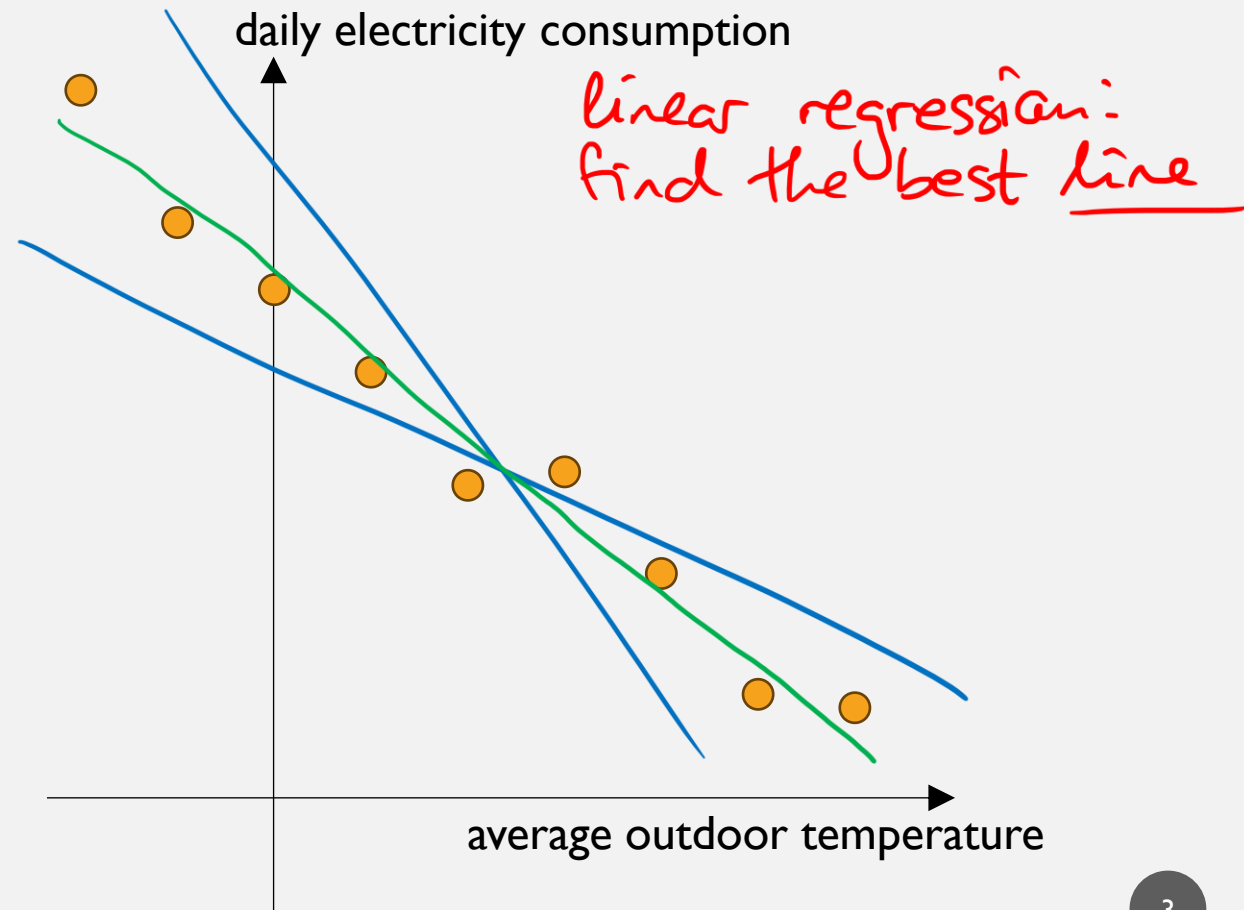# REGRESSION

Lecture 3

MAL1, 2025

1

# 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 line*

average outdoor temperature

# REGRESSION

| $x_1$ | y |
|---|---|
| $\beta_0 + -10 \times \beta_1 \approx$ | 46.5 |
| $\beta_0 + -5 \times \beta_1 \approx$ | 37.9 |
| $\beta_0 + 0 \times \beta_1 \approx$ | 33.2 |
| $\beta_0 + 5 \times \beta_1 \approx$ | 27.5 |
| $\beta_0 + 10 \times \beta_1 \approx$ | 20.3 |
| $\beta_0 + 15 \times \beta_1 \approx$ | 21.1 |
| $\beta_0 + 20 \times \beta_1 \approx$ | 14.2 |
| $\beta_0 + 25 \times \beta_1 \approx$ | 6.3 |
| $\beta_0 + 30 \times \beta_1 \approx$ | 5.6 |

9 equations, 2 unknowns

with $\hat{y} = \beta_0 + \beta_1 x_1$
find the "best" values
of $\beta_0$ and $\beta_1$

$\beta_0 = 33.72$    $\beta_1 = -1.009$

finding these numbers
= training the model

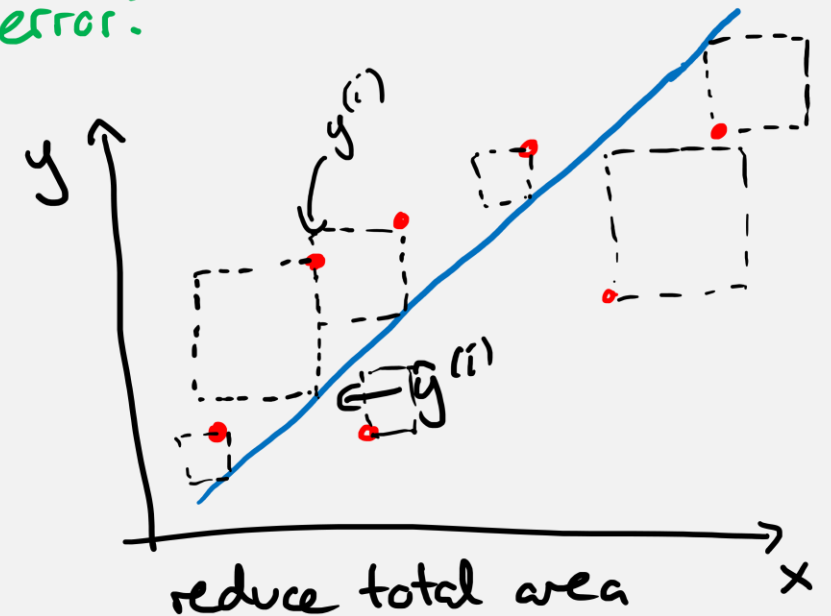$$\hat{y} = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_3 + \cdots = \underbrace{B^T x}_{\text{matrix form}}$$

prediction →

features

The best line reduces the error:

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

"sum of squared errors"



reduce total area

# FINDING BETA'S

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

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

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

summarize in matrix form:

$$2 \cdot 0 = 0 \qquad 2 \overbrace{X^T(XB - y)}^{=0} = 0$$

$$\Rightarrow X^T X B - X^T y = 0$$

$$\Rightarrow X^T X B = X^T y$$

$$\Rightarrow B = (X^T X)^{-1} X^T y \quad \leftarrow \text{"normal equation"}$$

# FINDING BETA'S

$$\hat{y} = \beta_0 \cdot 1 + \beta_1 \cdot x_1$$

new variable
matching $\beta_0$

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

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

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

$$X^T X = \begin{bmatrix} 9 & 90 \\ 90 & 2400 \end{bmatrix} \qquad (X^T X)^{-1} = \begin{bmatrix} \dfrac{8}{45} & -\dfrac{1}{150} \\ -\dfrac{1}{150} & \dfrac{1}{1500} \end{bmatrix}$$

$$X^T y = \begin{bmatrix} 212.6 \\ 612 \end{bmatrix}$$

pseudoinverse

$$B = \underbrace{(X^T X)^{-1} X^T}_{} y = \begin{bmatrix} 33.72 \\ -1.009 \end{bmatrix} \begin{array}{l} \leftarrow \beta_0 \text{ intercept} \\ \leftarrow \beta_1 \text{ slope} \end{array}$$

| $x_0$ | $x_1$ | 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 |

↳ design matrix

in practice, the pseudoinverse is computed using SVD

# THE DESIGN MATRIX

**One variable**

$$\left( \begin{array}{c} \phantom{x} \\ \vdots \\ x^{(i)} \\ \phantom{x} \end{array} \right)$$

**Multiple variables**

$$\left( \begin{array}{ccc} \phantom{x} & \phantom{x} & \phantom{x} \\ \vdots & x^{(i)}_1 & x^{(i)}_2 \; \text{---} \\ \phantom{x} & \phantom{x} & \phantom{x} \end{array} \right)$$

# REGRESSION

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

# SSE AND FRIENDS

- The smaller the SSE, the better the model …

but SSE depends on # data points

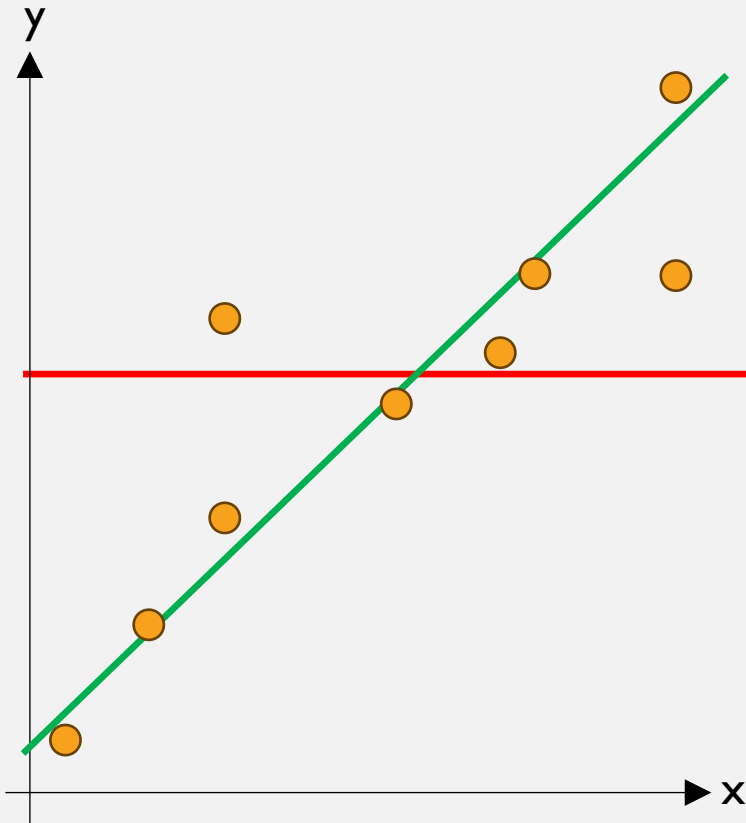$$\rightarrow \quad MSE = \frac{1}{n} SSE \qquad \text{(mean squared error)}$$

but MSE depends on the scale of response variable

… so what should we use as our performance metric?

# SSE AND FRIENDS



least squares model $\hat{y}$

$y$

unexplained part

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

base case average/mean $\bar{y}$

total

$$SST = \sum_i \left(y^{(i)} - \bar{y}\right)^2$$

(total sum of squares)

explained part

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

(regression sum of squares)

SST
SSE
SSR

# SSE AND FRIENDS



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

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

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

performance metric: $r^2$

$$r^2 = \frac{SSR}{SST} = \frac{\text{"explained"}}{\text{"total"}} \leq 1$$

=1 means
→ perfect
predictive
model

↳ the amount of variance is
the model able to explain

# CODE EXAMPLE



*Jupyter Notebook* **Regression - Hitters**

# REGRESSION

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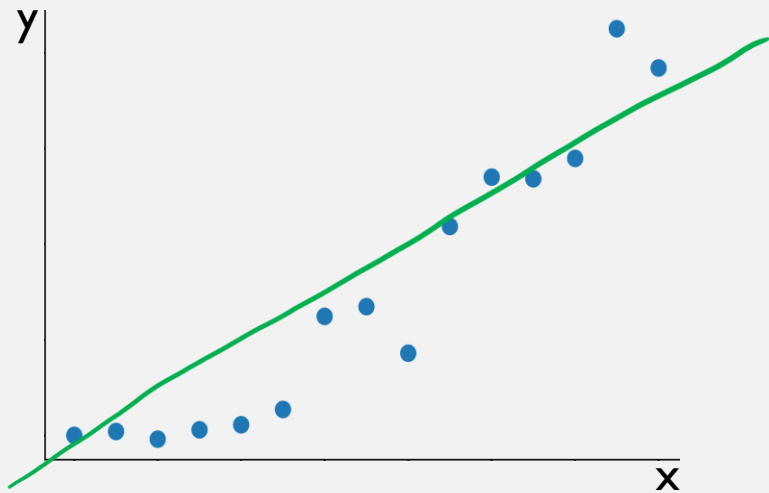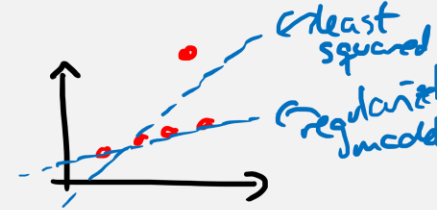- **Polynomial regression**
- Regularization

# POLYNOMIAL REGRESSION

$$\hat{y} = \beta_0 + \beta_1 x + \beta_2 x^2 + \beta_3 x^3 + \dots$$
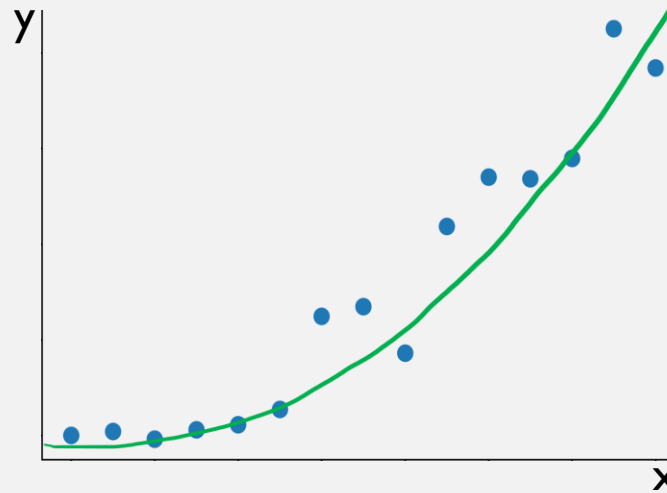
$x^2, x^3$ are just new features

$$X = \begin{bmatrix} | & | & | & | \\ | & x & x^2 & x^3 \dots \\ | & | & | & | \end{bmatrix}$$
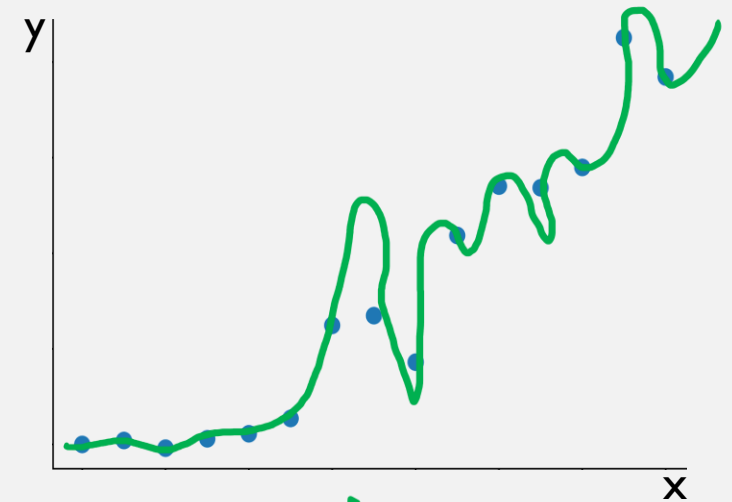
may be a good idea

# UNDERFITTING AND OVERFITTING



underfitting
High bias
Low variance

sweet spot

← bias-variance trade-off →

overfitting
Low biased
High variance

Bias : Inability to learn from data
Variance: Reliance on data

# REGRESSION

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

# REGULARIZATION

Tool to avoid overfitting

Idea: Penalize large coefficients

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

$\rightarrow$ regularization parameter

$\hookrightarrow$ penalty function of $\beta$

$\cup$ $R(\beta) = \sum_i \beta_i^2$ $\leftarrow$ $L_2$ regularization / Ridge regression

$\vee$ $R(\beta) = \sum_i |\beta_i|$ $\leftarrow$ $L_1$ regularization / Lasso regression

# THE OPTIMAL REGULARIZATION PARAMETER

# RIDGE VS LASSO REGRESSION

Ridge

drives coefficients to small values overall

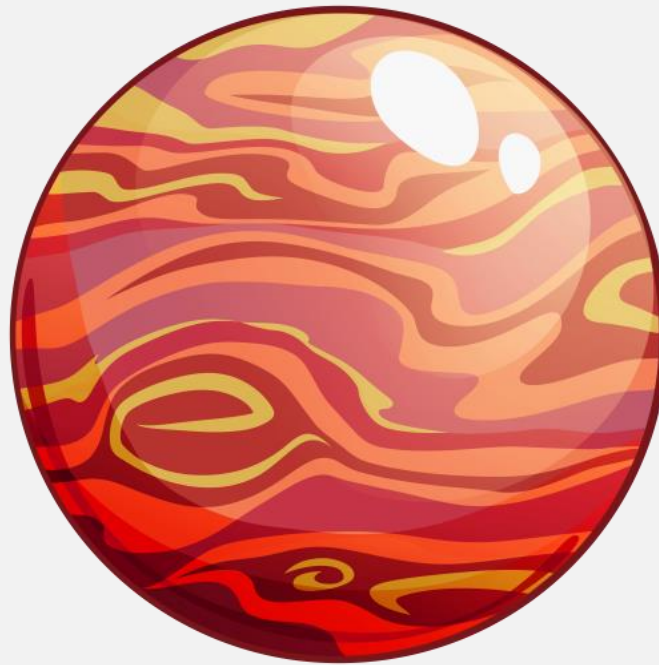Lasso

drives certain coefficients to zero

⟹ built-in feature selection

Elastic Net combines the two:

$$R(\beta) \sim \gamma \cdot \text{Lasso} + (1-\gamma) \cdot \text{Ridge}$$

# CODE EXAMPLE



*Jupyter Notebook* **Regression - Hitters**

- Explain what regression is, including OLS, Ridge, Lasso and Elastic Net regression

- Calculate and interpret relevant performance metrics

- Solve regression problems with sklearn