

Ensembles: Random Forests } Scikit learn  
 Boosting (XGBoost)

Bias-Variance Decomposition:

Assume we want to estimate some population parameter  $\theta$ , and we have a sample  $x_1, \dots, x_n$

Bias:  $E(\hat{\theta}) - \theta$

Variance:  $Var(\hat{\theta})$

concrete:

- $x_1, \dots, x_n \sim N(\mu, 1)$
- $\theta = \mu$
- $\hat{\theta}_1 = \bar{x}$
- $\hat{\theta}_2 = \text{sample median}$

Bias: how far from the truth are we?

Variance: how noisy our estimator is.

Variance over many datasets:

$$D_1 = x_{11}, \dots, x_{1n} \rightarrow \bar{x}_1 \hat{\theta}_1$$

$$D_2 = x_{21}, \dots, x_{2n} \rightarrow \bar{x}_2 \hat{\theta}_2$$

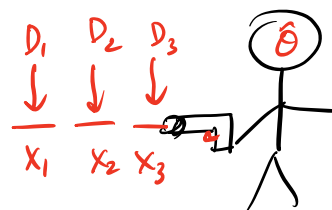
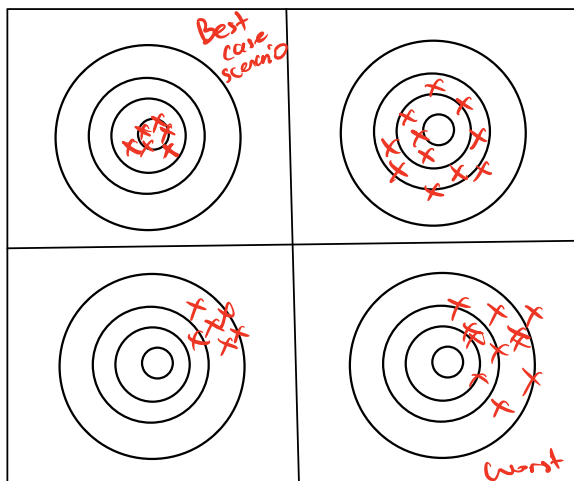
$$\vdots$$

$$\hat{\theta}_\infty$$

low Variance high

Bias

high



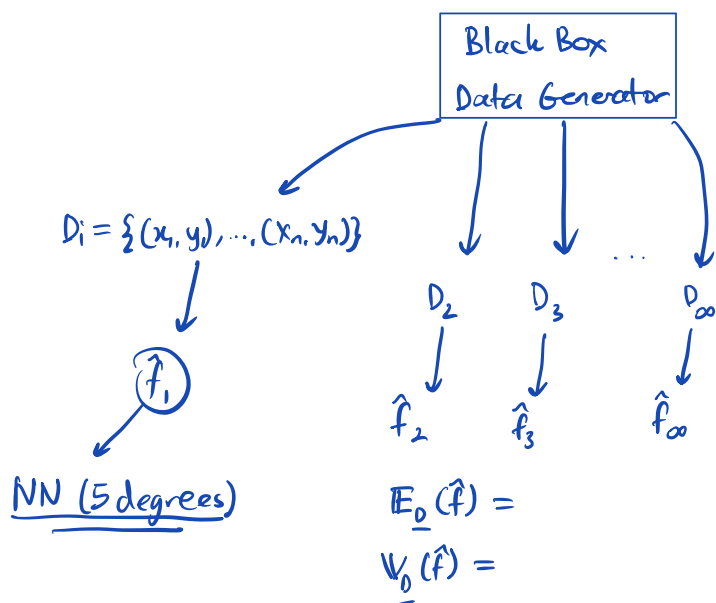
## Bias - Variance Decomposition:

- truth:  $f(x)$
- estimate:  $\hat{f}(x) = \hat{f}(x; D)$ ,  $D = \{(x_i, y_i) : i=1, \dots, n\}$   
 $\frac{1}{n} \sum (x-y)^2$
- Expected: MSE on a test point  $x_0$ :

$$\underbrace{\mathbb{E}_D (f(x_0) - \hat{f}(x_0; D))^2}_{\text{Expected MSE at } x_0} = \underbrace{\text{Bias}_D (f(x; D))^2}_{\text{Bias}^2} + \underbrace{\text{Var}_D (\hat{f}(x_0; D))}_{\text{Variance}} + \underbrace{\sigma^2}_{\text{irreducible error}}$$

$\sigma^2$ : irreducible error

reduce test MSE  $\rightarrow$  reduce MSE  
 reduce variance



## ML Pipeline:

1. Get a Dataset
2. Build a model
3. Evaluate model
4. Use model in wild

## Simulating Dataset

• true function:  $f$

$D_1 = ?$

$$\begin{matrix} x_1 & y_1 = f(x_1) + \epsilon_1 \\ x_2 & y_2 = f(x_2) + \epsilon_2 \end{matrix}$$

$$\vdots \quad \vdots$$

$$x_n \quad y_n = f(x_n) + \epsilon_n$$

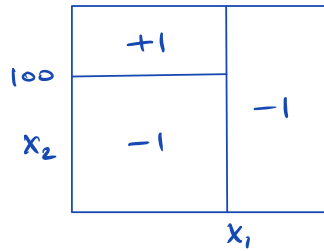
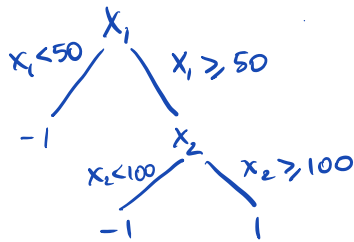
$$\epsilon_i \sim N(0, \sigma^2)$$

high Bias model: Strong assumptions shape of the fit is predetermined regardless of data.

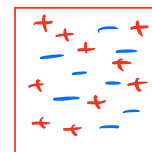
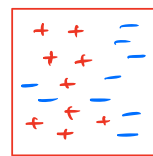
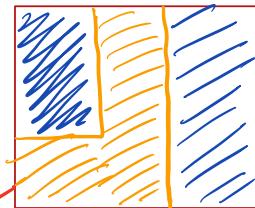
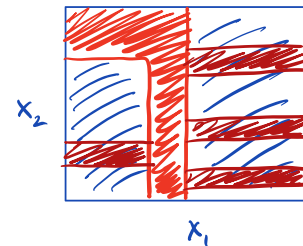
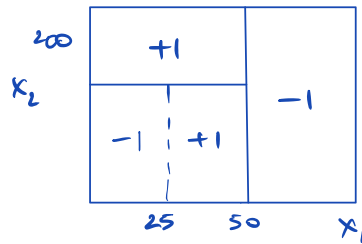
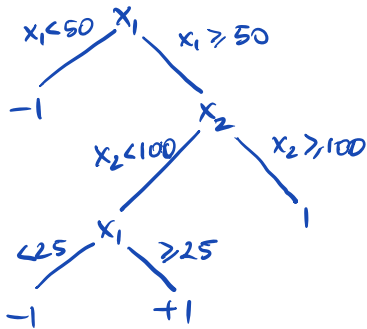
high Variance model: Error due to fitting randomness

Back to Ensembles  $\rightarrow$  how do ensembles approach B-V decomposition.

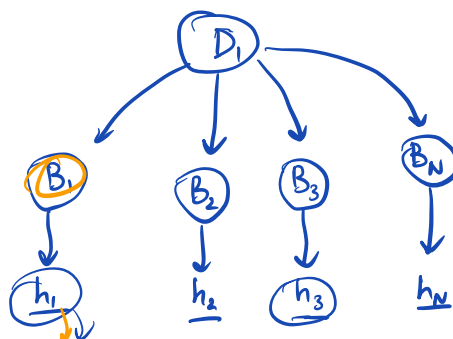
Decision Trees:



DTs: low bias  
high variance



Bagging : Bootstrap Aggregation



$B_i$  : sample from  $D_1$  with replacement  
m points

where  $m < n$

$D_1 = \{(x_1, y_1), \dots, (x_n, y_n)\}$

$h_i$  = tree built on  $B_i$

$x_1 = (x_{11}, x_{12}, \dots, x_{1d})$   
 $\downarrow$   
 $\tilde{x}_{1B1} = (x_{11}, x_{12}, x_{16}, x_{18})$

Ensemble → Bagging Classifier

Average over many trees.

new point  $\underline{C(x)} = \begin{cases} +1 & \text{if } \frac{1}{n} \sum_{i=1}^n h_i(x) > 0 \\ -1 & \text{otherwise} \end{cases}$  majority vote

Why?

Intuition: Ensembles: Wisdom of Crowds versus traditional building expert.

Randomisation induces independence

1. Randomise across samples
2. Randomise across features

(Random Forest)

Why does model averaging work?

$Y_1, \dots, Y_n$  Independent  $(\mu, \sigma^2)$

$$E(Y_i) = \mu$$

$$V(Y_i) = \sigma^2$$

$$\text{bias}(f) = E(f) - f$$

$$= b$$

$$E\left(\frac{1}{n} \sum_{i=1}^n Y_i\right) = \frac{1}{n} \sum_{i=1}^n E(Y_i) = \frac{1}{n} \sum_{i=1}^n \mu = \underline{\mu}$$

$$V\left(\frac{1}{n} \sum_{i=1}^n Y_i\right) = \frac{1}{n^2} V\left(\sum_{i=1}^n Y_i\right) + \text{circled O}$$

$$= \frac{1}{n^2} \sum_{i=1}^n V(Y_i)$$

$$= \frac{1}{n^2} \sum_{i=1}^n \sigma^2$$

$$= \frac{n\sigma^2}{n^2}$$

$$= \frac{\sigma^2}{n} < \sigma^2$$

so: average has the same mean and lower variance than individual values.

Average of  $N$  models with bias  $b$  is  $b$ , but the variance is lower.

Appeal of Ensembles:

- Base classifier (Decision Trees)  
is really simple/fast to train
- RF →  $N$  DTs,
- Deep learning

Disadvantage:

- DTs are really simple to interpret
- RFs lose any easy interpretability