

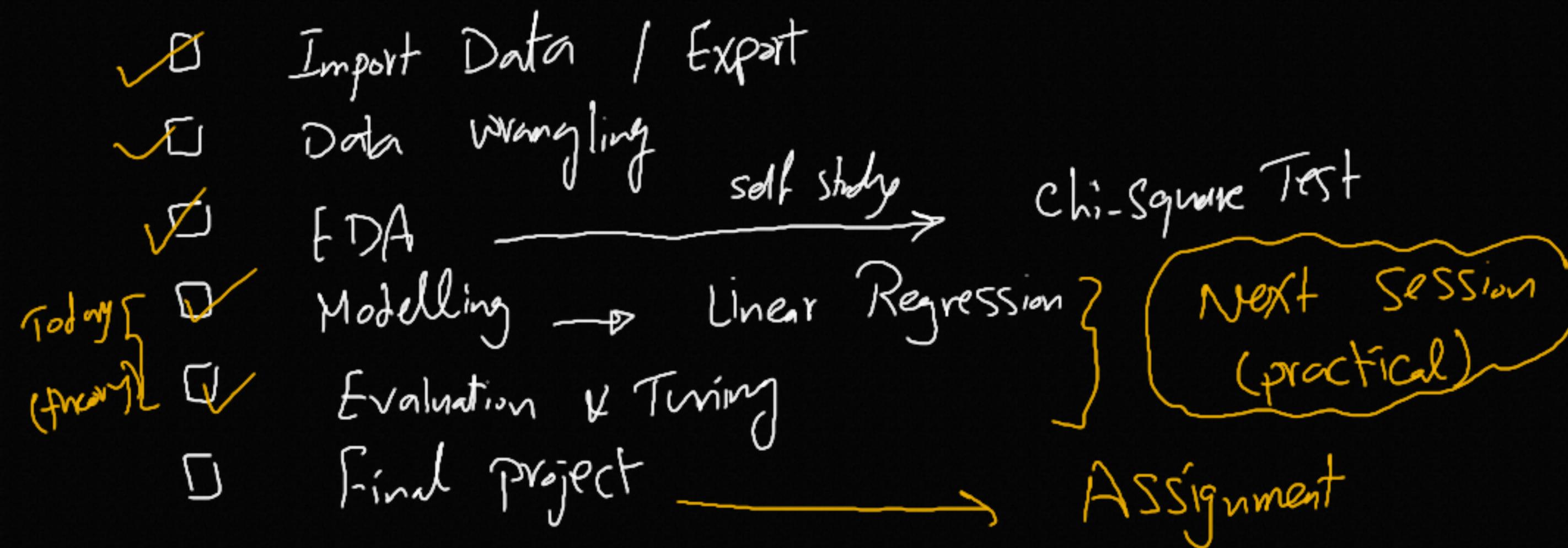
Exploratory Data Analysis

EDA

(Data Wrangling)

- ① Read Data
- ② Explore Data
 - pandas profiling
 - ydata profiling
- ③ Duplicated values, Missing values, outliers, Scaling, Categorical Handling
- ④ Correlation, Visualization → Histogram, Boxplot
 - ↳ Barplot, Countplot
 - ↳ Hue
- ⑤ Feature Engineering :
mpg → L / Kg , lat, long pickup distination
distance
↳ pricing
- ⑥ Communication → Conclusion | Insights
↳ output

- * Search for "Steps of EDA"
 - * Car price
 - * laptop price
 - * Insurance Cost
 - * ... other
- * Data Analysis with python



Business
Model

Modelling

→ Relation
↓

$$y = f(x)$$



① $y = 5x$

② $y = 3x^2$

③ $y = e^x$

dependent var.
independent var.

...
Many
Many
Relations

$$\text{Speed} = \frac{\text{Distance}}{\text{Time}}$$

Known Relation
R ↗

$$V = IR$$

Data Modelling

$$y = f(x)$$

Target

Label

Output

Features

Input

Linear

Model

Car Data

features

city year Brand Km mpg fuel type

- - - - - - - -

$$\text{price} = B_1 \leftarrow \begin{matrix} 100000 \\ \text{Brand} \end{matrix} + B_2 \leftarrow \begin{matrix} 20 \rightarrow 300K \\ 20000 \\ 300 - \text{Suv} \\ 500000 \end{matrix} \leftarrow \begin{matrix} \text{gas} \\ 20000 \\ \text{fuel} \\ \text{type} \\ \text{Electric} \\ 50000 \end{matrix}$$

= $\frac{200000}{\equiv}$

Linear Model

one variable



$$y = 5x \rightarrow$$

$$\text{price} = 50 \times \text{horse power}$$

Target

one feature

$$y = 50000 + 5x \rightarrow$$

$$\text{price} = \frac{100000}{\text{Const}} + 25 \times \text{hp}$$

Const

$$\left. \begin{array}{l} \text{Salary} = \frac{6000}{\text{Const}} + \dots \\ \text{Const} = \dots \end{array} \right\}$$



$$y = b_0 + b_1 x_1 + b_2 x_2 + \dots + b_n x_n$$

find
coeff.

multiple
Features



$$y = b_0 + b_1 x_1 + b_2 x_1^2 + b_3 x_2 + b_4 x_1 x_2 + b_5 x_2^2 - \dots$$

X	y
- -	- -
- -	- -
- -	- -
- -	- -

Historical
Data

Polynomial

$$Y = f(X)$$

- ↳ Simple Linear Regression
- ↳ Multiple Linear Regression
- ↳ Polynomial Regression

$$y = b_0 + b_1 X$$

$$y = b_0 + \sum_{i=1}^n b_i X_i$$

↓
①

$b_0 X_0$

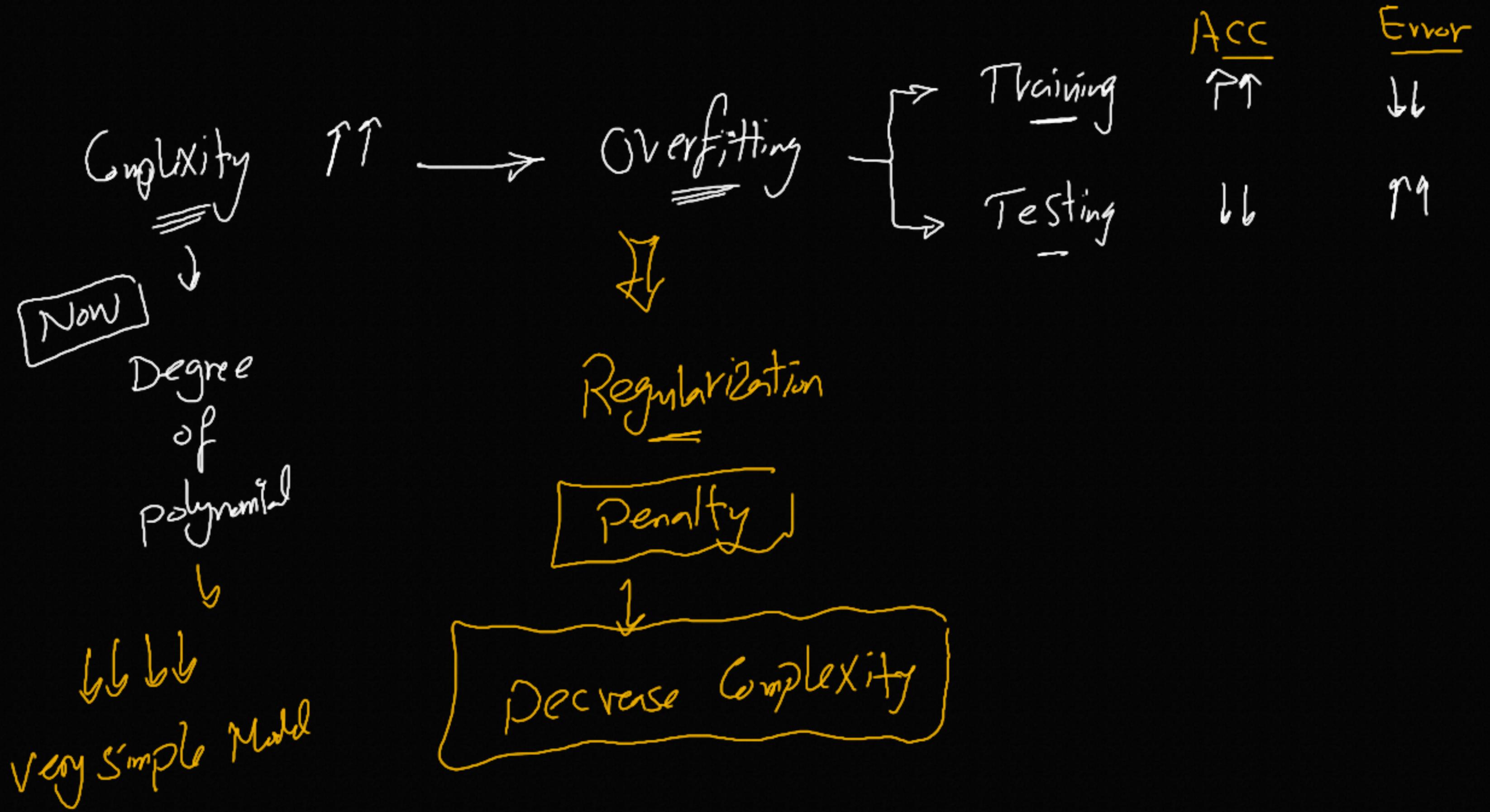
$$y = \sum_{i=0}^n b_i X_i$$

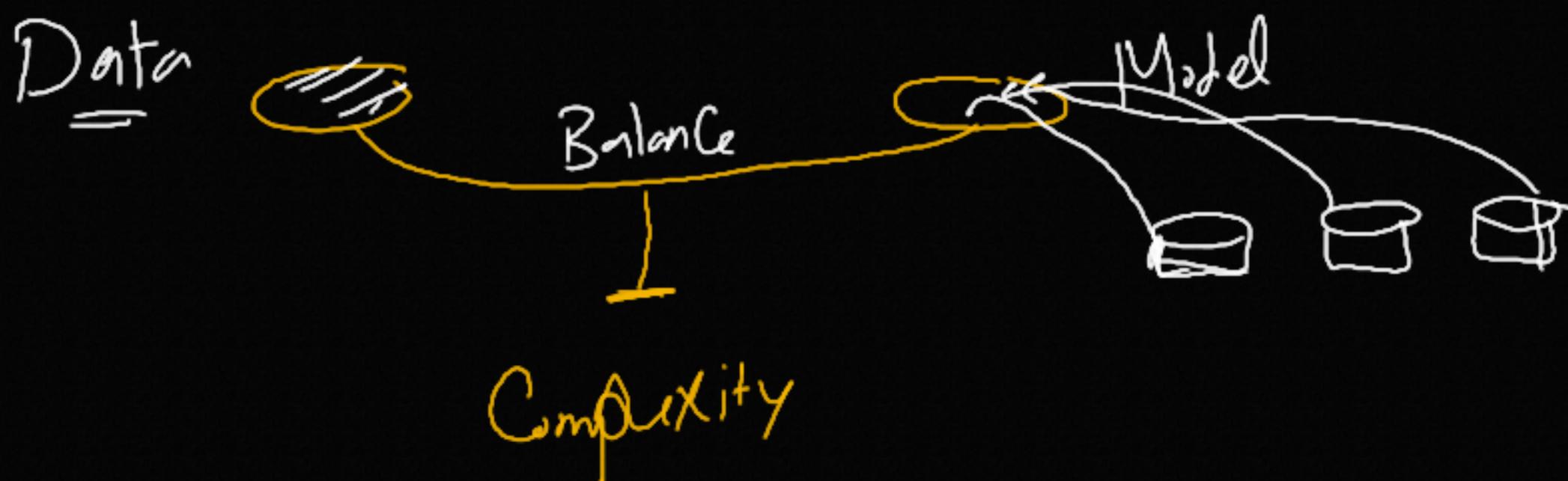
$$X_1 X_2^2$$

Degrees ↑

$$X_1^2 X_2$$

Complexity ↑





Complexity → Model >> Data → Overfitting
Complexity → Model << Data → Underfitting
Complexity → Model ~ Data → Good fit

→ High Correlation

→ Linear Model

↳ Best fit Model

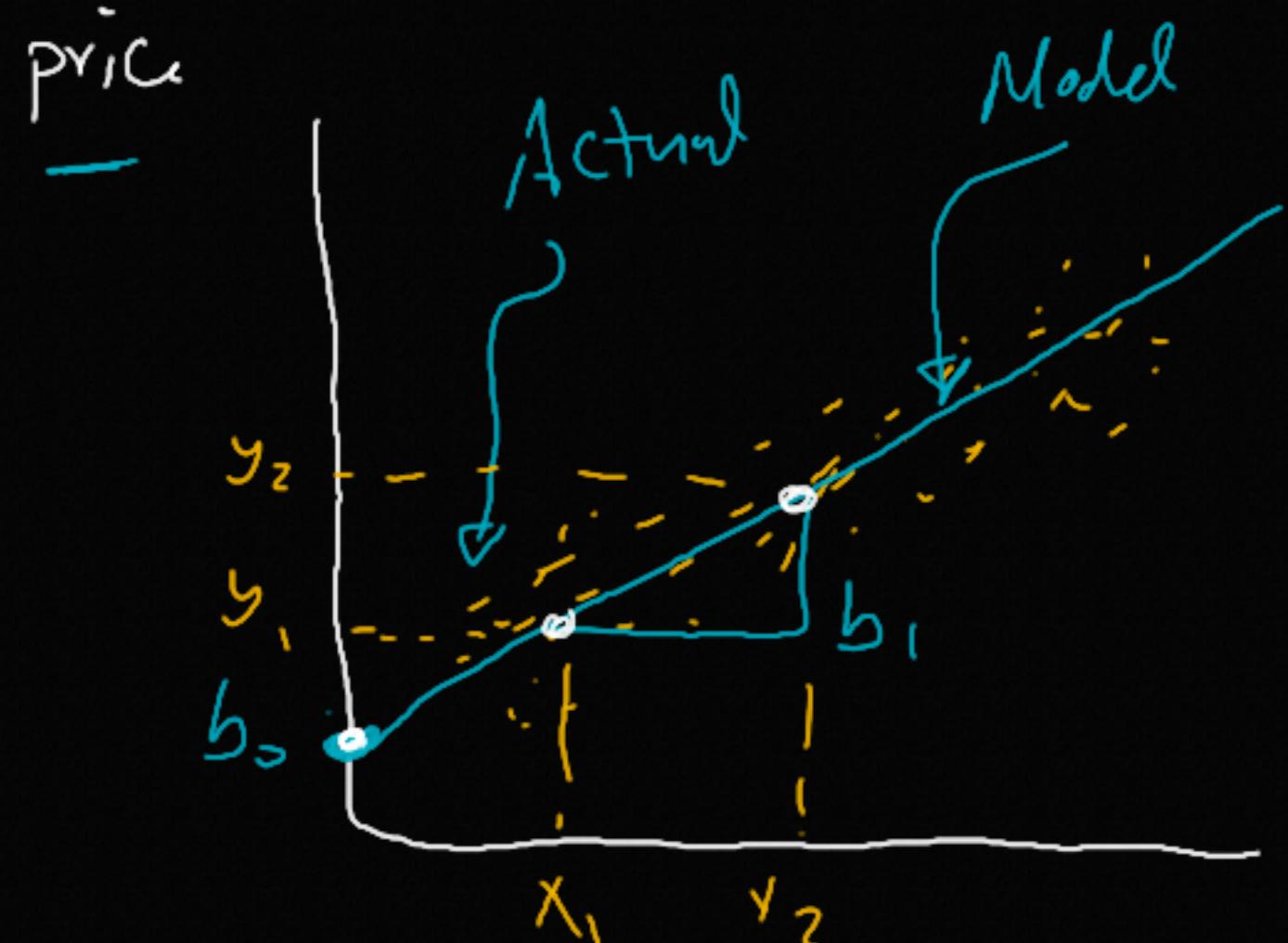
→ Actual & Model (prediction)

→ Measure Error

→ Min. Error

Every line → b_0, b_1

↳ Best fit line



$$\text{price} = b_0 + b_1 \text{ hp}$$

↑ ↑
intercept Geff
(Bias) (Slope) $\frac{y_2 - y_1}{x_2 - x_1}$

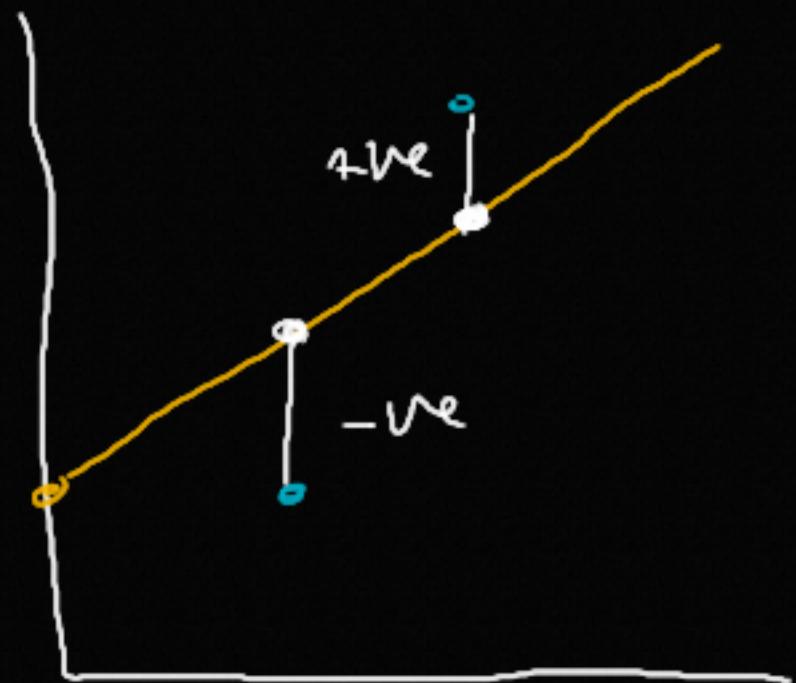
Get Parameters b_0, b_1

→ Min. Error

Cost fn

$$\frac{1}{m} \sum_{i=1}^m (y_{\text{pred}} - y_{\text{actual}})^2$$

Squared Error



↳ Mean Squared Error

↑
Minimize

MSE

L1 200

L2 50

L3 13 ✓

Get parameters
(weights)

Analytical
—
Close loop

Normal Equation vectorized

$$\rightarrow y = f(w, x)$$

$$\rightarrow w = f(y, x)$$

$\{ \}$
known

Numerical
open loop

Gradient Descent

→ Try α error (with steps)

→ Set weights random

→ Error

→ change weights $\begin{cases} \nearrow \\ \searrow \end{cases}$

→ Error

→ change weights \rightarrow 

Gradient Descent

① Set Random weights

new weight \rightarrow

$f(w)$

$$\textcircled{1} \quad J = \frac{1}{m} \sum_{i=1}^m (\hat{y}_i - y_i)^2$$

$w_0 + w_1 x$

② $\frac{\partial J}{\partial w}$ \rightarrow gradient
 Contribution of weight to error

$$\textcircled{3} \quad w^{\text{new}} = w^{\text{old}} - \alpha \frac{\partial J}{\partial w}$$

definit by
 $w^{\text{new}} \leftarrow w^{\text{old}} - \text{learning rate} \cdot \frac{\partial J}{\partial w}$

Normal Equation

x

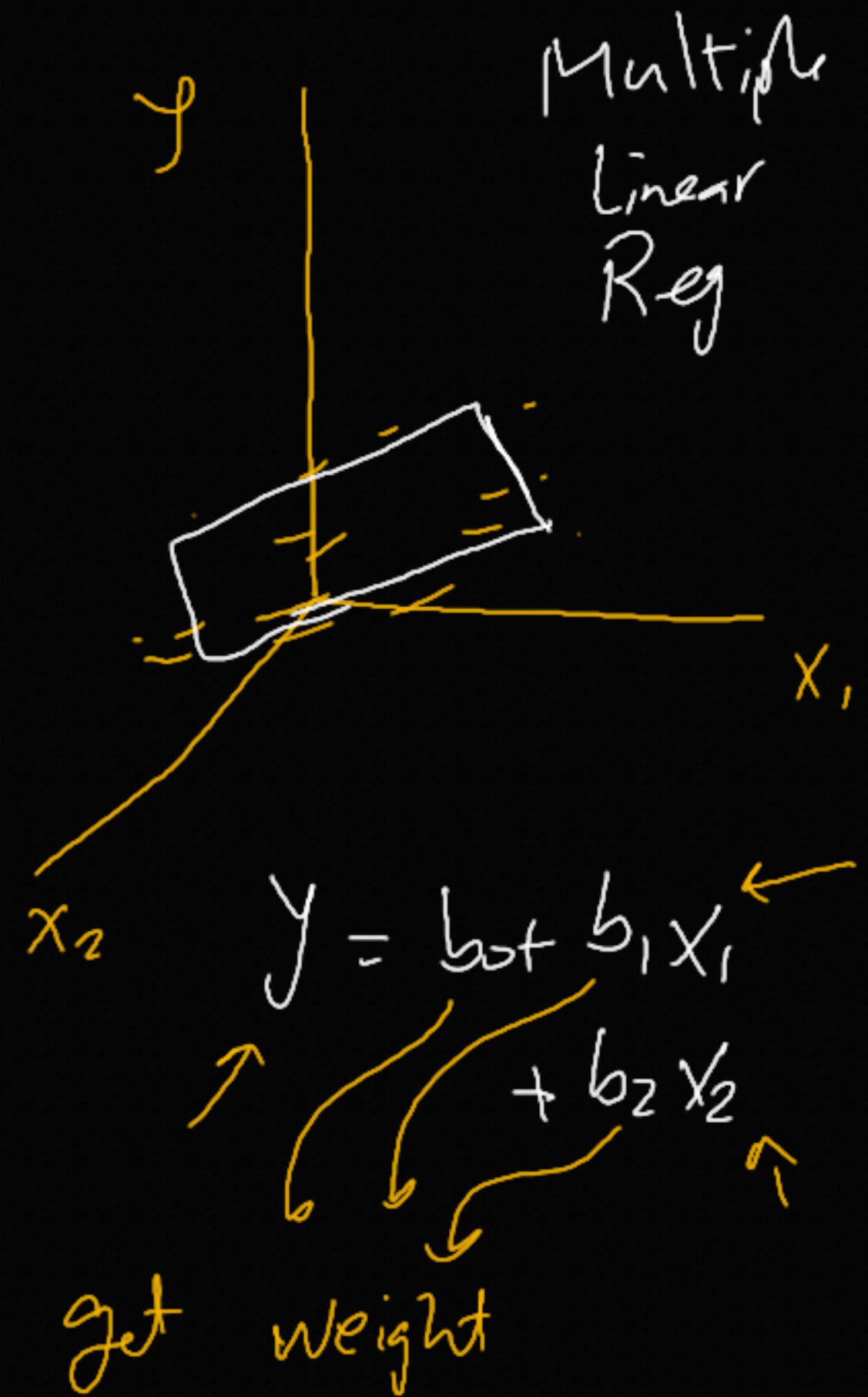
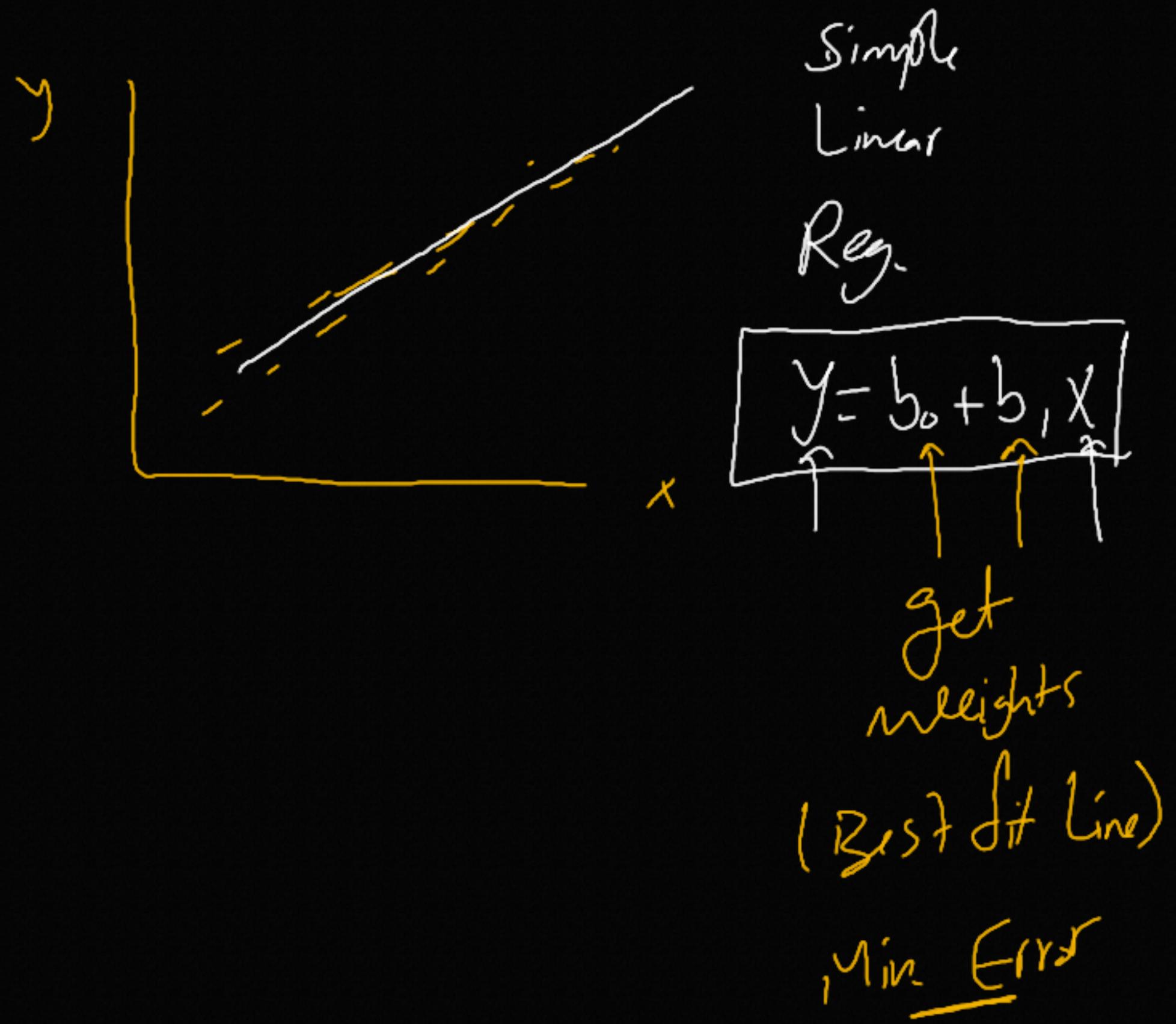
y

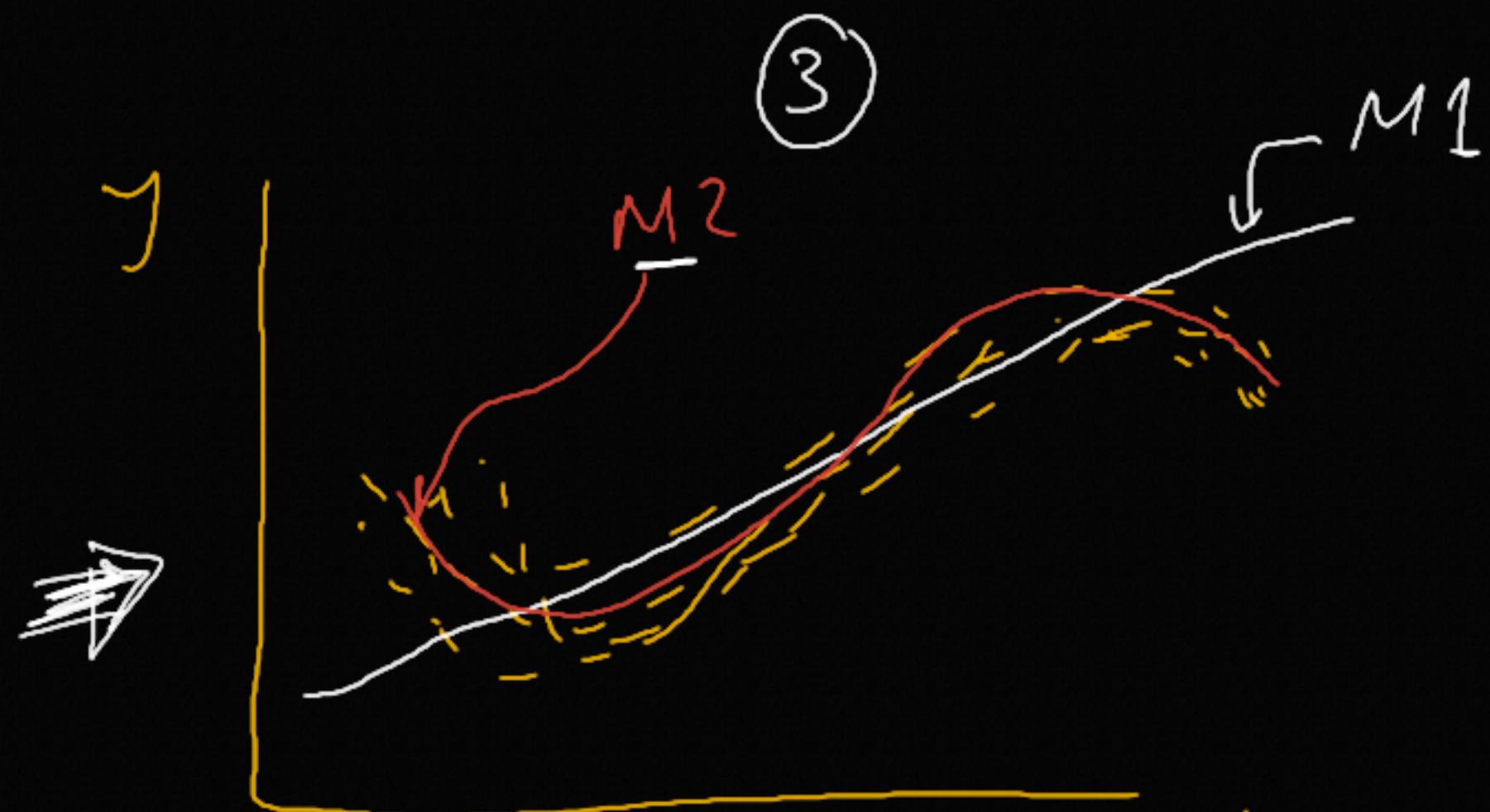
$$\begin{bmatrix} & & & \\ & & & \\ & & & \\ & & & \\ & & & \end{bmatrix} \quad \begin{bmatrix} & \\ & \\ & \\ & \\ & \end{bmatrix}$$

$$y = w^T x$$

$$w = (X^T X)^{-1} X^T y$$

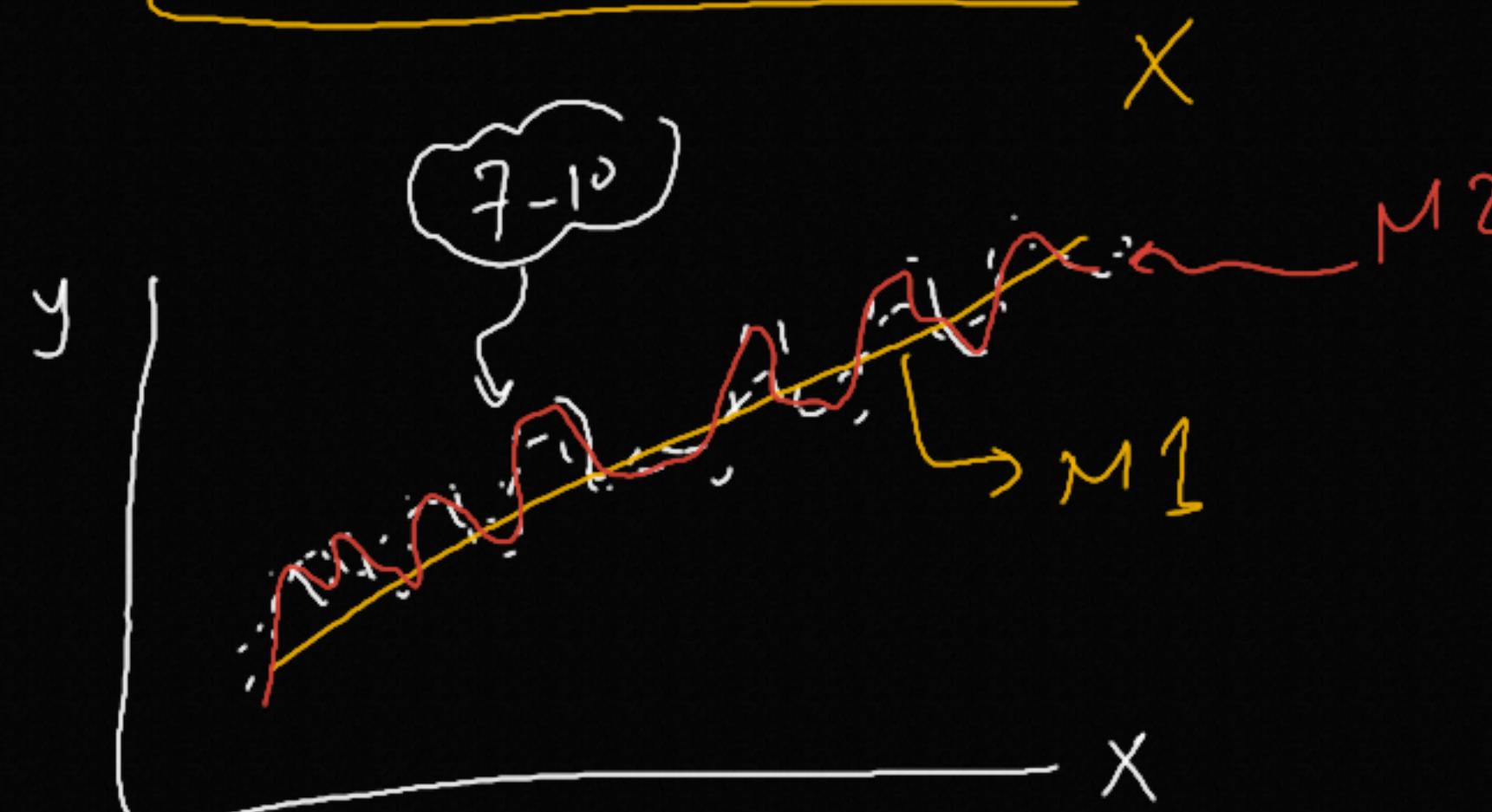
$$\frac{\partial J}{\partial w} = 0 \quad (\text{to prove})$$





Error

M₁ >> M₂ Error



Error

M₁ >> M₂ Error

overfitting

Non Linear Data \longrightarrow Non Linear Model

1

Polynomial Regression

Polynomial features

1

Linear Regression

2

Multivar Linear ~~Regressi~~

$$\begin{array}{ccc} x & y & \rightarrow x & x^2 & y \\ . & . & & . & . \\) & / & &) & / \\ , & (& & , & (\end{array}$$

$$df[x_2] \subseteq df[x] \text{ for } 2$$

$$Y = b_0 + b_1 X_1 + b_2 X_2$$

x²

Linear
Reg
 L_2 penalty

Complexity
(Degree) ↑↑ → overfitting

Control ← Regularization Term

Ridge

L_2 Penalty

Lasso

L_1 penalty

$$\lambda \sum |w_i|$$

$$J = \text{MSE} + \lambda \sum_{i=1}^n w_i^2$$

Regularization

penalty for Higher Degrees

$$w_2^{\text{new}} = w_2^{\text{old}} - \alpha \frac{\partial J}{\partial w_2}$$

Learning rate

$$y = b_0 + b_1 x + b_2 x^2 + b_3 x^3 + b_4 x^4 + b_5 x^5$$

0.00001

stop

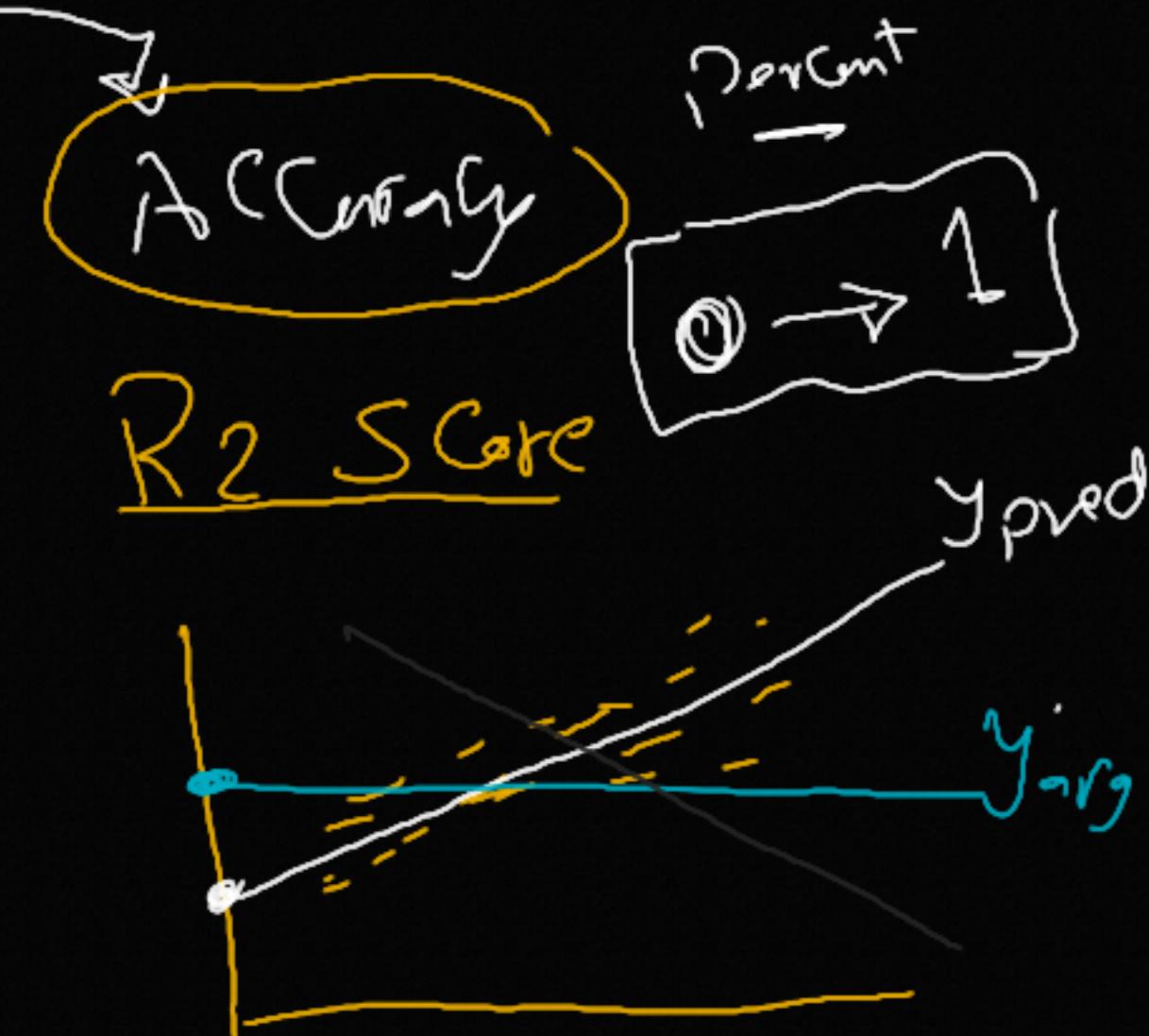
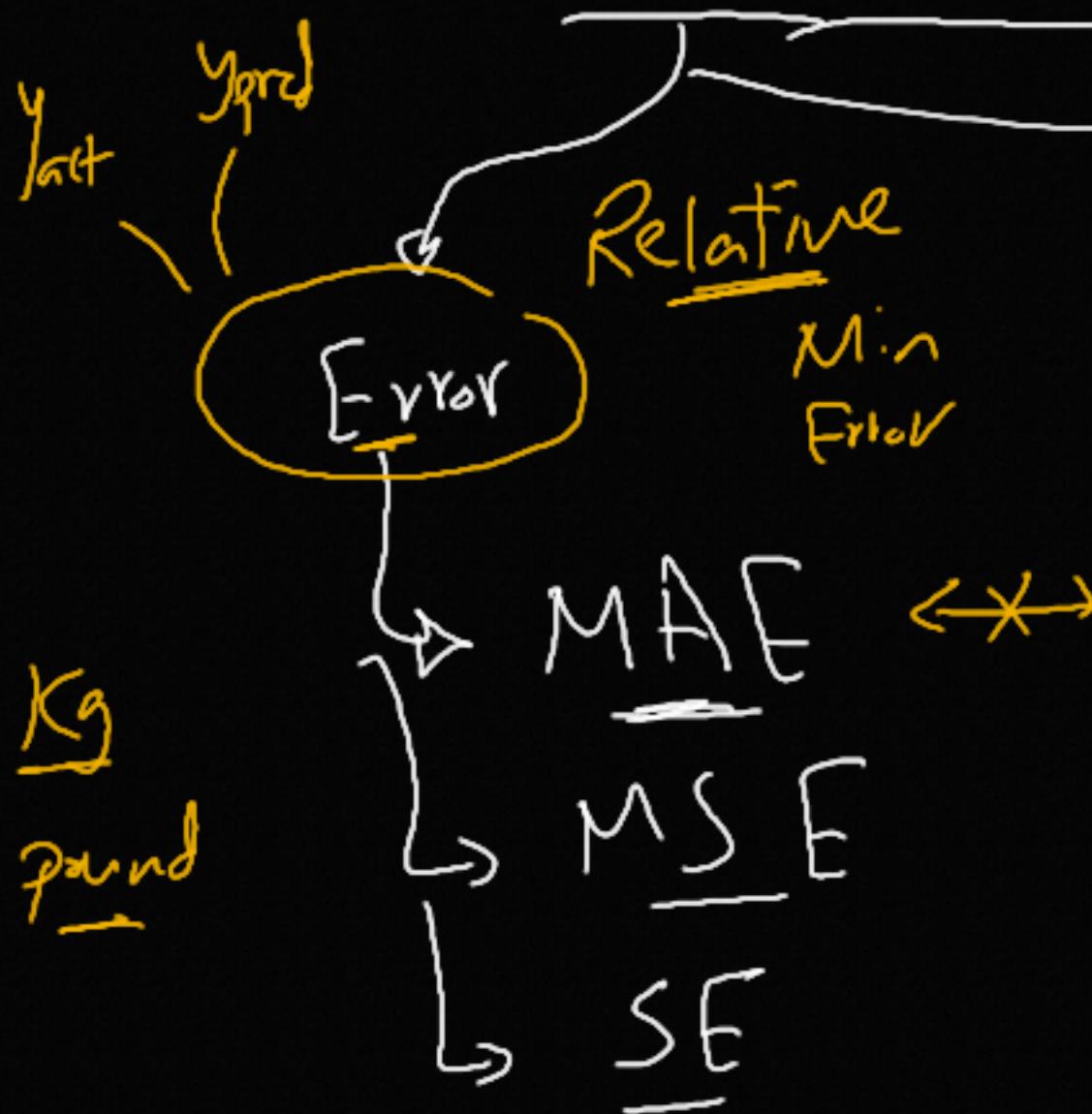
L_1, L_2 penalty \rightarrow Elastic Net
Linear Reg + L_1, L_2 penalty

- Linear Regression
- Evaluation & Tuning

$\boxed{100,000}$ Normal Eqn

(α) Gradient Desc

Evaluation on Regression Models



$$R^2 = 1 - \frac{\sum (y_{pred} - y_{act})^2}{\sum (y_{act} - \bar{y}_{act})^2}$$

Model Error

M1 M2 M3

MAE₍₁₎ MAE₍₂₎ MAE₍₃₎

MSE₍₁₎ MSE₍₂₎ MSE₍₃₎

SG
10
15
12.5

10000
15000
125000

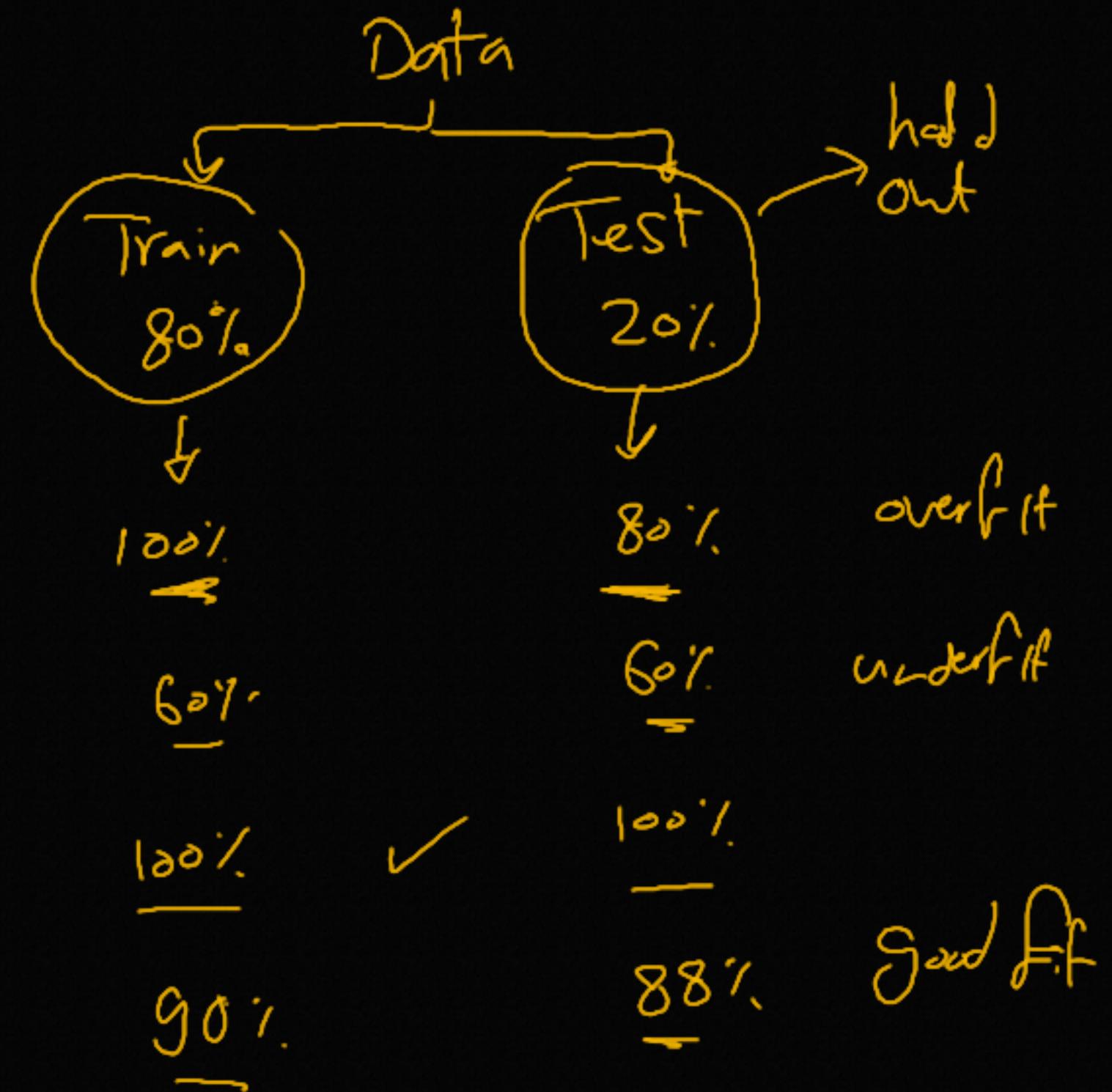
Tuning

$$LR, g_d = f(\alpha)$$

learning
Rate

		Hyperparameter
M1	$\alpha = 0.001$	0.80
M2	$\alpha = 0.01$	0.7 Tuning
M3	$\alpha = 0.1$	0.75
M4	$\alpha = 1$	0.65

Same Model



Ridge \rightarrow L₂ penalty (λ)