

# INTRODUCTION TO MACHINE LEARNING COMPSCI 4ML3

LECTURE 4

HASSAN ASHTIANI

# MATRIX FORM OLS

$$\bullet \Delta = \begin{pmatrix} \Delta_1 \\ \dots \\ \Delta_n \end{pmatrix} = \begin{pmatrix} x_1^1 & \dots & x_d^1 \\ \vdots & \ddots & \vdots \\ x_1^n & \dots & x_d^n \end{pmatrix} \begin{pmatrix} w_1 \\ \dots \\ w_d \end{pmatrix} - \begin{pmatrix} y^1 \\ \dots \\ y^n \end{pmatrix}$$

$$\min_{W \in \mathbb{R}^{d \times 1}} \sum_{i=1}^n (\Delta_i)^2 = \min_{W \in \mathbb{R}^{d \times 1}} \|\Delta\|_2^2 =$$

$$\min_{W \in \mathbb{R}^{d \times 1}} \|XW - Y\|_2^2$$

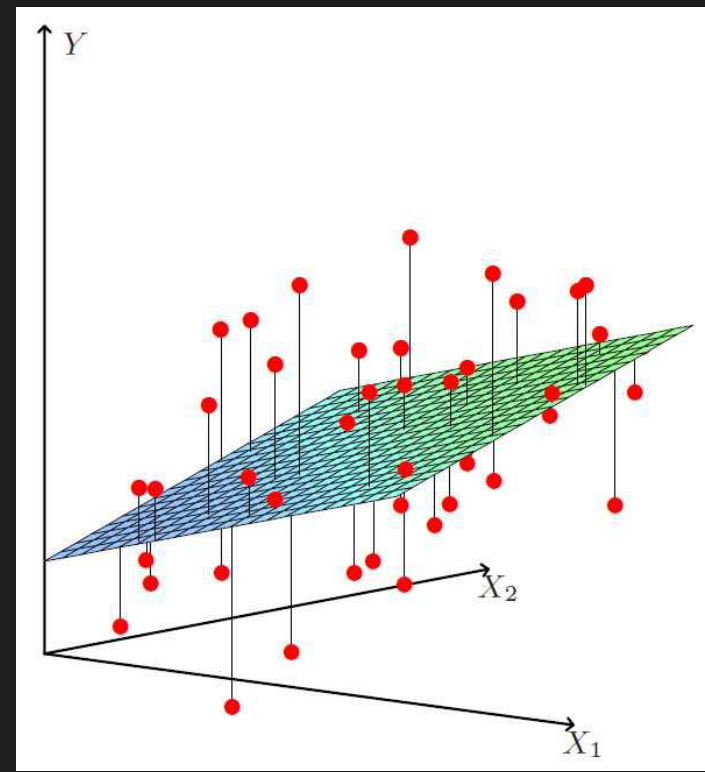
$$W^{LS} = (X^T X)^{-1} X^T Y$$

# BIAS/INTERCEPT TERM

WE ARE MISSING THE BIAS TERM ( $w_0$ )

$$\text{MIN}_{w_0, w_1, \dots, w_d \in \mathbb{R}} \sum_{i=1}^n (w_1 x_1^i + \dots + w_d x_d^i + w_0 - y^i)^2$$

$$\text{MIN}_{w_0 \in \mathbb{R}, W \in \mathbb{R}^{d \times 1}} \left\| XW + \begin{pmatrix} w_0 \\ w_0 \\ \dots \\ w_0 \end{pmatrix} - Y \right\|_2^2$$



# BIAS/INTERCEPT TERM

- ADD A NEW AUXILIARY DIMENSION TO THE DATA

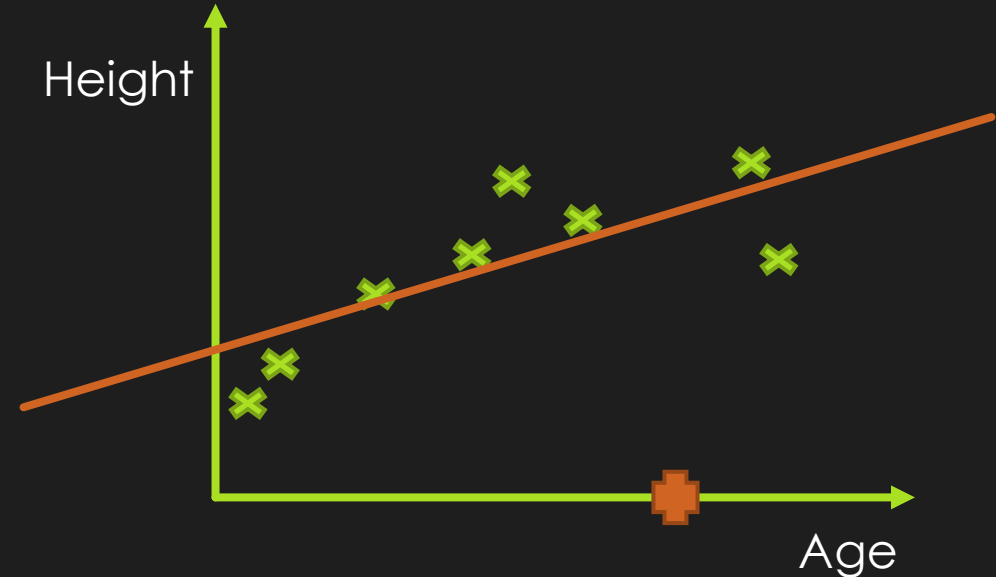
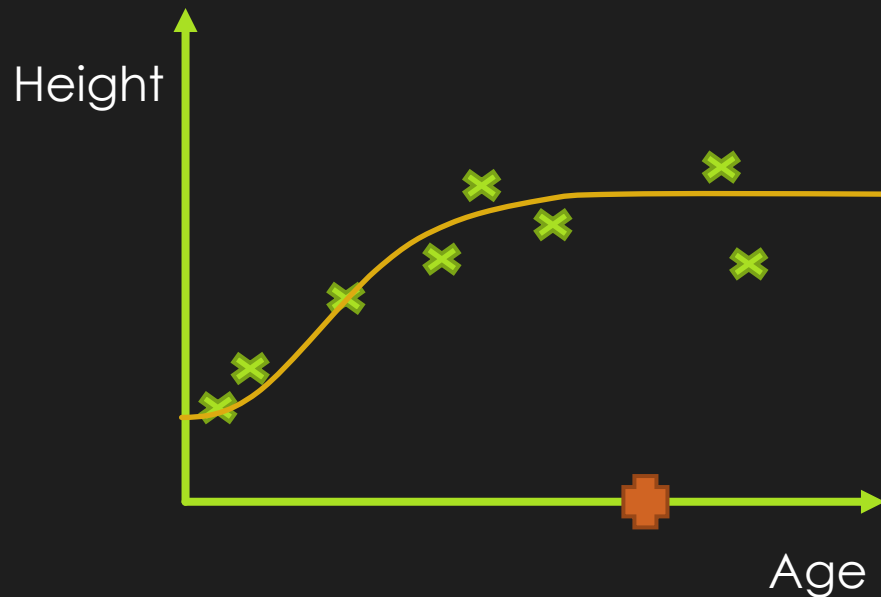
- $$X_{n \times (d+1)} = \begin{pmatrix} x_1^1 & \cdots & x_d^1 & 1 \\ \vdots & \ddots & \vdots & 1 \\ x_1^n & \cdots & x_d^n & 1 \end{pmatrix}, W_{(d+1) \times 1} = \begin{pmatrix} w_1 \\ \vdots \\ w_d \\ w_0 \end{pmatrix}$$

- SOLVE OLS: 
$$\min_{W \in \mathbb{R}^{(d+1) \times 1}} \|XW - Y\|_2^2$$

- $w_0$  WILL BE THE BIAS TERM!

# “NON-LINEAR” DATA?

- FOR EXAMPLE, WHAT IS THE BEST DEGREE 2 POLYNOMIAL?



- HOW CAN WE REUSE THE “LEAST-SQUARES MACHINERY”?

# IDEA: DATA TRANSFORMATION

- WE INCREASED THE FLEXIBILITY OF OUR PREDICTOR BY A FORM OF DATA TRANSFORMATION/AUGMENTATION

- $X'_{n \times (d+1)} = \begin{pmatrix} x_1^1 & \dots & x_d^1 & 1 \\ \vdots & \ddots & \vdots & 1 \\ x_1^n & \dots & x_d^n & 1 \end{pmatrix}$

- CAN WE USE THE SAME IDEA TO MAKE OUR PREDICTOR EVEN MORE FLEXIBLE (NON-LINEAR)?

# EXAMPLE

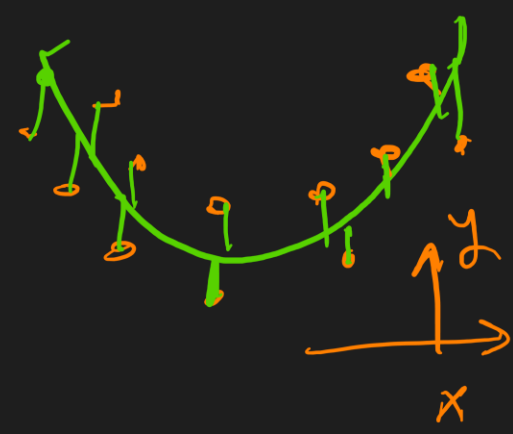
$$g \approx \underline{a}x^2 + \underline{b}x + \underline{c}$$

$$X = \begin{pmatrix} x^1 \\ \vdots \\ x^n \end{pmatrix}$$

$$X' = \begin{pmatrix} x^1 & (x^1)^2 & 1 \\ \vdots & \vdots & \vdots \\ x^n & (x^n)^2 & 1 \end{pmatrix} (X')^T$$

$$W' = \begin{pmatrix} b \\ a \\ c \end{pmatrix}$$

$$(X')^T W' = b \underline{x^1} + a \underline{(x^1)^2} + \underline{c}$$



$$\rightarrow \binom{2+1}{1} = 3 \quad \begin{matrix} d=1 \\ M=2 \end{matrix}$$

is  $\hat{y}$  a linear function  
of  $x^2$  and  $x$

$$\hat{y} = ax^2 + bx + c$$

$\log \hat{y}$   
↘



# LEAST-SQUARES FOR POLYNOMIALS

- IDEA:  $ax^2 + bx + c$  IS STILL LINEAR WITH RESPECT TO THE PARAMETERS! (W.R.T.  $a, b$  AND  $c$ )

- INSTEAD OF  $X_{n \times 1} = \begin{pmatrix} x^1 \\ \dots \\ x^n \end{pmatrix}$  USE  $X'_{n \times 3} = \begin{pmatrix} x^1 & (x^1)^2 & 1 \\ \dots & \dots & \dots \\ x^n & (x^n)^2 & 1 \end{pmatrix}$

- TREAT  $X_{n \times 3}$  AS IF IT WAS YOUR ORIGINAL INPUT DATA
- WE CAN EXTEND THIS TO HIGHER DEGREE POLYNOMIALS SIMILARLY, E.G.,  $ax^3 + bx^2 + cx + d$
- NOTEBOOK EXAMPLE

# MULTIVARIATE POLYNOMIALS

$d = 2$   
↓  
 $d = 6$   
2

- HOW ABOUT WHEN  $x$  IS MULTIVARIATE ITSELF?
  - $w_1x_1 + w_2x_2 + w_3x_1x_2 + w_4(x_1)^2 + w_5(x_2)^2 + w_6$
  - INSTEAD OF  $(x_1, x_2)$  USE  $(x_1 \quad x_2 \quad x_1x_2 \quad (x_1)^2 \quad (x_2)^2 \quad 1)$
- TREAT THE NEW  $X$  AS (A HIGHER-DIMENSIONAL) INPUT

- INPUT DIMENSION:  $d$
- DEGREE OF POLYNOMIAL:  $M$
- NUMBER OF TERMS (MONOMIALS) OF DEGREE AT MOST  $M \approx$

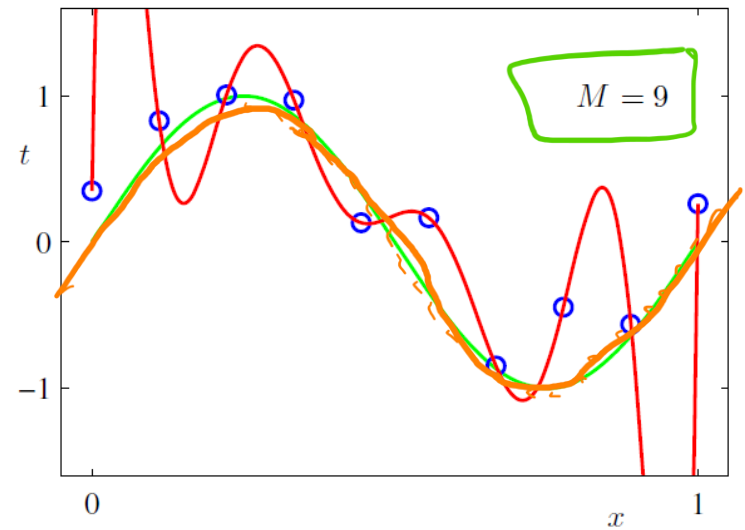
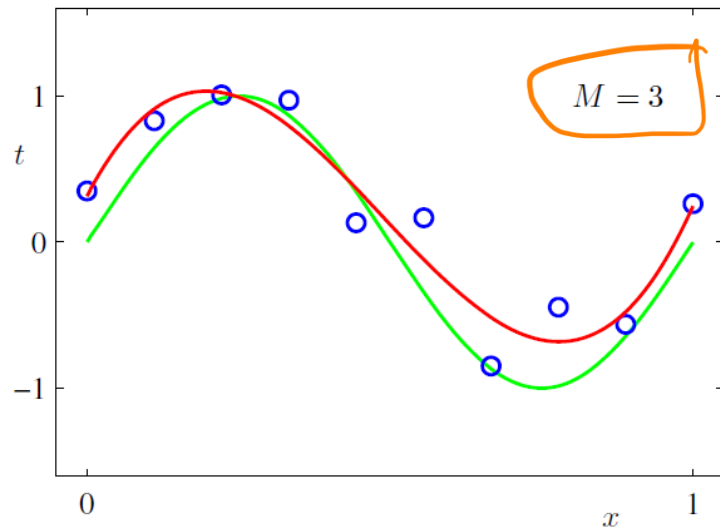
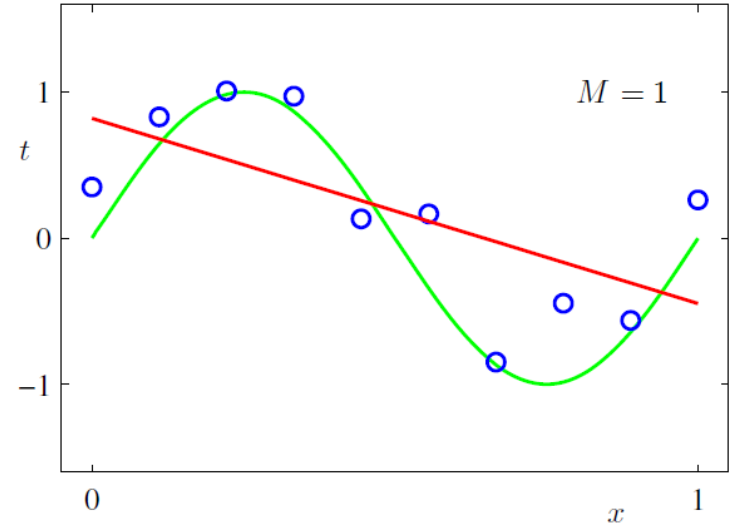
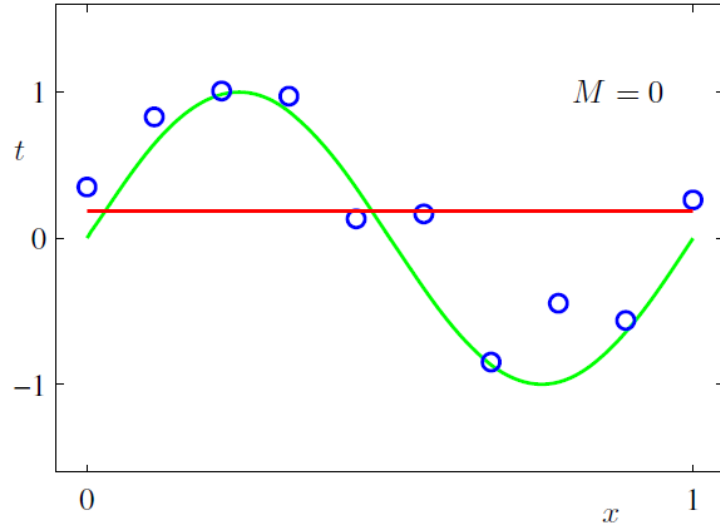
$$\binom{M+d}{d} = \binom{M+d}{M}$$

$$d=2 \rightarrow \binom{4}{2} = 6 \quad \checkmark$$

$$M=2$$


 n balls

# OVERFITTING

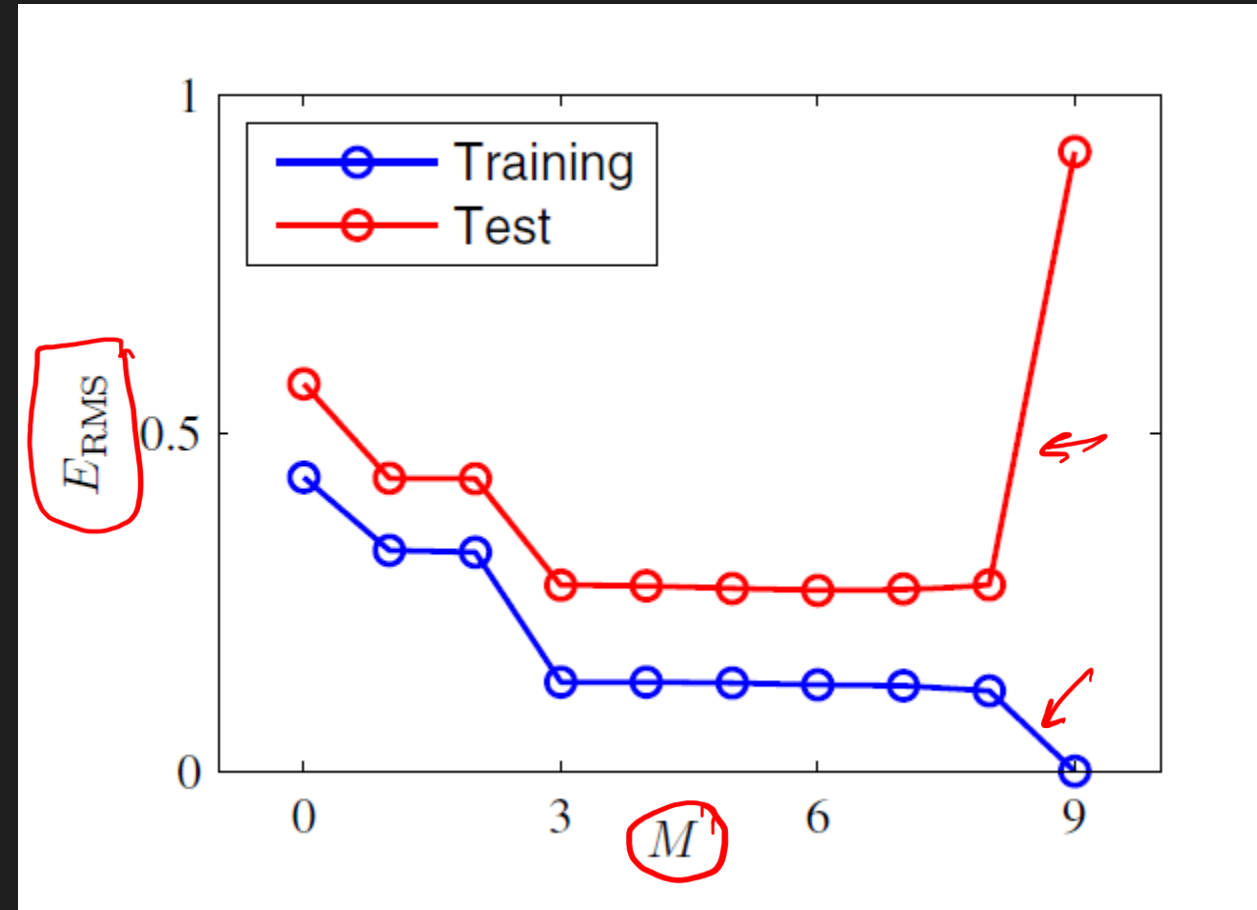


# OVERFITTING

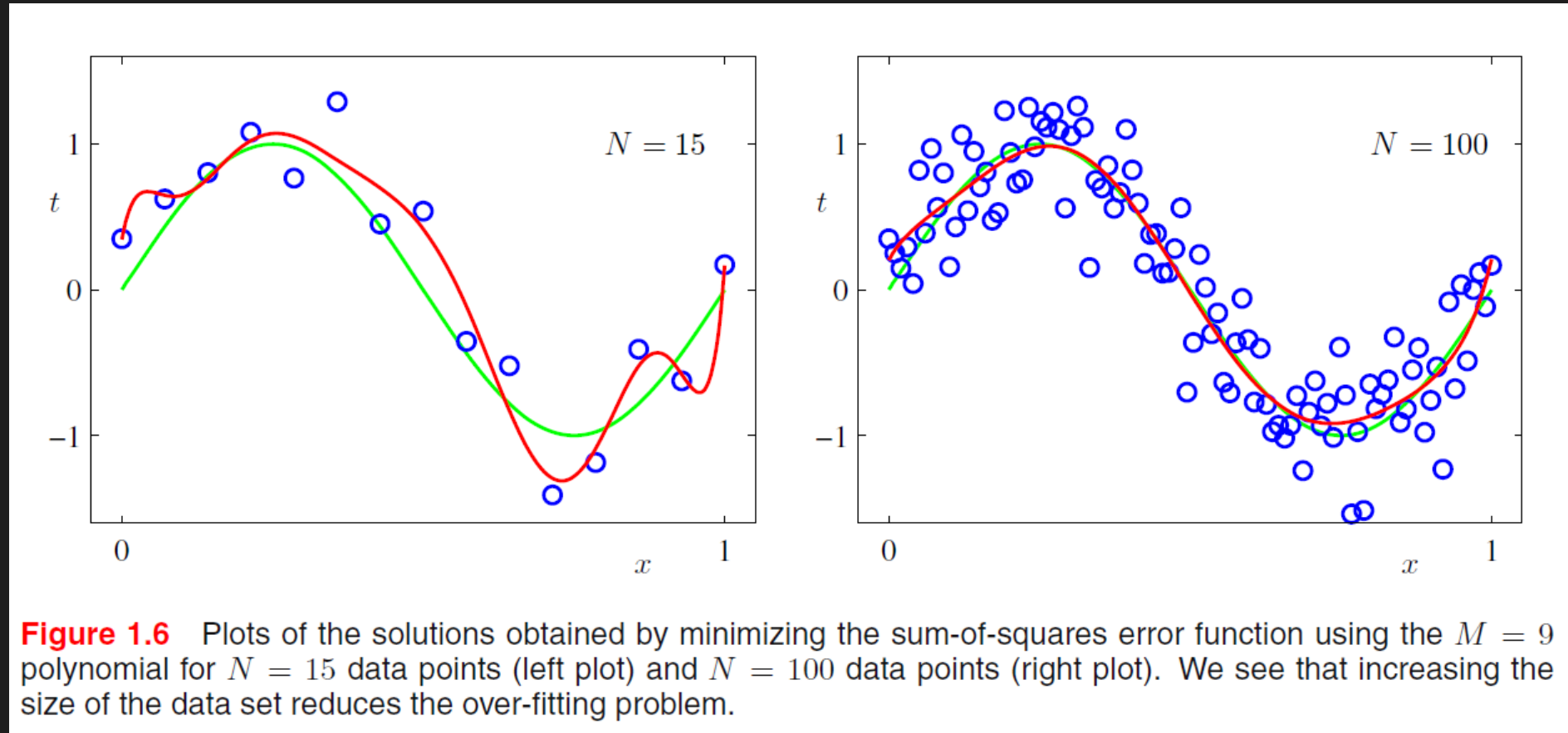
- DIVIDE THE DATA RANDOMLY TO “TRAIN” AND “TEST” SETS

- ROOT-MEAN-SQUARE ERROR FOR EACH SET:

$$\sqrt{\frac{\|\hat{Y} - Y\|_2^2}{n}} = \sqrt{\frac{\sum_{i=1}^n (y_i - \hat{y}_i)^2}{n}}$$



# MORE DATA, LESS OVER-FITTING



# THE TRADE-OFF

- A POWERFUL/FLEXIBLE CURVE-FITTING METHOD
  - ➔ • SMALL TRAINING ERROR
  - ➔ • REQUIRES MORE TRAINING DATA TO GENERALIZE
    - OTHERWISE LARGE TEST ERROR
- A LESS FLEXIBLE CURVE-FITTING METHOD
  - ➔ • LARGER TRAINING ERROR
  - ➔ • REQUIRES LESS TRAINING DATA
  - ➔ • SMALLER DIFFERENCE BETWEEN TRAINING AND TEST ERROR
- THE SO-CALLED “BIAS-VARIANCE” TRADE-OFF

# THE CASE OF MULTIVARIATE POLYNOMIALS

- ASSUME  $M \gg d$
- NUMBER OF TERMS (MONOMIALS):  $\approx \left(\frac{M}{d}\right)^d$
- #TRAINING SAMPLES  $\approx$  #PARAMETERS  $\approx \left(\frac{M}{d}\right)^d$ 
  - #TRAINING SAMPLES SHOULD INCREASE EXPONENTIALLY WITH  $d$
  - SUSCEPTIBLE TO OVER-FITTING...
  - AN EXAMPLE OF CURSE OF DIMENSIONALITY!
- WE CAN SAY SAMPLE COMPLEXITY OF LEARNING MULTIVARIATE POLYNOMIALS IS EXPONENTIAL IN  $d$ 
  - ORTHOGONAL TO COMPUTATIONAL COMPLEXITY

$$d \leq 1$$

$$\binom{m+d}{m}$$

$$\binom{m+d}{1}$$

$$\approx m+d$$



