,  $Pr_{x \sim P}[\hat{f}(x) \neq f(x)]$

We need good **True Error**

## Evaluation of hypothesis under limited sample data

- Split the samples into train and test



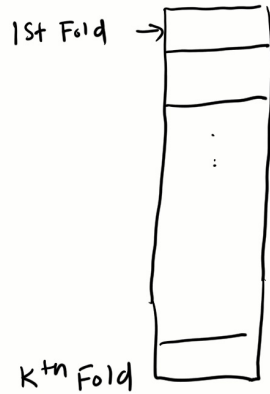
What if train set is too small

- we have small samples size.
- Can't touch test set
- Can we use data in validation for training too?

yes. Procedure is called k-fold cross validation

## K - Fold Cross Validation

- partition data into K subset



→ For  $i$  in  $\{1, 2, \dots, K\}$

+ use Fold  $i$  for validation

+ use rest of the Folds

for training

+ calculate  $Error_i$  on Fold  $i$

$$\text{Final Validation Error} = \frac{\sum_{i=1}^K Error_i}{K}$$