# INTRODUCTION TO MACHINE LEARNING COMPSCI 4ML3

## Lecture 4

### Hassan Ashtiani

# MATRIX FORM OLS

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

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

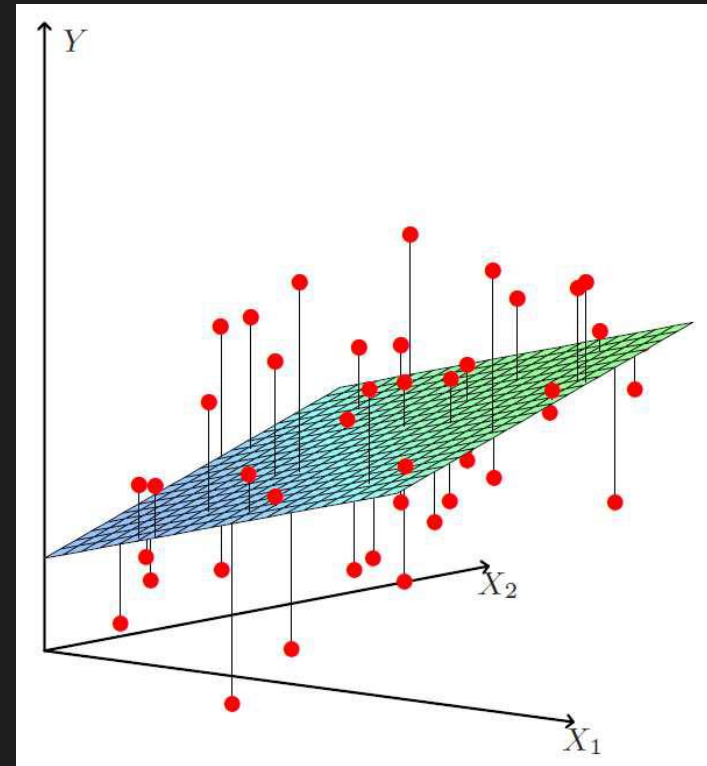$$\underset{W \in \mathbb{R}^{d \times 1}}{min} \|XW - Y\|_2^2$$

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

# BIAS/INTERCEPT TERM

*WE ARE MISSING THE BIAS TERM $(w_0)$*

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

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

# BIAS/INTERCEPT TERM

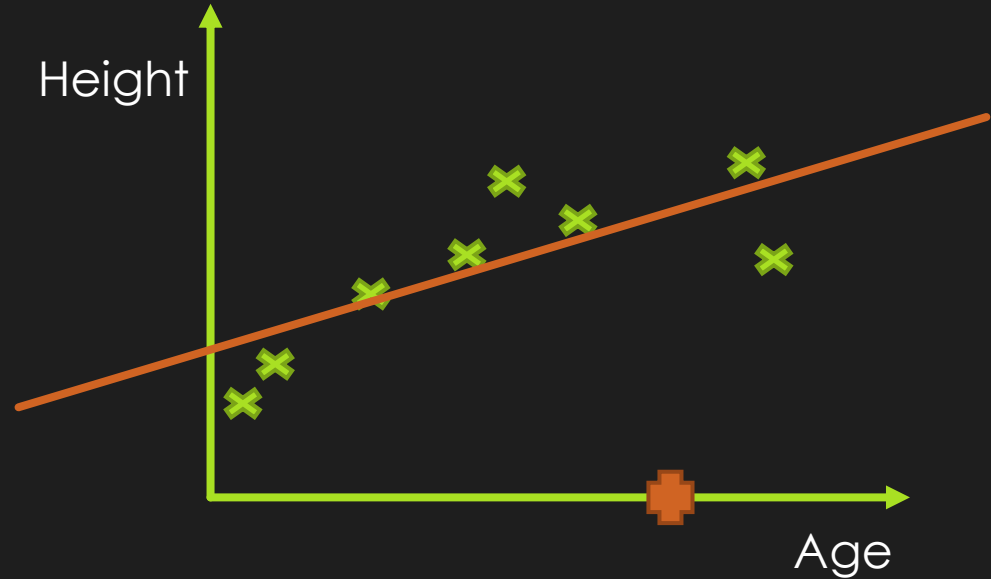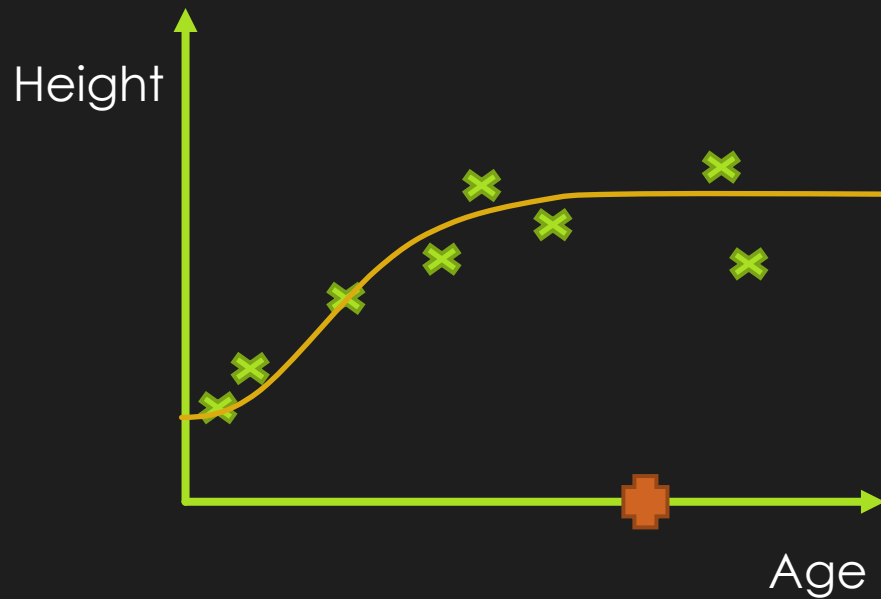- *Add a new auxiliary dimension to the data*

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

- *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

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

- Can we use the same idea to make our predictor even more flexible (non-linear)?

# EXAMPLE

# 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 $X_{n \times 1} = \begin{pmatrix} x^1 \\ ... \\ x^n \end{pmatrix}$ USE $X'_{n \times 3} = \begin{pmatrix} x^1 & (x^1)^2 & 1 \\ ... & ... & ... \\ x^n & (x^n)^2 & 1 \end{pmatrix}$
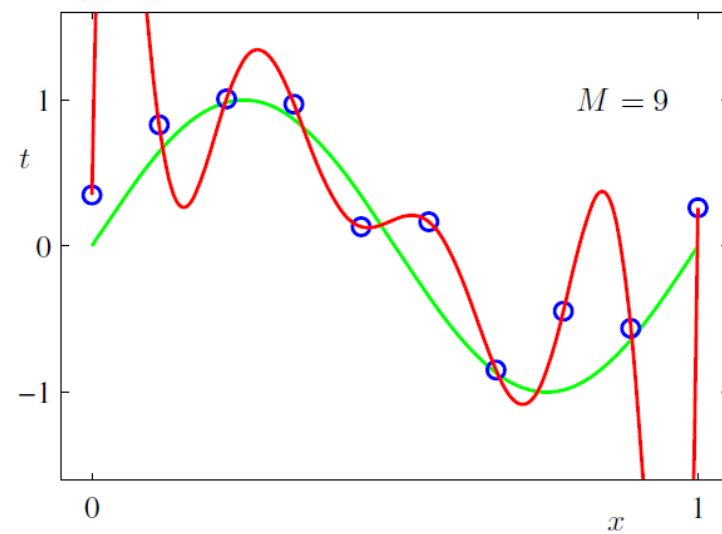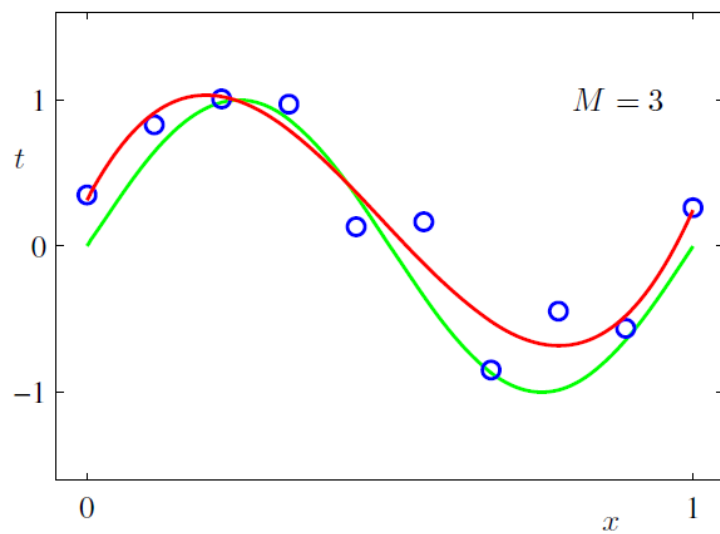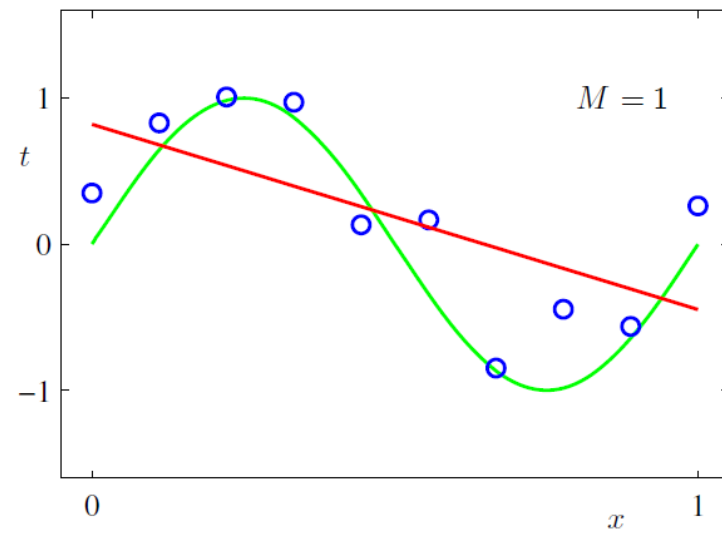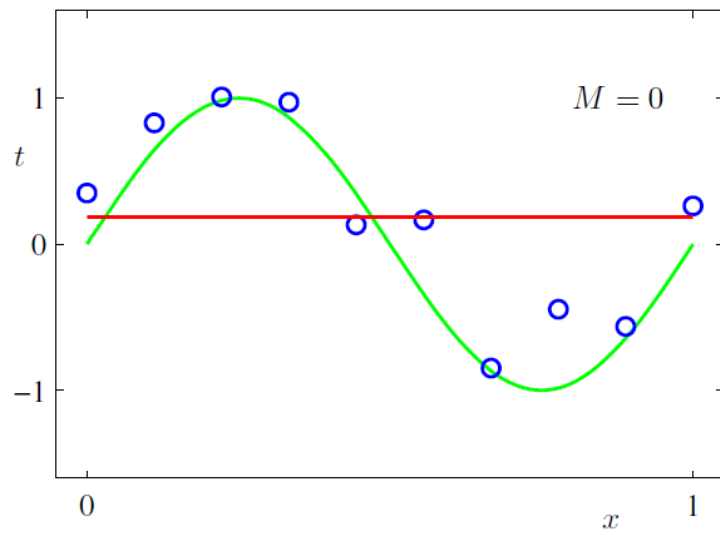
- TREAT $X_{n \times 3}$ AS IF IT WAS YOUR ORIGINAL INPUT DATA

- WE CAN EXTEND THIS TO HIGHER DEGREE POLYNOMIALS SIMILARLY, E.G., $ax^3 + bx^2 + cx + d$

- NOTEBOOK EXAMPLE

# MULTIVARIATE POLYNOMIALS

- How about when $x$ is multivariate itself?

  - $w_1 x_1 + w_2 x_2 + w_3 x_1 x_2 + w_4 (x_1)^2 + w_5 (x_2)^2 + w_6$

  - Instead of $(x_1, x_2)$ use $(x_1 \quad x_2 \quad x_1 x_2 \quad (x_1)^2 \quad (x_2)^2 \quad 1)$

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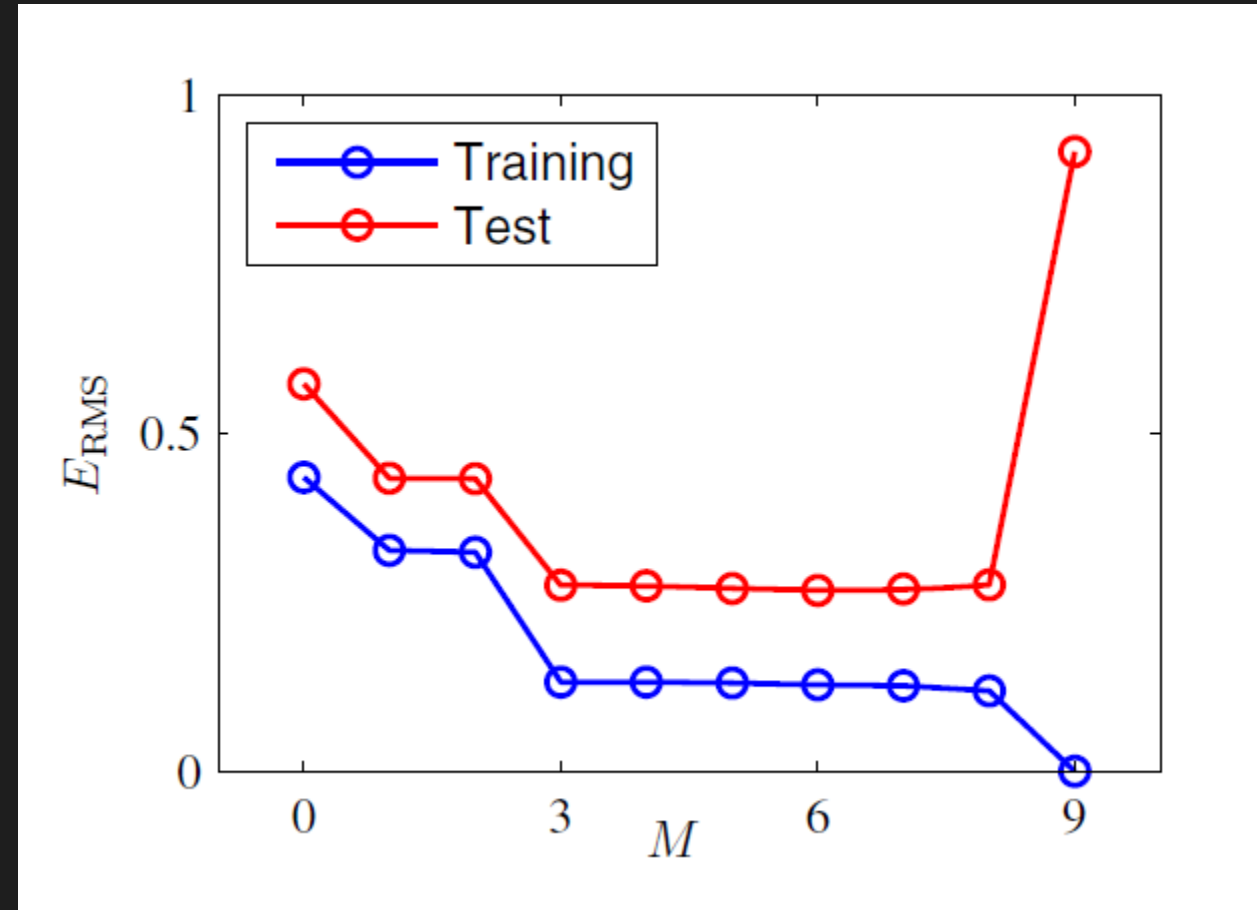$$\binom{M + d}{d} = \binom{M + d}{M}$$

# OVERFITTING

# OVERFITTING

- Divide the data randomly to "train" and "test" sets

- Root-mean-square error for each set:

- $$\sqrt{\frac{\|\hat{Y}-Y\|_2^2}{n}} = \sqrt{\frac{\sum_{i=1}^{n}(y_i-\hat{y_i})^2}{n}}$$
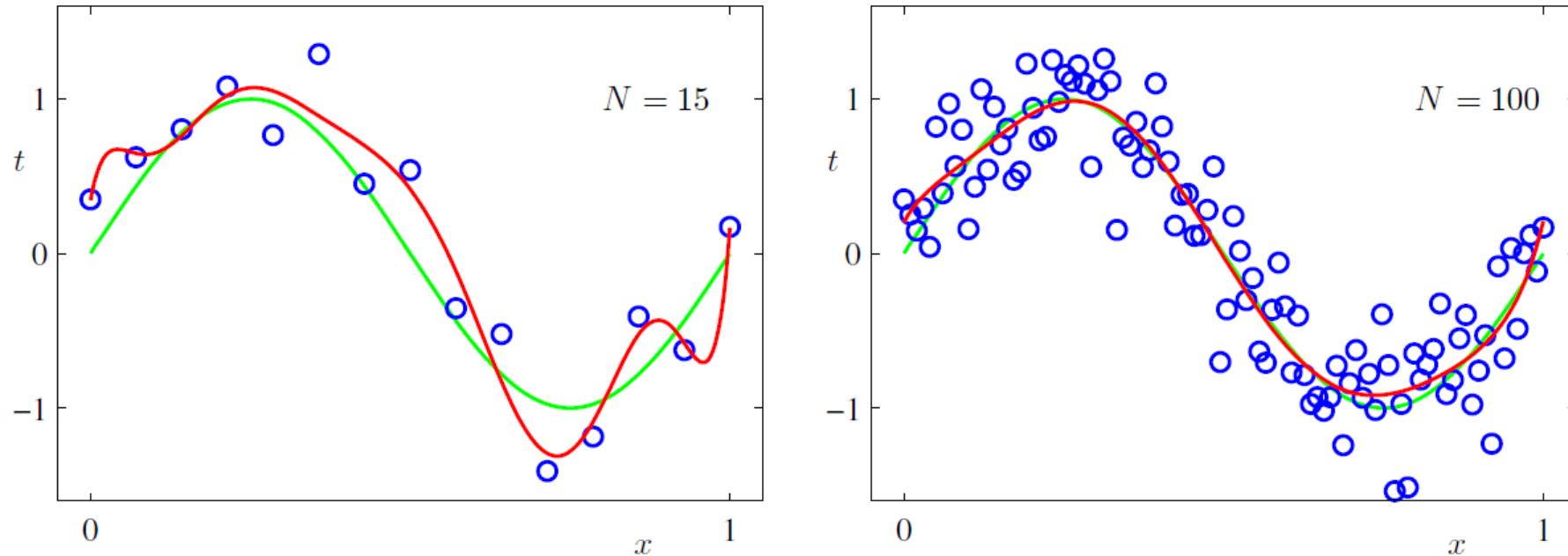
# MORE DATA, LESS OVER-FITTING



**Figure 1.6** Plots of the solutions obtained by minimizing the sum-of-squares error function using the $M = 9$ polynomial for $N = 15$ data points (left plot) and $N = 100$ data points (right plot). We see that increasing the size of the data set reduces the over-fitting problem.

# 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): $\approx (\frac{M}{d})^d$

- #training samples $\approx$ #parameters $\approx (\frac{M}{d})^d$

  - #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