



# Applied Data Science fall 2017

## Session 4: Multi-variate linear regression

***Instructor: Prof. Stanislav Sobolevsky***  
***Course Assistants: Tushar Ahuja, Maxim Temnogorod***

# Bi-variate to multi-variate

$$y \sim x$$

$$y = w_1 x + w_0 + \varepsilon$$

$$x = (x_1, x_2, \dots, x_n)$$

$$y = \sum_{j=1}^n w_j x_j + w_0 + \varepsilon$$

**Intercept:**  $x_{n+1} = 1, w_{n+1} = w_0$

$$y = \sum_{j=1}^{n+1} w_j x_j + \varepsilon$$

$$y = w^T x + \varepsilon$$

# Matrix form

$$y = x_1 + 2x_2 + 3x_3$$

$$x_3 = 1$$

$$w_1 = 1, \quad w_2 = 2, \quad w_3 = 3$$

$$y = w_1x_1 + w_2x_2 + w_3x_3 = (w_1 \ w_2 \ w_3) \begin{pmatrix} x_1 \\ x_2 \\ x_3 \end{pmatrix} = w^T x$$

# Matrix form

$$y = x_1 + 2x_2 + 3$$

$$x_1 = 1, x_2 = 1, y = 6.2$$

$$x_1 = 2, x_2 = 0, y = 4.9$$

$$x_1 = 3, x_2 = -1, y = 4.1$$

$$x_1 = 4, x_2 = -2, y = 2.9$$

$$x_3 = 1$$

$$w = \begin{pmatrix} 1 \\ 2 \\ 3 \end{pmatrix}$$

$$Y = \begin{pmatrix} 6.2 \\ 4.9 \\ 4.1 \\ 2.9 \end{pmatrix} = Xw + E$$

$$X = \begin{pmatrix} 1 & 1 & 1 \\ 2 & 0 & 1 \\ 3 & -1 & 1 \\ 4 & -2 & 1 \end{pmatrix}$$

$$E = Y - Xw = \begin{pmatrix} 0.2 \\ -0.1 \\ 0.1 \\ -0.1 \end{pmatrix}$$

# Least-square estimate

$$y = w^T x + \varepsilon$$

$$X = \{(x_j^i), j = 1..n, i = 1..N\}, Y = \{(y^i), i = 1..N\}$$

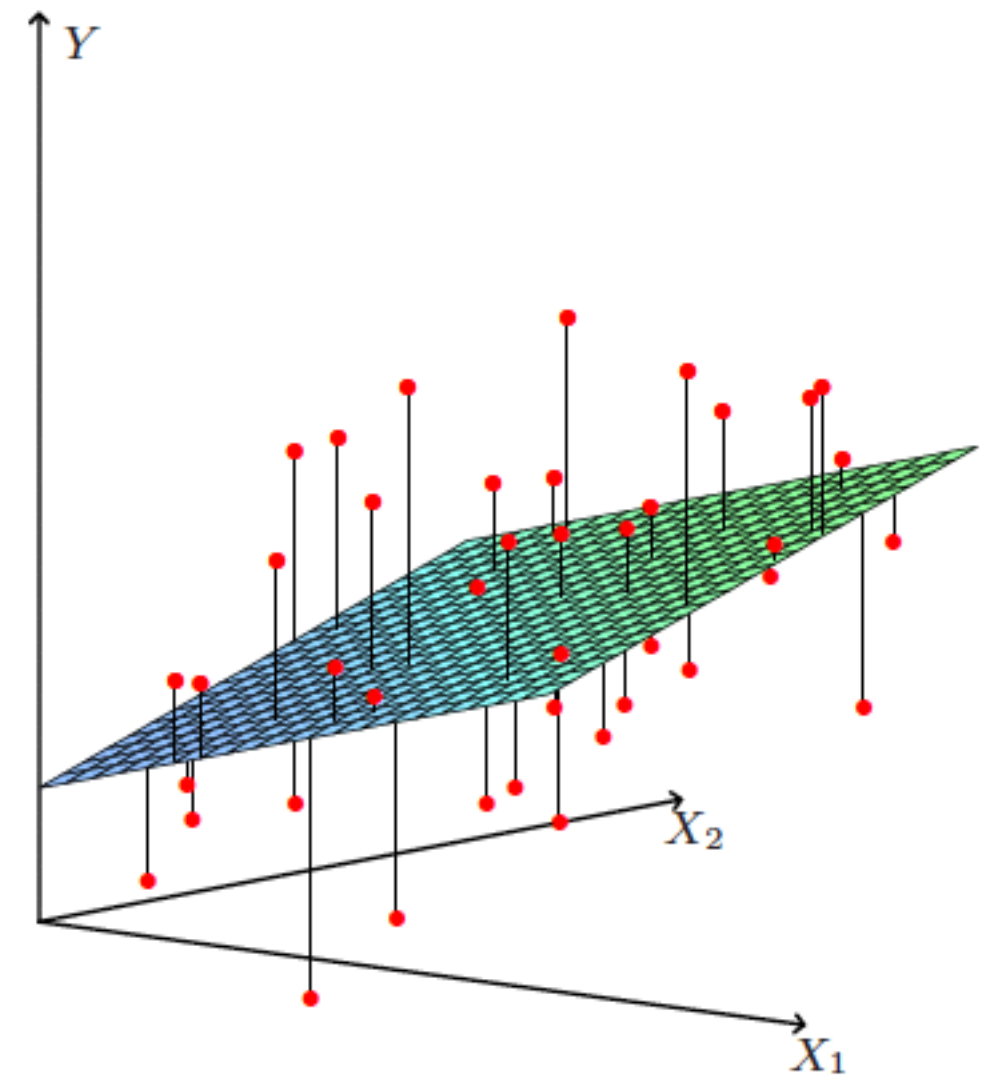
$$E = (\varepsilon_i, i = 1..N)^T = Y - Xw$$

$$RSS(w) = \sum_i \varepsilon_i^2 = \sum_i (y^i - w^T x^i)^2$$

$$RSS(w) = E^T E$$

$$RSS(w) = (Y - Xw)^T (Y - Xw)$$

$$\hat{w} = \operatorname{argmin}_w RSS(w)$$



Hastie, *et al.*, The Elements of Statistical Learning, Data Mining, Inference and Prediction, 2<sup>nd</sup> Edition, Springer. (Free: [http://web.stanford.edu/~hastie/local.ftp/Springer/OLD/ESLII\\_print4.pdf](http://web.stanford.edu/~hastie/local.ftp/Springer/OLD/ESLII_print4.pdf))  
These materials are included under the fair use exemption and are restricted from further use.

# Least-square estimate

$$RSS(w) = (Y - Xw)^T(Y - Xw) \quad \hat{w} = \operatorname{argmin}_w RSS(w)$$

$$0 = \frac{\partial RSS(\hat{w})}{dw} = -2X^T(Y - X\hat{w})$$

$$X^T Y = (X^T X)\hat{w}$$

$$\hat{w} = (X^T X)^{-1} X^T Y$$

# Least-squares: geometric sense

$$y = \sum_j w_j x_j$$

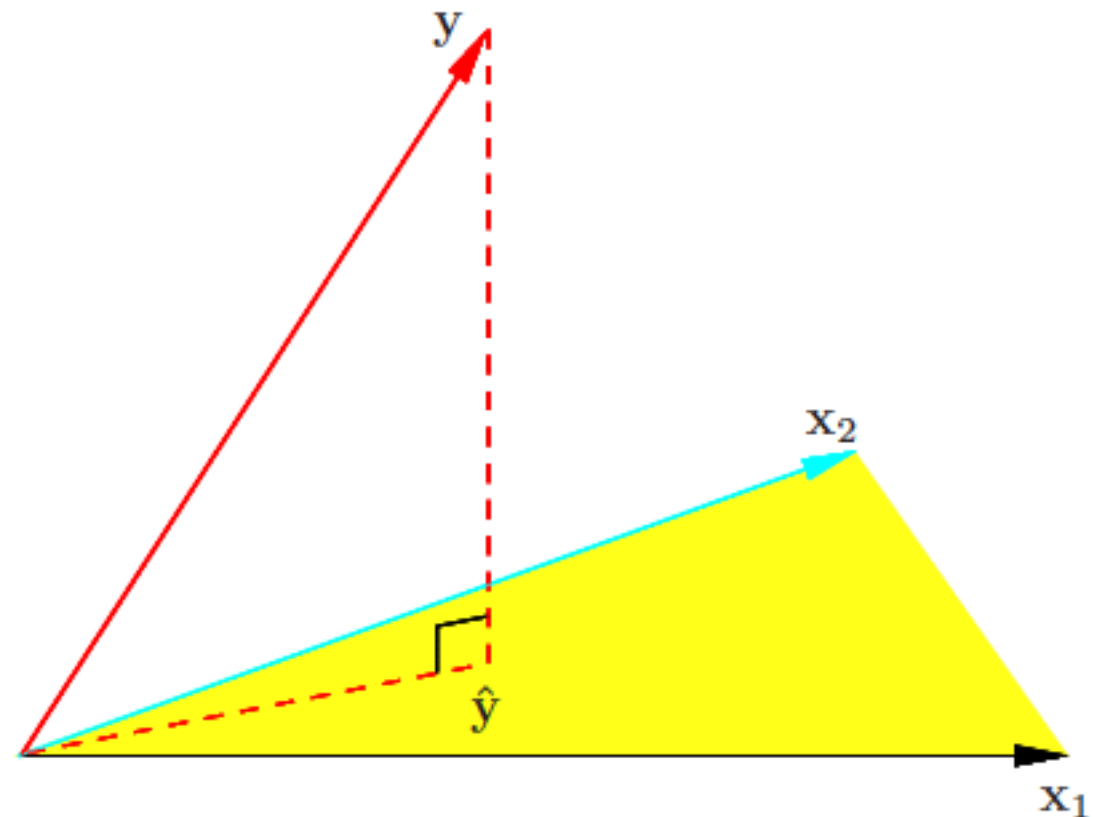
$$Y \sim \sum_j w_j X_j$$

$$\hat{Y} = X\hat{w} = X(X^T X)^{-1} X^T Y$$

$$H = X(X^T X)^{-1} X^T$$

$$\hat{Y} = HY$$

$$\hat{w} = (X^T X)^{-1} X^T Y$$



Hastie, *et al.*, The Elements of Statistical Learning, Data Mining, Inference and Prediction, 2<sup>nd</sup> Edition, Springer. (Free: [http://web.stanford.edu/~hastie/local.ftp/Springer/OLD/ESLII\\_print4.pdf](http://web.stanford.edu/~hastie/local.ftp/Springer/OLD/ESLII_print4.pdf))  
These materials are included under the fair use exemption and are restricted from further use.

# Linear Model - R-squared

$$R^2 = 1 - \frac{RSS}{\sum_i (y_i - \bar{y})^2} = \frac{\sum_i (\hat{y}_i - \bar{y})^2}{\sum_i (y_i - \bar{y})^2},$$





# Linear Model - Orthogonal regressors

If you know  $y \sim w_j x_j$

do you know  $y \sim x = (x_1, x_2, \dots, x_n)$  ?

$$Y = \begin{pmatrix} 5 \\ 4 \\ 3 \end{pmatrix} \quad X = \begin{pmatrix} 2 & 0 & 1 \\ 1 & 1 & -2 \\ 0 & 2 & 0 \end{pmatrix}$$

$$y \sim 2.8x_1$$

$$y \sim 2x_2$$

$$y \sim -\frac{3}{5}x_3$$

$$y \sim 2.8x_1 - 0.6x_3$$

$$y \sim 2.5x_1 + 1.5x_2$$

# Linear Model - Orthogonal regressors

If you know  $y \sim w_j x_j$

you do know  $y \sim x = (x_1, x_2, \dots, x_n)$  if

$$\langle X_j, X_k \rangle = X_j^T X_k = \sum_i x_j^i x_k^i = 0.$$

# Linear Model - Feature scaling

$$x_j^* = \frac{x_j - \bar{x}_j}{\sigma_j}$$

$$\bar{x}_j = E[X_j] \quad \sigma_j = std[X_j]$$

# Linear Model - Orthogonal regressors

$$\langle X_j, X_k \rangle = X_j^T X_k = \sum_i x_j^i x_k^i = 0$$

$$E[X_j] = 0$$

$$0 = \text{corr}[X_j, X_k] = \frac{\text{Cov}[X_j, X_k]}{\text{std}[X_j] \text{std}[X_k]} = \frac{\frac{\langle X_j, X_k \rangle}{N} - E[X_j]E[X_k]}{\text{std}[X_j] \text{std}[X_k]} = \frac{\langle X_j, X_k \rangle}{\text{std}[X_j] \text{std}[X_k]}$$

$$y \sim x = (x_1, x_2, \dots, x_n)$$

$$y \sim \text{corr}[X_j, Y] x_j$$

$$y \sim \sum_j \text{corr}[X_j, Y] x_j$$

# Multicollinearity

$$\hat{w} = (X^T X)^{-1} X^T Y$$

$$x_1 = \sum_{j=2}^n k_j x_j$$

$$\det(X^T X) \sim 0$$

# Multicollinearity example

$$y = 2x_1$$

$$x_2 = x_1$$

$$y = 2x_2$$

$$y = x_1 + x_2$$

$$y = 0.5x_1 + 1.5x_2$$

$$y = -1000x_1 + 1001x_2$$

$$y = kx_1 + (2 - k)x_2$$



# Overfitting

Area	Month	Price
1010	Jan	100000
2000	Feb	200000
2990	March	300000

$$\text{Price} = 100.1418 * \text{Area}$$

Area	Month	Price
1500	Dec	150000

$$\text{Price} = 100.1418 * 1500 = 150212.7$$



# Overfitting

Area	Month	Price
1010	1	100000
2000	2	200000
2990	3	300000

$$\text{Price} = 100.000 * \text{Month}$$

Area	Month	Price
1500	12	150000

$$\text{Price} = 12 * 100.000 = 1.200.000?$$





# Overfitting

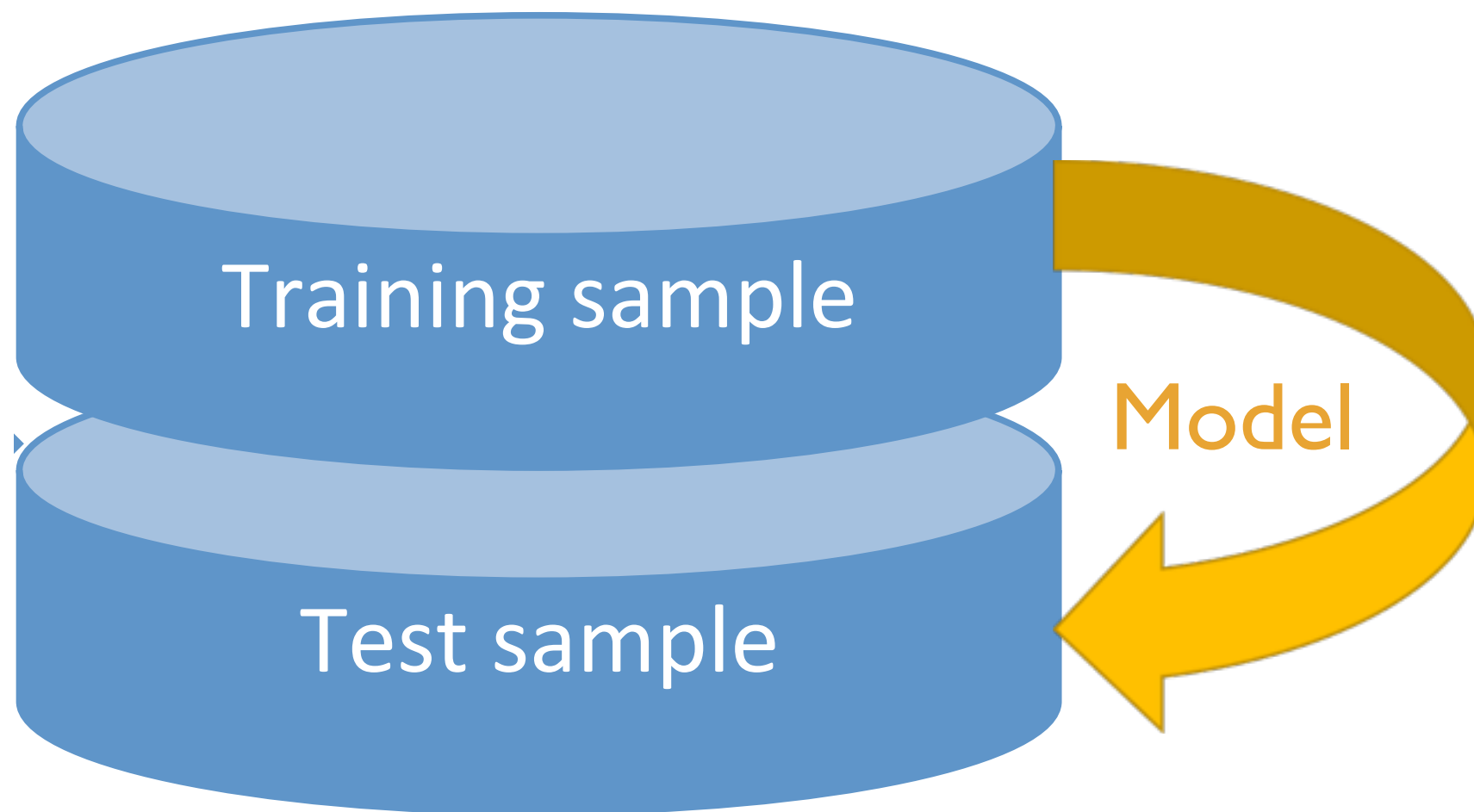
Area	Month	Price
1010	1	90000
2000	2	200000
2990	3	310000

$$\text{Price} = -1000 * \text{Area} + 1.100.000 * \text{Month}$$

Area	Month	Price
1500	12	150000

$$\text{Price} = -1000 * 1500 + 1.1\text{M} * 12 = 11.7\text{M?}$$

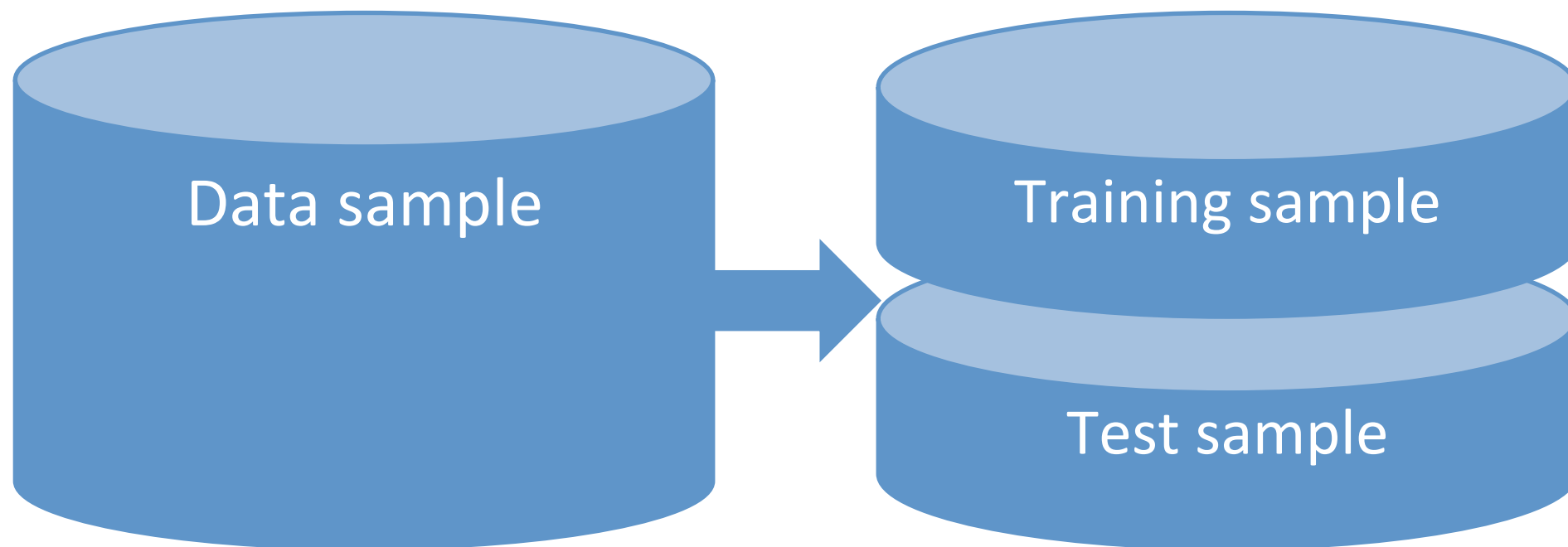
# Validating regression



# Fighting overfitting

No future data?

Make it!



# Adding more features is not always better

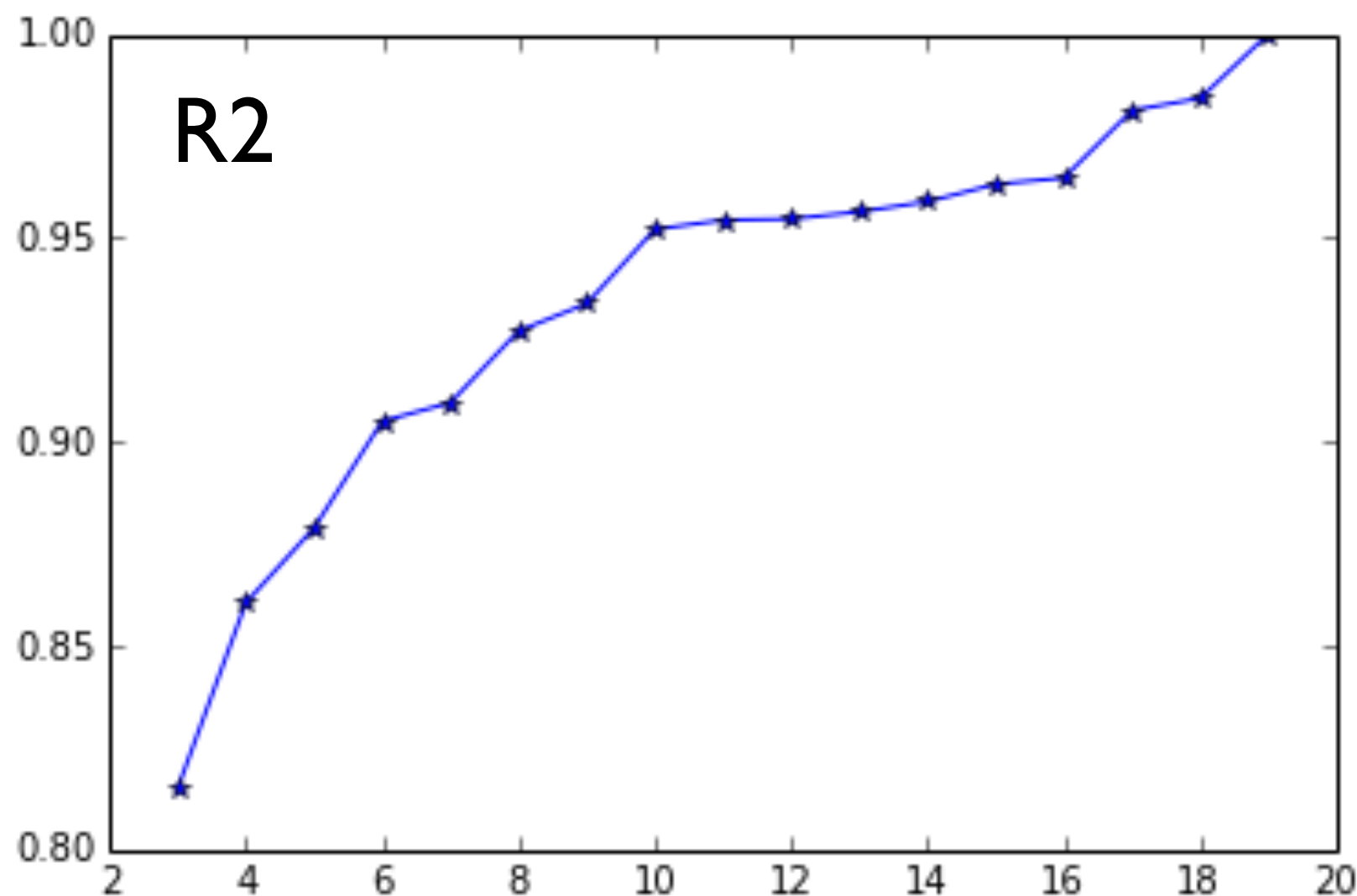
$$Y \sim X_1 + X_2 + X_3 + \mathcal{N}(0, 1.1)$$

$$Y \sim X_1 + X_2 + X_3$$

$$X_j \sim \mathcal{N}(0, 1), j = 4, 5, 6, \dots?$$

# Adding more features is not always better

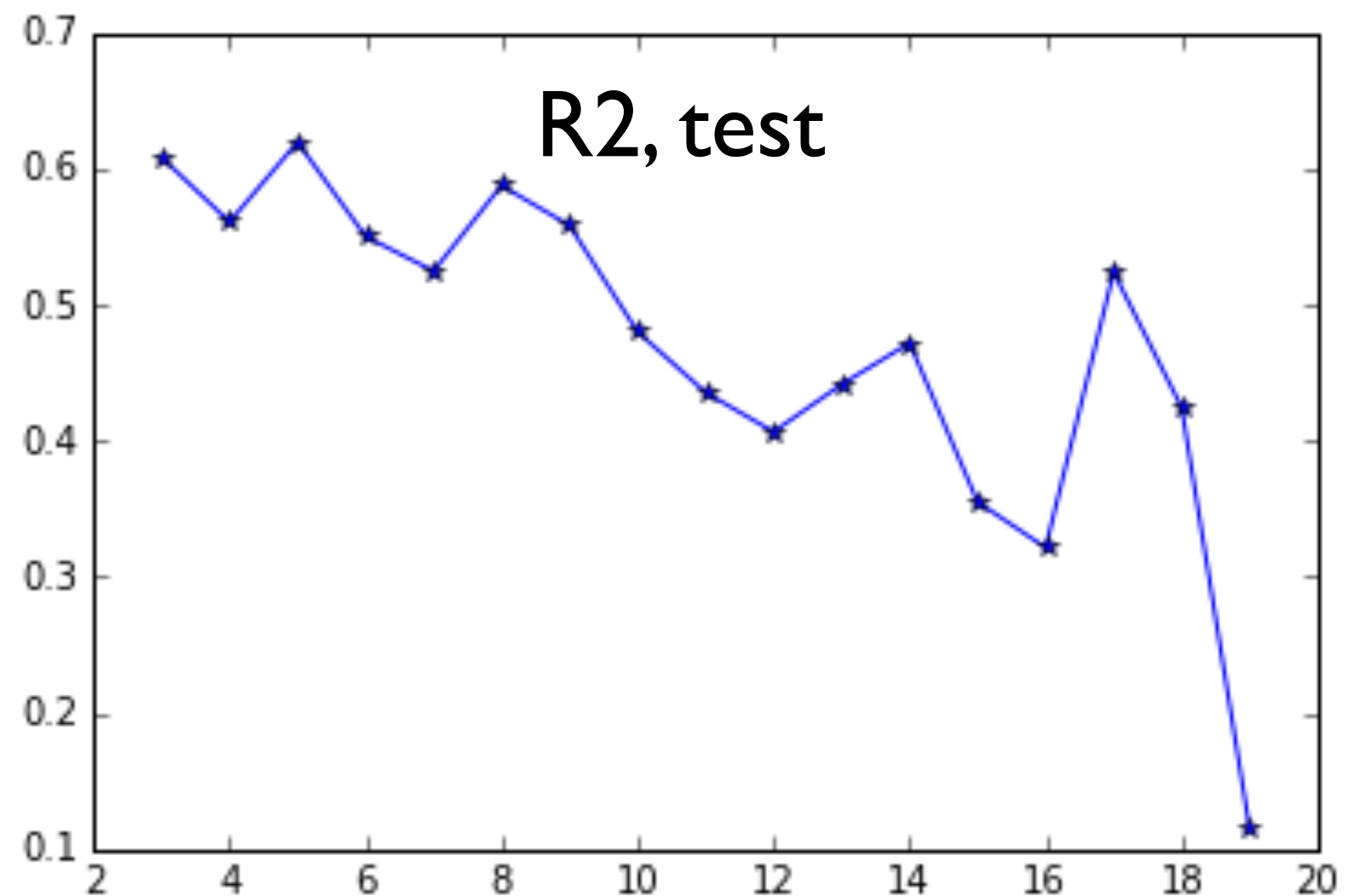
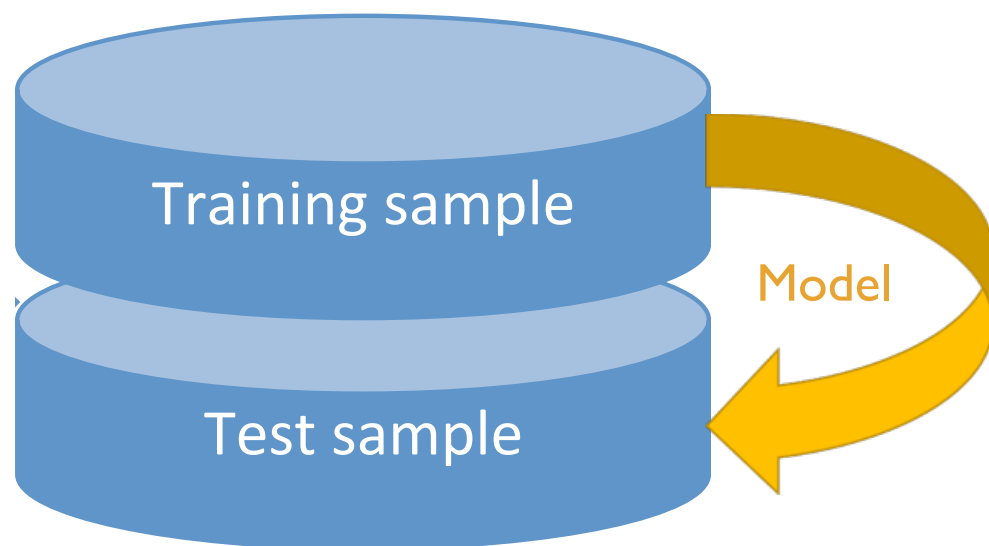
$$Y \sim \sum_{j=1}^k X_j, k = 3, 4, \dots, 20 \quad X_j \sim \mathcal{N}(0, 1), j = 4, 5, 6, \dots?$$



# Adding more features is not always better

$$Y \sim \sum_{j=1}^k X_j, k = 3, 4, \dots, 20$$

$$X_j \sim \mathcal{N}(0, 1), j = 4, 5, 6, \dots?$$



# Polynomial Models

$$y = w_2x^2 + w_1x + w_0$$

$$x_1 = x^2, x_2 = x, x_3 = 1$$

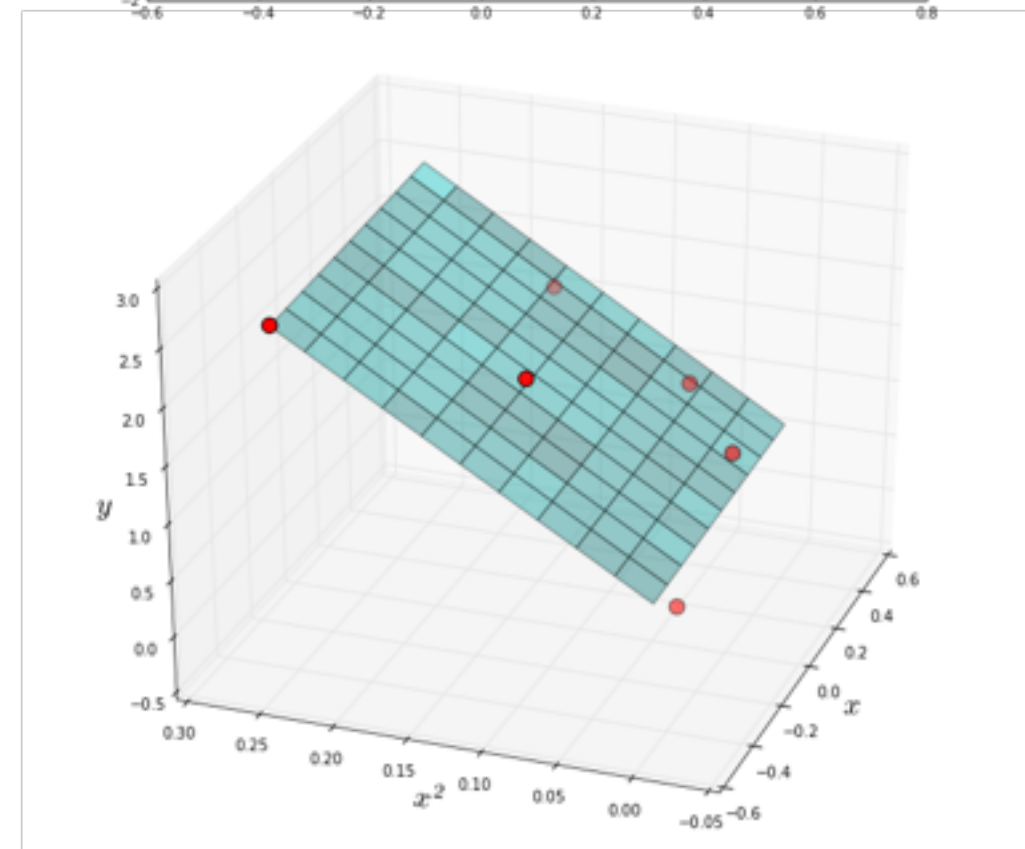
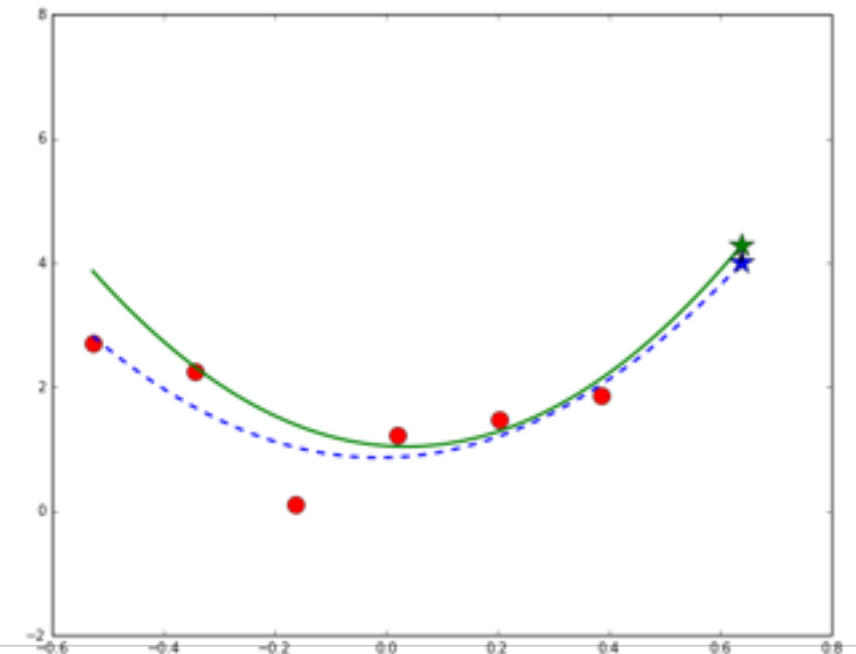
$$y \sim x_1, x_2, x_3$$

$$y = w_mx^m + w_{m-1}x^{m-1} + \dots + w_1x + w_0$$

$$y \sim x^m, x^{m-1}, \dots, x, 1$$

$$y \sim w_{2,0}x_1^2 + w_{1,1}x_1x_2 + w_{2,0}x_2^2 + w_{1,0}x_1 + w_{0,1}x_2 + w_{0,0}$$

$$y \sim 1, x_1, x_2, x_1^2, x_2^2$$





# Polynomial Models - overfitting example

