* Agenda and Recap
    * Normality of Residuals
    * Homoscedasticity
    * Auto-collinearity
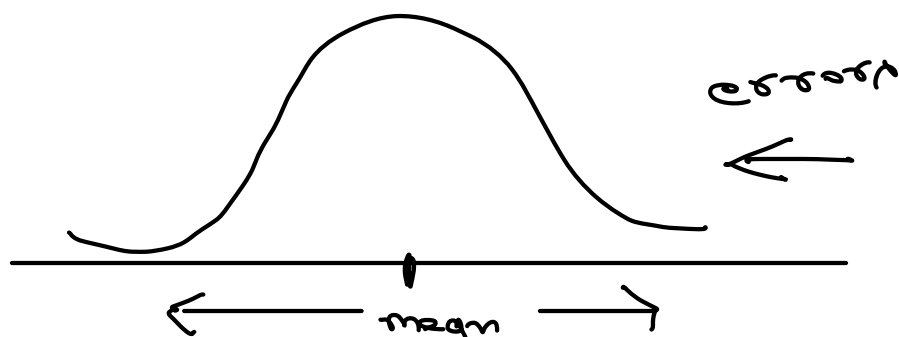
3

* Gradient Descent Variants

* Polynomial Regression

* Generalization and Occam's razor
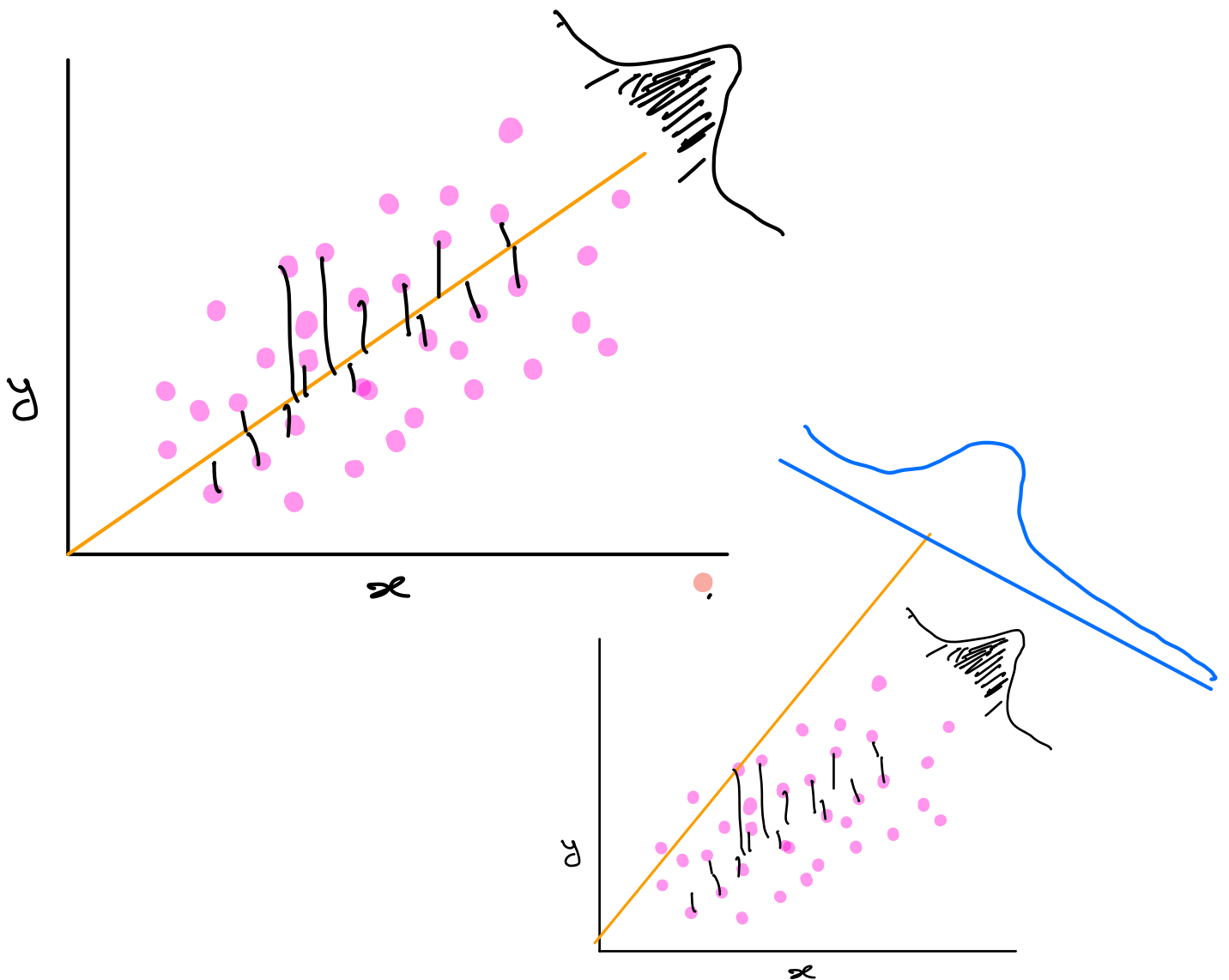
* Underfitting and overfitting

* Bias Variance Tradeoff

---

Assumption 3: Errors are Normally distributed



Step 1 : Build Model

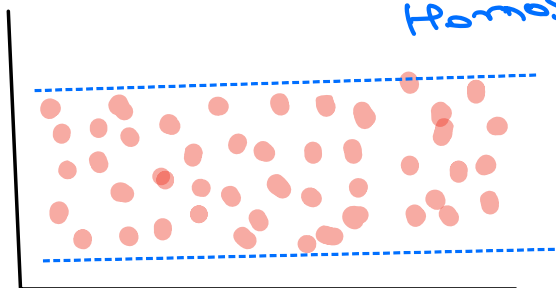Step 2 : Calculate Error

Step 3 : Plot Errors with Histogram

→ Variance of residuals $(\epsilon_i)$ vs predictions $(\hat{y}_i)$ should be constant

$y_i$ vs $\epsilon_i$

Homoscedasticity

Error $(y_i - \hat{y}_i)$

$\hat{y}_i$

Heteroscedasticity

Error $(y_i - \hat{y}_i)$

$\hat{y}_i$

* Goldfeld Quandt Test

⮕ Null Hypothesis: Dataset has Homoscedasticity
* If P value ≤ significant level Threshold
⮕ reject Null Hypothesis

* To mitigate:

⮕ Remove outlier

⮕ Perform Non-Linear Transformation such as box-cox

(correlate between Errors)

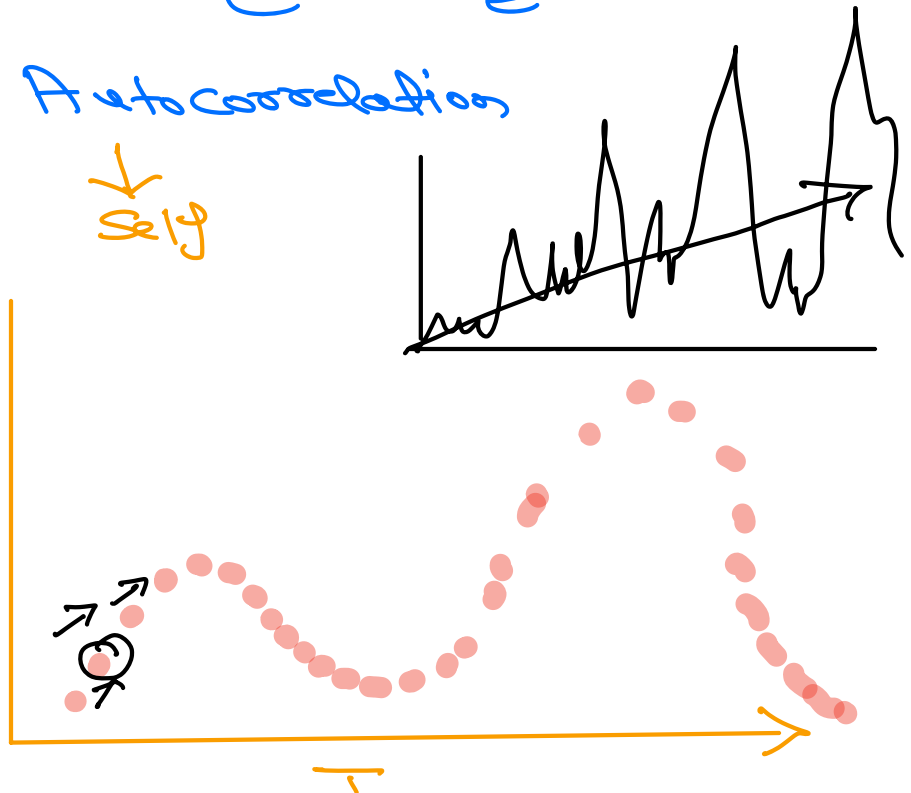Assumption 5: No Autocorrelation

self

$\epsilon_i$ vs T

⮕ Time-series Model

$\epsilon_i$

T

(Just 1 update per iteration)

$$W_j^{new} = W_j^{old} - \eta \frac{\partial L}{\partial W_j^{old}}$$

$$\frac{2}{3} \sum_{i=1}^{M} (\hat{y}_i - y_i) \times x_{ij}$$

$1000 \Rightarrow 1000$ updates
iter

For One Single iteration/update

↦ 1 million Errors and Gradients

⇒ Update Weights

1 million

1 Million

$k \Rightarrow$ between 1 and N

let's $k = 256$

Batches

$$\frac{2}{K} \sum_{i=1}^{M} \frac{K}{K} (\hat{y}_i - y_i) \times x_{ij}$$

$$W_j^{new} = W_j^{old} - \eta \frac{\partial L}{\partial W_j^{old}}$$

Batches ⇒ $\left( \dfrac{1 \text{ million}}{256} \right)$  $k$ Updates ←

(updates per iteration)

$1000$ iters → $1000 \times$ Batches

→ 1 (Stochastic Gradien Descent (SGD))

$k$ → N (Batch GD)

(Weight update) in Single Epoch

→ (1, N) (Mini-Batch GD) MGD

| 1 million rows of Data | shuffle $\Rightarrow$ | 512 | $B_1$ |
|---|---|---|---|
| | | 512 | $B_1$ |
| | | 512 | ... |
| | | ... | ... |
| | | | $B_{1700}$ |

| $\theta_+$ update $B_1$ |
|---|
| $\theta_+$ update $B_2$ |
| $\theta_+$ update $B_3$ |
| ... |

Single iteration (Epoch)

## Comparison: Minibatch vs Batch

Mini-batch

Loss

iterations

Batch

Loss

iterations

# Polynomial Regression

John



$$y = \omega_1 x_1 + \omega_0$$

Height

Weight



Height

Weight

$$y = \omega_1 \rho_1 + \omega_2 \rho_1^2 + \omega_0$$

New feature

\* How Does this Look

→ In Data

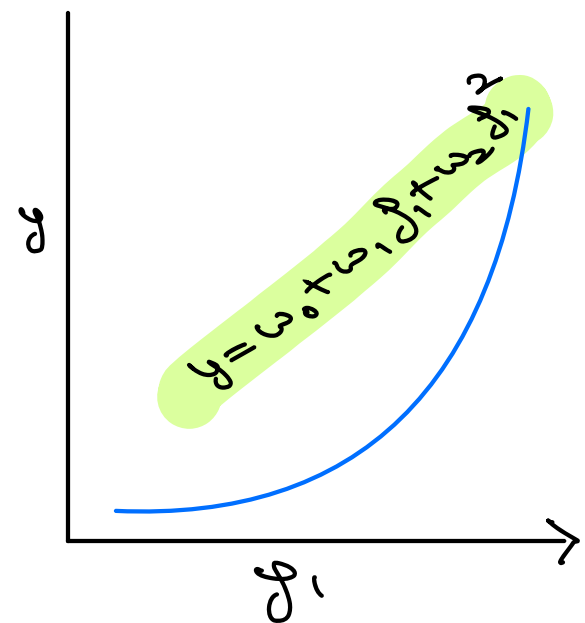| x | y |
|---|---|
| $x_{11}$ | $y_1$ |
| $x_{21}$ | $y_2$ |
| $x_{31}$ | $y_3$ |
| $x_{41}$ | $y_4$ |
| $x_{n1}$ | $y_n$ |

$\Rightarrow$

| $x^2$ | $x_1$ | y |
|---|---|---|
| $x_{11}^2$ | $x_{11}$ | $y_1$ |
| $x_{21}^2$ | $x_{21}$ | $y_2$ |
| $x_{31}^2$ | $x_{31}$ | $y_3$ |
| $\vdots$ | $x_{41}$ | $y_4$ |
| $x_{n1}^2$ | $x_{n1}$ | $y_n$ |

$x_1^2$   $x_1$

$x_2 \Rightarrow x_2^2$

$x_1 \times x_2$

→ In plot.





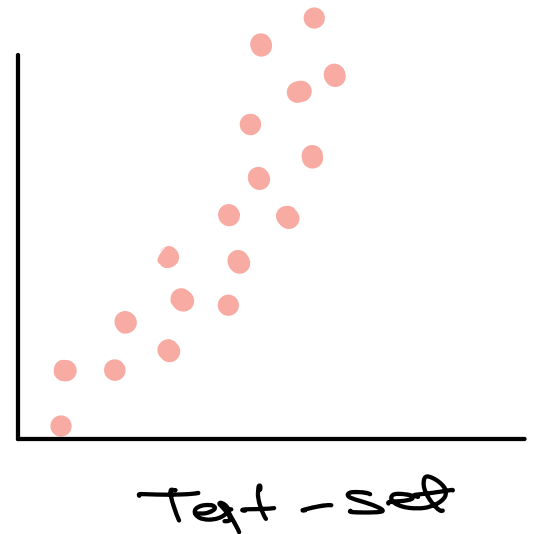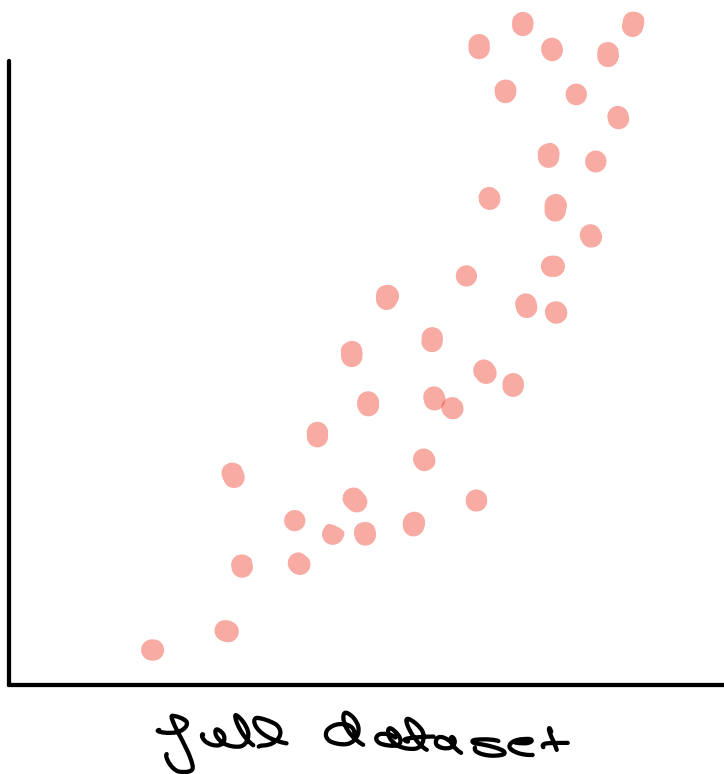On the right plot: $y = w_0 + w_1 x_1 + w_2 x_1^2$

\* What about Multi-collinearity?

✗   $x_2 = \alpha x_1 + \beta$   ← Linear Relationship

✓   $x_2 = x_1 \times x_1$   ← Non Linear

Full dataset

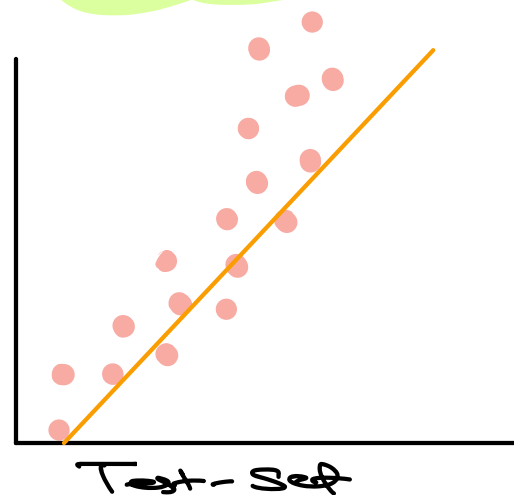→

Train Set

Test - set

Model - 1

Underfitting

Train Set

$r2\text{-}score = 0.65$

Test - Set

$r2\text{-}score = 0.60$

Model 2

Train Set

$r2$-score $= 0.90$

$r2$-score $= 0.85$

Test - Set

Model 3

Good Fit

Train Set

$r2$-score $= 0.85$

$r2$-score $= 0.84$

Test - Set

5) M3

Simplest
↑
|
|
|
|
↓
Complex

| Model | Training | Testing |
|---|---|---|
| Underfit | bad | bad |
| Perfect | Decent | Decent |
| Overfit | Excellent | Decent or bad |

Principles for picking the Best Model

Generalization

① Pick model which performs the best on Unseen Data (Generalizes)

Occam's Razor

② If There are many Solutions Pick the Simplest One.

M2 ↪ M3
(Complex)  (Simple)

* Higher the degree More Complex the Model

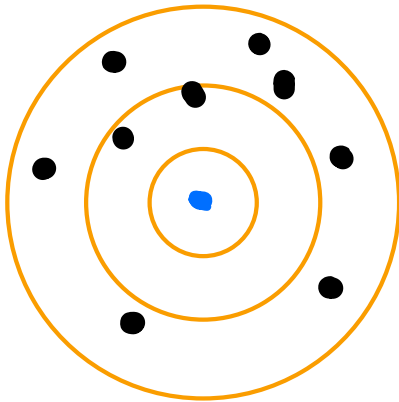(Two types of Error)

Low ← → High

Bias ← → Variance

Vivek is preparing for Olympic archery Competition

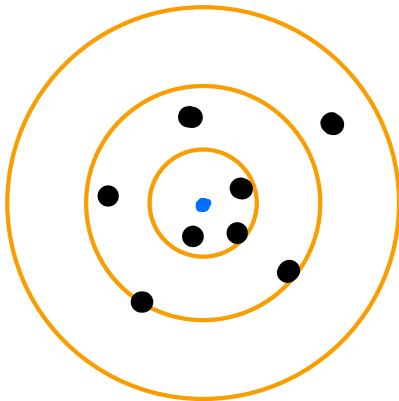## Case-1



→ High Bias ⇒ The shots are far away from Bull's eye
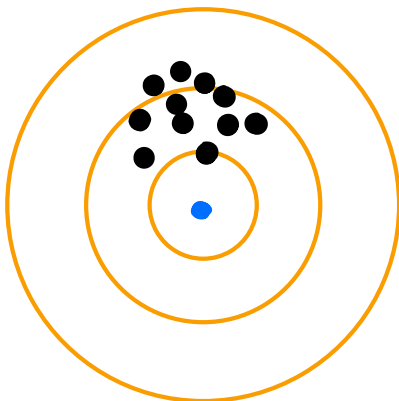
→ High Variance ⇒ Not Consistent (all over the board)
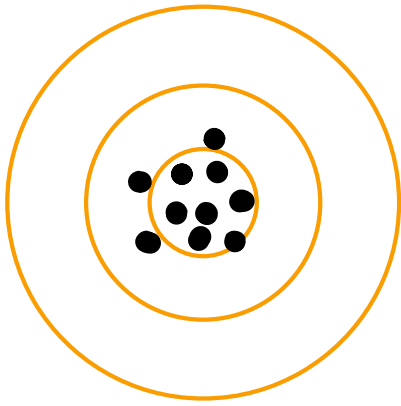
## Case-2



Low Bias

High Variance

## Case 3:



High Bias

Low Variance (Consistent)

# Case-4



Low Bias
Low Variance

* Conclusion:

1) High bias leads to high Errors

2) High Variance leads to high Errors

## Bias vs Variance Tradeoff

—Bias
— Variance



← Good Spot

Prediction Error

Complexity Degree (3)