# REGRESSION

Lecture 8

MAL1, 2024

1

MACHINE
LEARNING

CLUSTERING

CLASSIFICATION

UNSUPERVISED LEARNING

SUPERVISED LEARNING

DIMENSIONALITY REDUCTION

MACHINE LEARNING

REGRESSION

REINFORCEMENT LEARNING

CLASSIFICATION

SUPERVISED LEARNING
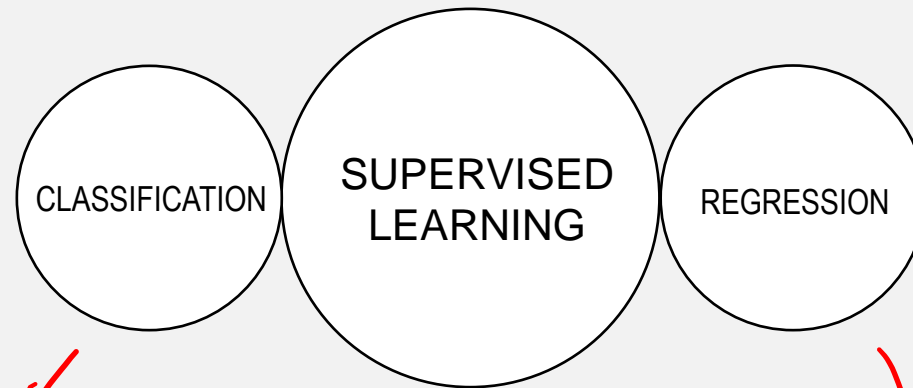
REGRESSION

response variable is a class

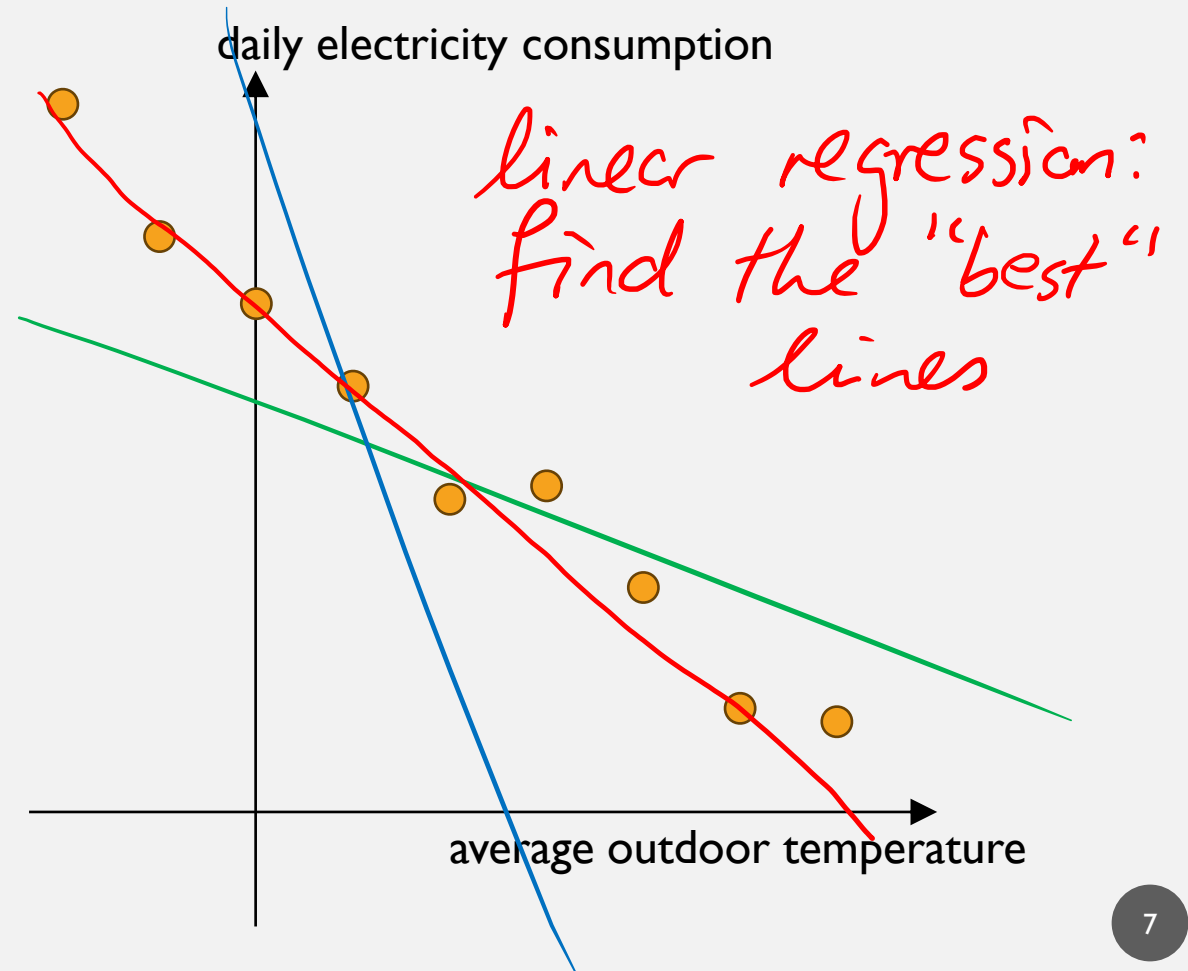response variable is a number

# REGRESSION

- **Linear regression**
- Performance metrics
- Polynomial regression
- Regularization

# REGRESSION

| average outdoor temperature (°C) | daily electricity consumption (kWh) |
|:---:|:---:|
| -10 | 46.5 |
| -5 | 37.9 |
| 0 | 33.2 |
| 5 | 27.5 |
| 10 | 20.3 |
| 15 | 21.1 |
| 20 | 14.2 |
| 25 | 6.3 |
| 30 | 5.6 |

feature

response variable

daily electricity consumption

linear regression: find the "best" lines

average outdoor temperature

# REGRESSION

| $x_1$ | y |
|---|---|
| $\theta_0 + -10 \times \theta_1 \approx$ | 46.5 |
| $\theta_0 + -5 \times \theta_1 \approx$ | 37.9 |
| $\theta_0 + 0 \times \theta_1 \approx$ | 33.2 |
| $\theta_0 + 5 \times \theta_1 \approx$ | 27.5 |
| $\theta_0 + 10 \times \theta_1 \approx$ | 20.3 |
| $\theta_0 + 15 \times \theta_1 \approx$ | 21.1 |
| $\theta_0 + 20 \times \theta_1 \approx$ | 14.2 |
| $\theta_0 + 25 \times \theta_1 \approx$ | 6.3 |
| $\theta_0 + 30 \times \theta_1 \approx$ | 5.6 |

with $\hat{y} = \theta_0 + \theta_1 x_1$
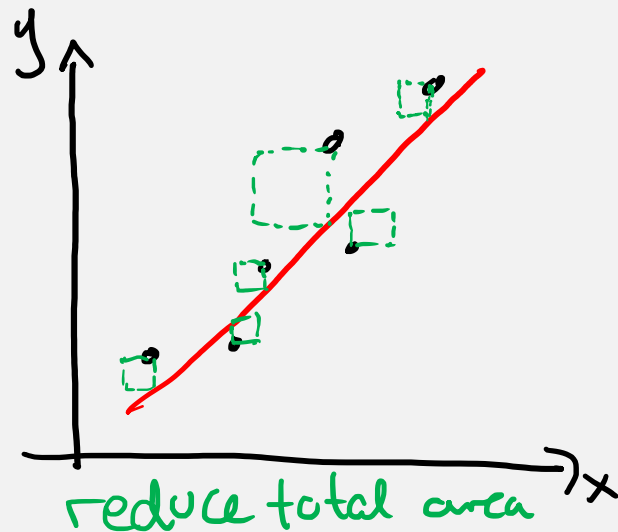find the "best" values of
$\theta_0$ and $\theta_1$

$\theta_0 = 33.72$     $\theta_1 = -1.009$

finding these numbers
=
training the model

$$\hat{y} = \Theta_0 + \Theta_1 x_1 + \Theta_2 x_2 + \Theta_3 x_3 + \cdots = \Theta^T x$$

prediction

features

matrix notation

The "best" line reduces/minimizes

$$SSE = \sum_i \left( y^{(i)} - \hat{y}^{(i)} \right)^2$$

sum of squared errors

observation     prediction

reduce total area

# FINDING THETA'S

$$\text{SSE} = \sum_i \left(y^{(i)} - \hat{y}^{(i)}\right)^2 = \sum_i \left(y^{(i)} - \theta^T x^{(i)}\right)^2$$

minimize → take the derivative wrt. all $\theta$'s and set equal to zero:

$$\frac{\partial}{\partial \theta_j} \text{SSE} = \frac{\partial}{\partial \theta_j} \sum_i \left(y^{(i)} - \theta^T x^{(i)}\right)^2 = 2 \sum_i \left(y^{(i)} - \theta^T x^{(i)}\right) x_j^{(i)} = 0$$

summarize in matrix form:

$$2\boldsymbol{X}^T(\boldsymbol{X}\boldsymbol{\theta} - \boldsymbol{y}) = 0$$

$$X^T X \theta = X^T y$$

$$\theta = (X^T X)^{-1} X^T y \qquad \text{"normal equation"}$$

$$\boldsymbol{\theta} = (\boldsymbol{X}^T\boldsymbol{X})^{-1}\boldsymbol{X}^T\boldsymbol{y}$$

$$X = \begin{bmatrix} 1 & -10 \\ 1 & -5 \\ 1 & 0 \\ 1 & 5 \\ 1 & 10 \\ 1 & 15 \\ 1 & 20 \\ 1 & 25 \\ 1 & 30 \end{bmatrix} \quad y = \begin{bmatrix} 46.5 \\ 37.9 \\ 33.2 \\ 27.5 \\ 20.3 \\ 21.1 \\ 14.2 \\ 6.3 \\ 5.6 \end{bmatrix}$$

$$X^T = \begin{bmatrix} 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 \\ -10 & -5 & 0 & 5 & 10 & 15 & 20 & 25 & 30 \end{bmatrix}$$

$$X^TX = \begin{bmatrix} 9 & 90 \\ 90 & 2400 \end{bmatrix}$$

$$(X^TX)^{-1} = \begin{bmatrix} \dfrac{8}{45} & -\dfrac{1}{150} \\ -\dfrac{1}{150} & \dfrac{1}{1500} \end{bmatrix}$$

$$(X^TX)^{-1}X^T = \begin{bmatrix} \dfrac{11}{45} & \dfrac{19}{90} & \dfrac{8}{45} & \dfrac{13}{90} & \dfrac{1}{9} & \dfrac{7}{90} & \dfrac{2}{45} & \dfrac{1}{90} & -\dfrac{1}{45} \\ -\dfrac{1}{75} & -\dfrac{1}{100} & -\dfrac{1}{150} & -\dfrac{1}{300} & 0 & \dfrac{1}{300} & \dfrac{1}{150} & \dfrac{1}{100} & \dfrac{1}{75} \end{bmatrix}$$

*pseudoinverse*

$$\boldsymbol{\theta} = (\boldsymbol{X}^T\boldsymbol{X})^{-1}\boldsymbol{X}^T\boldsymbol{y} = \begin{bmatrix} 33.72 \\ -1.009 \end{bmatrix}$$

← $\theta_0$ intercept
← $\theta_1$ slope

| $x_0$ | $x_1$ | y |
|---|---|---|
| 1 | -10 | 46.5 |
| 1 | -5 | 37.9 |
| 1 | 0 | 33.2 |
| 1 | 5 | 27.5 |
| 1 | 10 | 20.3 |
| 1 | 15 | 21.1 |
| 1 | 20 | 14.2 |
| 1 | 25 | 6.3 |
| 1 | 30 | 5.6 |

X          y

$O(n^2)$ instead of $O(n^3)$

in practice, the pseudoinverse is computed using SVD

# REGRESSION

- Linear regression
- **Performance metrics**
- Polynomial regression
- Regularization

# SSE AND FRIENDS
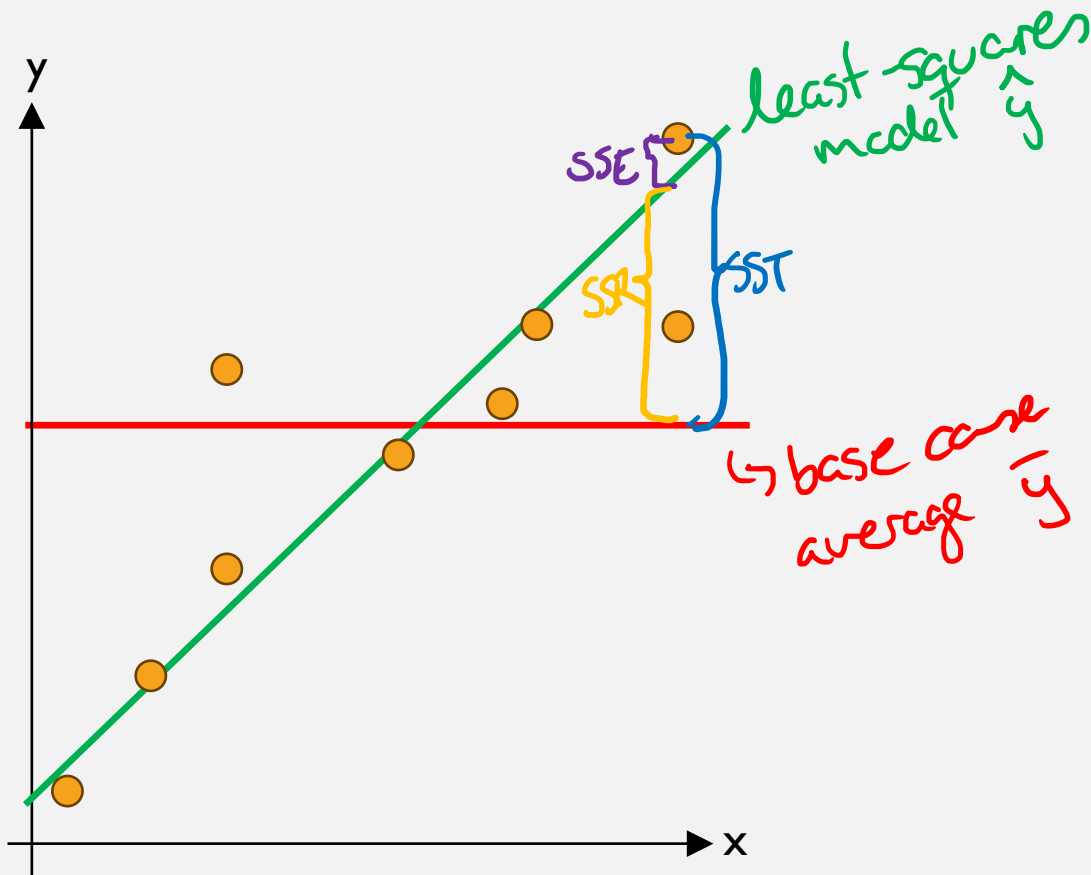
- The smaller the SSE, the better the model …

but SSE depends on # of data points

$\rightarrow$ MSE $= \frac{1}{n}$ SSE          mean squared error

but MSE depends on scale of response variable

… so what should we use as our performance metric?

# SSE AND FRIENDS



$$SST = \sum_i (y_i - \bar{y})^2$$

total deviation from mean
(total sum of squares)

$$SSE = \sum_i (y_i - \hat{y}_i)^2$$

unexplained

$$SSR = \sum_i (\hat{y}_i - \bar{y})^2$$

explained (sum of squares due to regression)

least squares model $\hat{y}$

base case average $\bar{y}$

$$SST = \sum_i \left( y^{(i)} - \bar{y} \right)^2 \qquad \text{total deviation from mean}$$

$$SSE = \sum_i \left( y^{(i)} - \hat{y}^{(i)} \right)^2 \qquad \text{unexplained part}$$

$$SSR = \sum_i \left( \hat{y}^{(i)} - \bar{y} \right)^2 \qquad \text{explained part}$$

$$r^2 = \frac{SSR}{SST} = \frac{\text{"explained"}}{\text{"total"}}$$

the amount of variance the model is able to explain

$r^2 \leq \boxed{1} - $ perfect predictive model

# CODE EXAMPLE



*Jupyter Notebook* **Regression - Hitters**

# REGRESSION

- Linear regression
- Performance metrics
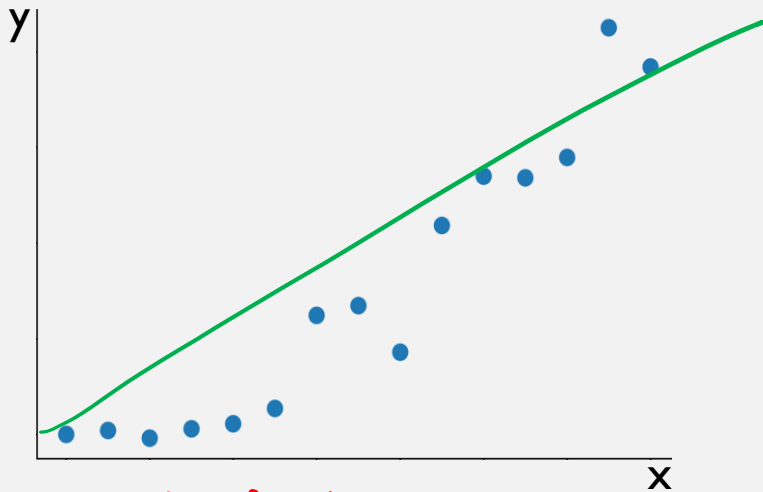- **Polynomial regression**
- Regularization

# POLYNOMIAL REGRESSION

$$\hat{y} = \Theta_0 + \Theta_1 x + \Theta_2 x^2 + \Theta_3 x^3 + \ldots$$

$x^2, x^3, \ldots$ are just new features

$$X = \begin{bmatrix} | & x & x^2 & \\ | & x & x^2 & \\ | & x & x^2 & \\ | & x & x^2 & \cdots \\ | & x & x^2 & \\ | & x & x^2 & \\ \vdots & \vdots & \vdots & \end{bmatrix}$$
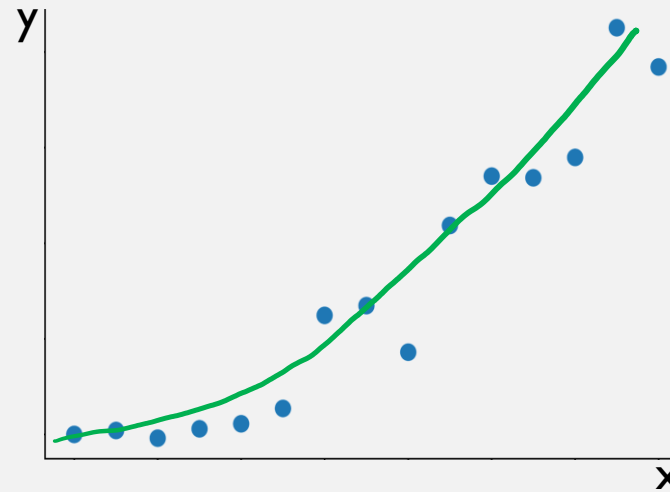
may be a good idea –
depending on the underlying relationship

# UNDERFITTING AND OVERFITTING
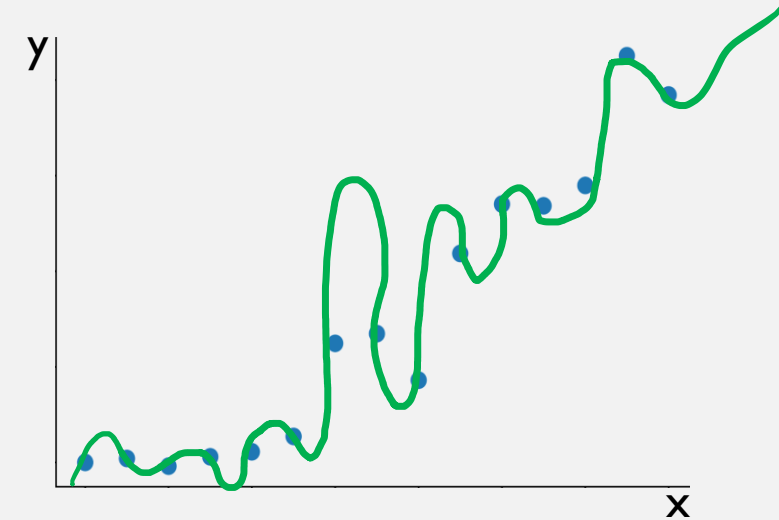


underfitting

high bias
low variance

sweet spot

bias-variance
trade-off

overfitting

low bias
high variance

Bias = "Inability to learn from data"
Variance = "Reliance on data"

19

# REGRESSION

- Linear regression
- Performance metrics
- Polynomial regression
- **Regularization**

# REGULARIZATION

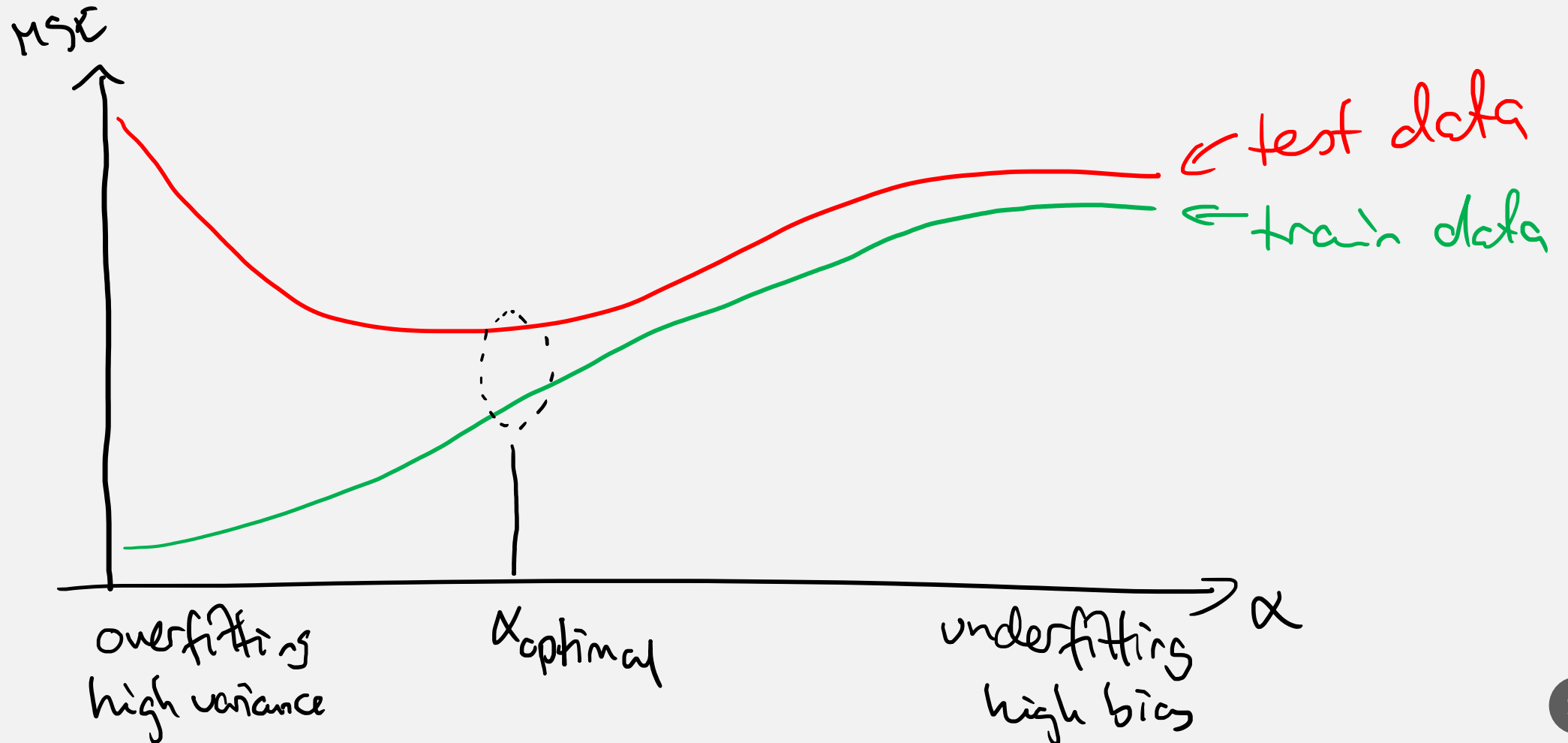Tool to avoid overfitting

Idea: Penalize large coefficient

loss function $\mathcal{L} = MSE + \alpha R(\theta)$

→ regularization (hyperparameter)

$\hookrightarrow$ penalty function of $\theta$

common choices

$\begin{cases} R(\theta) = \sum_i \theta_i^2 & L_2 \text{ regularization} \to \text{Ridge regression} \\ R(\theta) = \sum_i |\theta_i| & L_1 \text{ regularization} \to \text{Lasso regression} \end{cases}$

# THE OPTIMAL REGULARIZATION PARAMETER

# RIDGE VS LASSO REGRESSION

## Ridge

Drives coefficients
to small values
overall

## Lasso

Drives as many
coefficients as
possible to zero

## Elastic Net

combination

$$penalty = \beta \cdot Lasso + (1-\beta) \cdot Ridge$$

# CODE EXAMPLE



*Jupyter Notebook* **Regression - Hitters**