# DATA SCIENCE & MACHINE LEARNING COURSE

https://www.facebook.com/diceanalytics/
https://pk.linkedin.com/company/diceanalytics
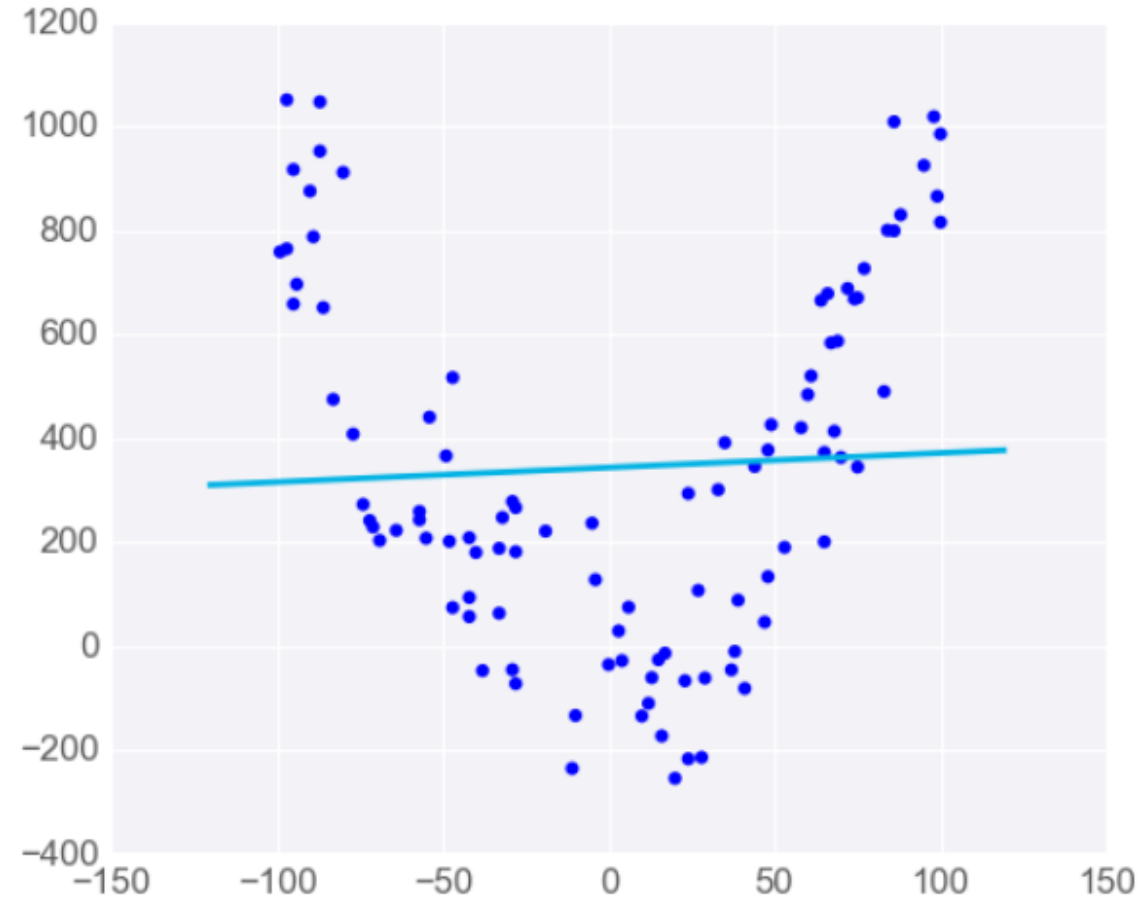
# Regression

# Polynomial Regression

Linear Regression Works Best When the Data is Linear

# Polynomial Regression



$$\hat{y} = w_1 x^3 + w_2 x^2 + w_3 x + w_4$$

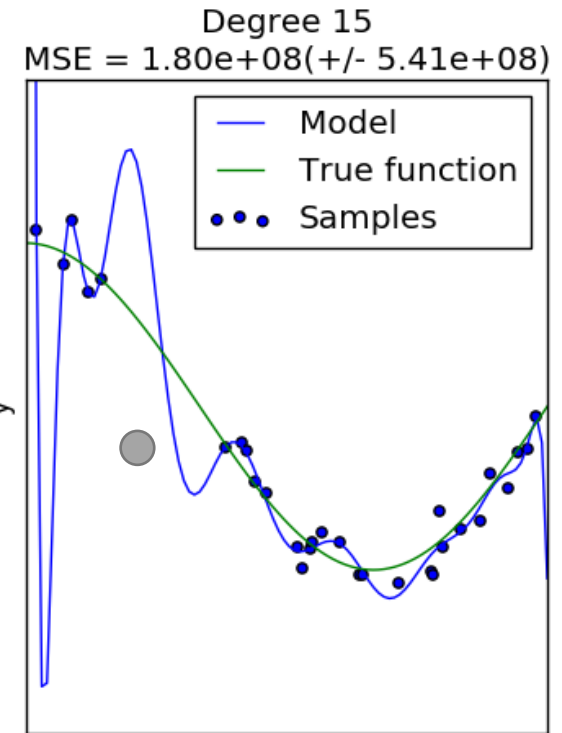# Model Selection



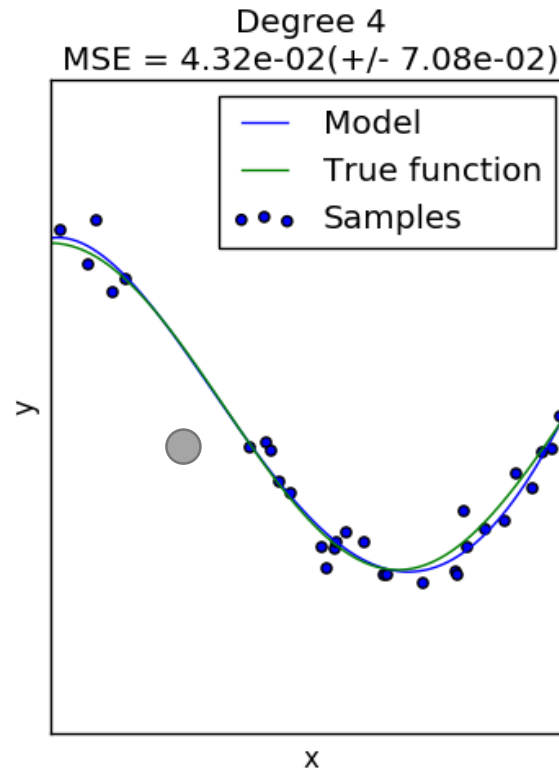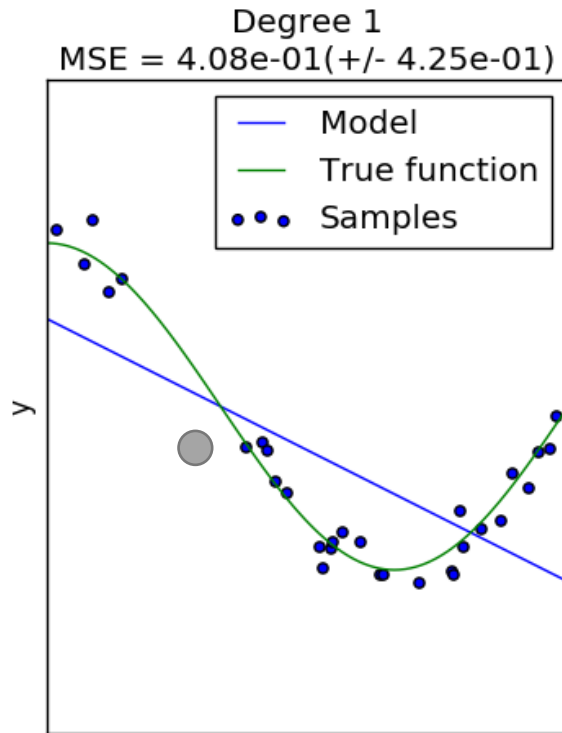Simple Problem                    Complex Solution

# Model Selection



Complex Problem                    Simple Solution

# Under-fitting & Over-fitting

Degree 1
MSE = 4.08e-01(+/- 4.25e-01)

Degree 4
MSE = 4.32e-02(+/- 7.08e-02)

Degree 15
MSE = 1.80e+08(+/- 5.41e+08)
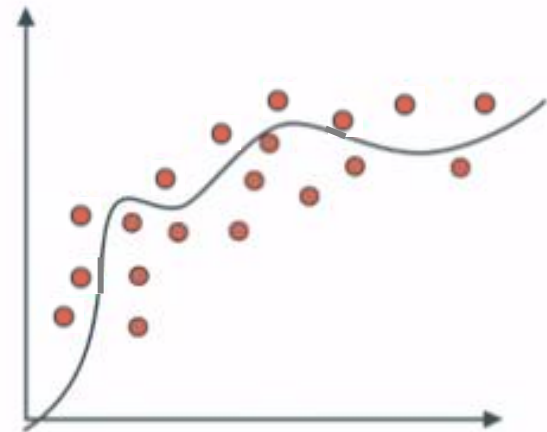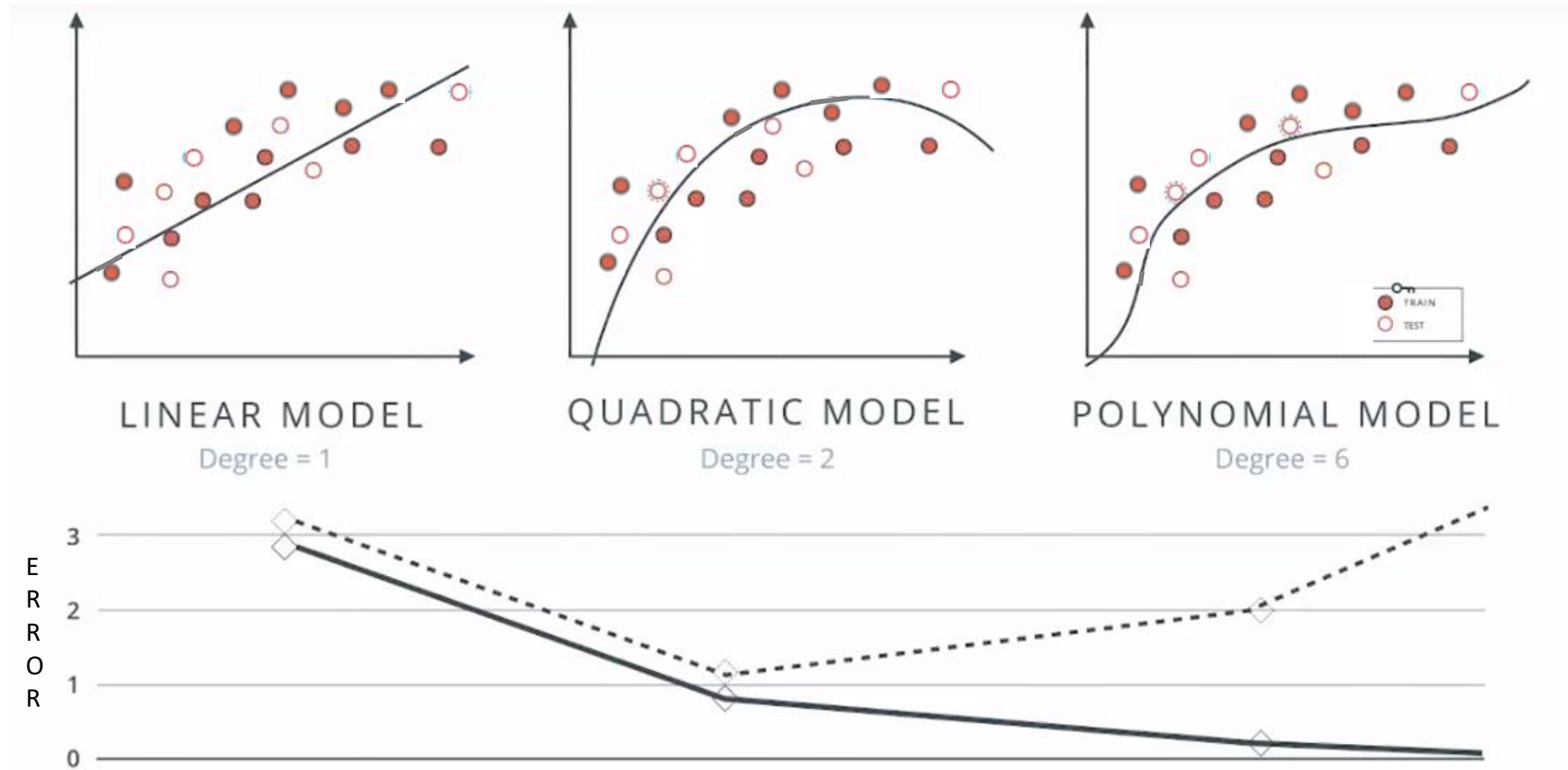
# Model Complexity



HIGH BIAS
Degree = 1

POLYNOMIAL
Degree = 2

HIGH VARIANCE
Degree = 6

# Model Complexity



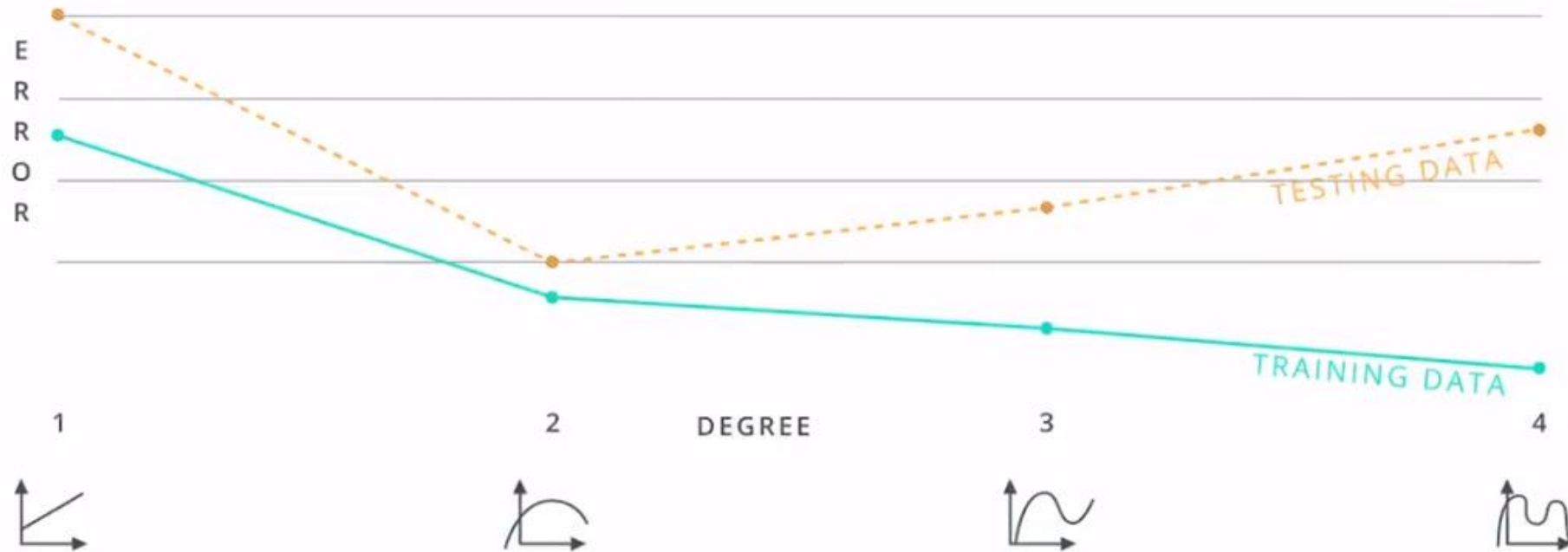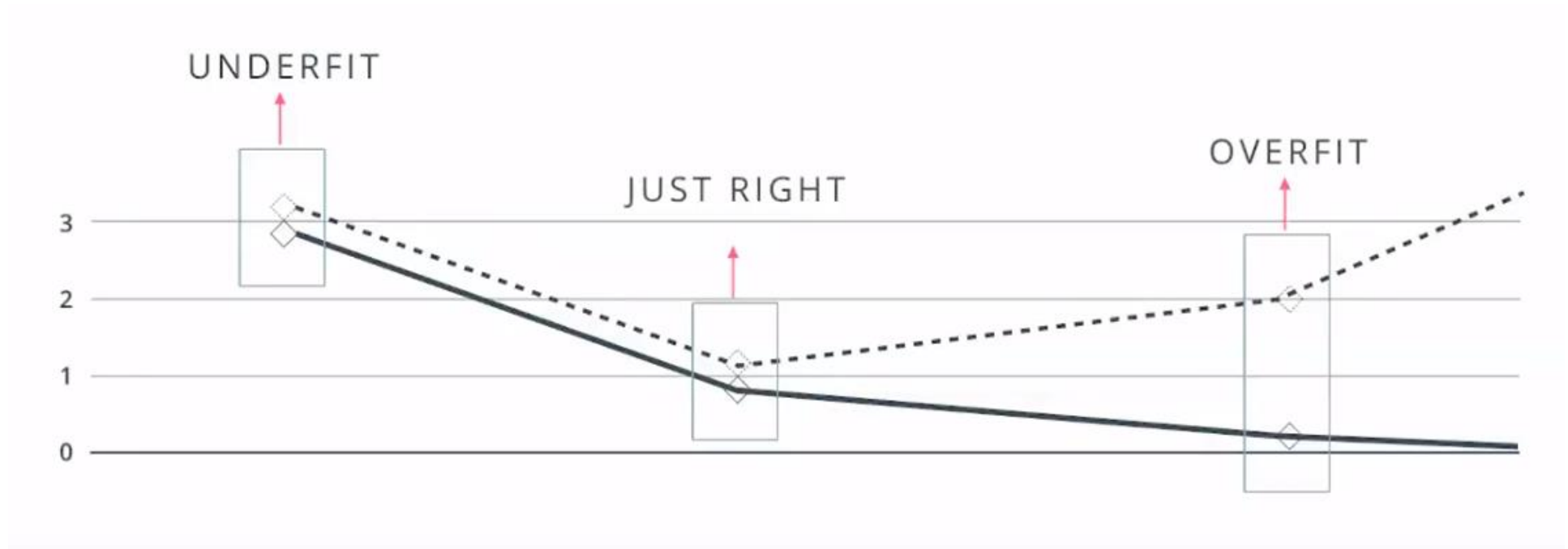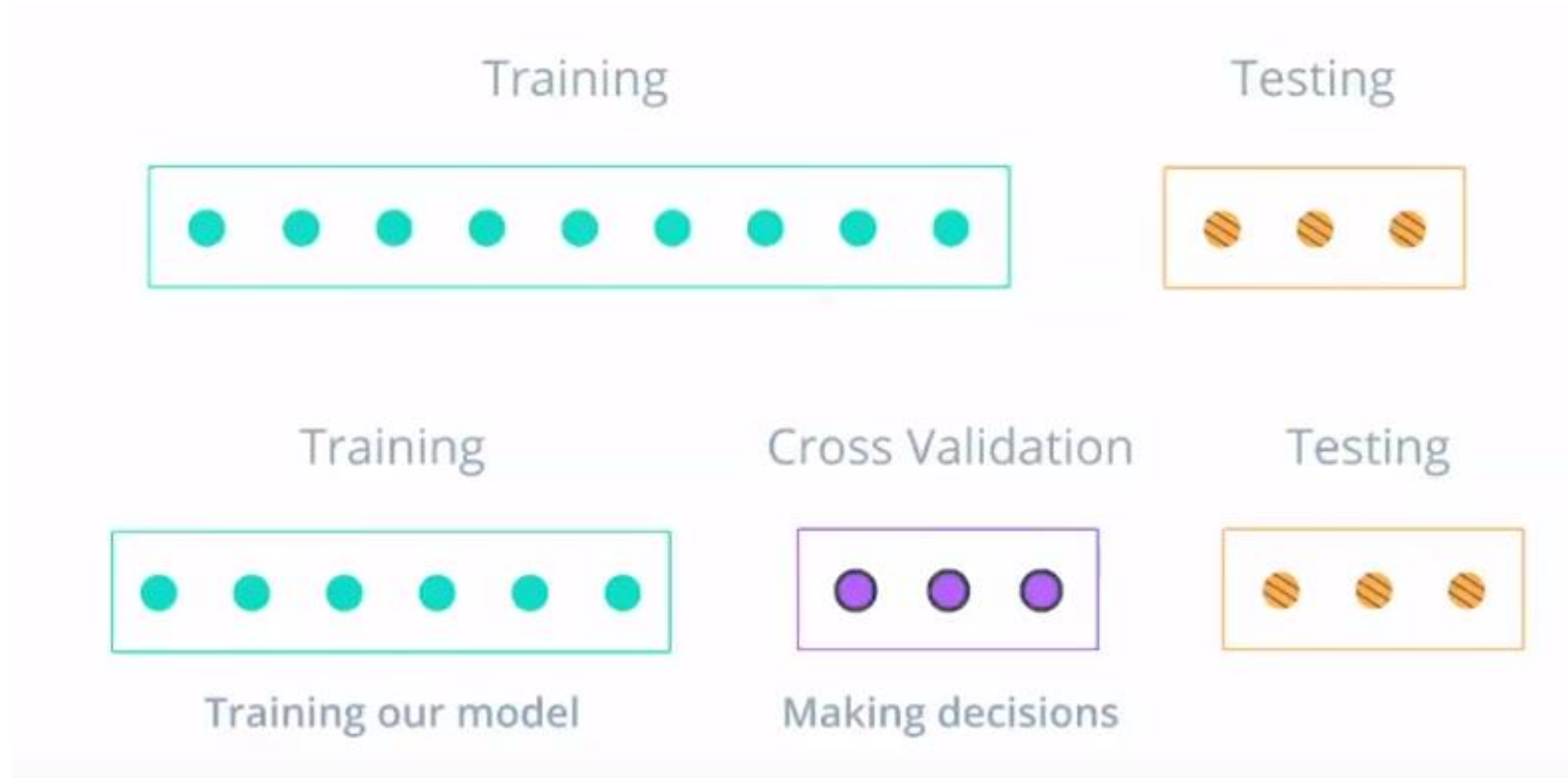LINEAR MODEL
Degree = 1

QUADRATIC MODEL
Degree = 2

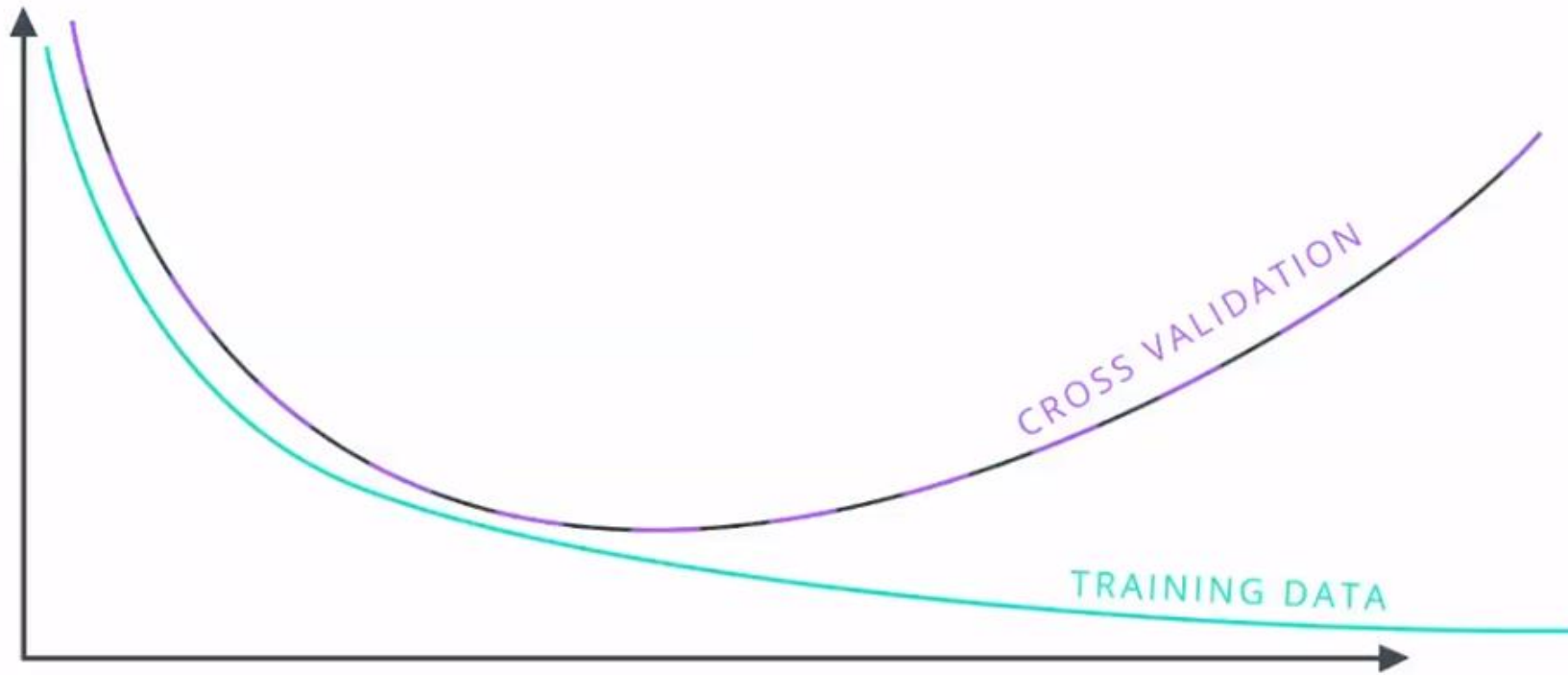POLYNOMIAL MODEL
Degree = 6

# Model Complexity

# Model Complexity
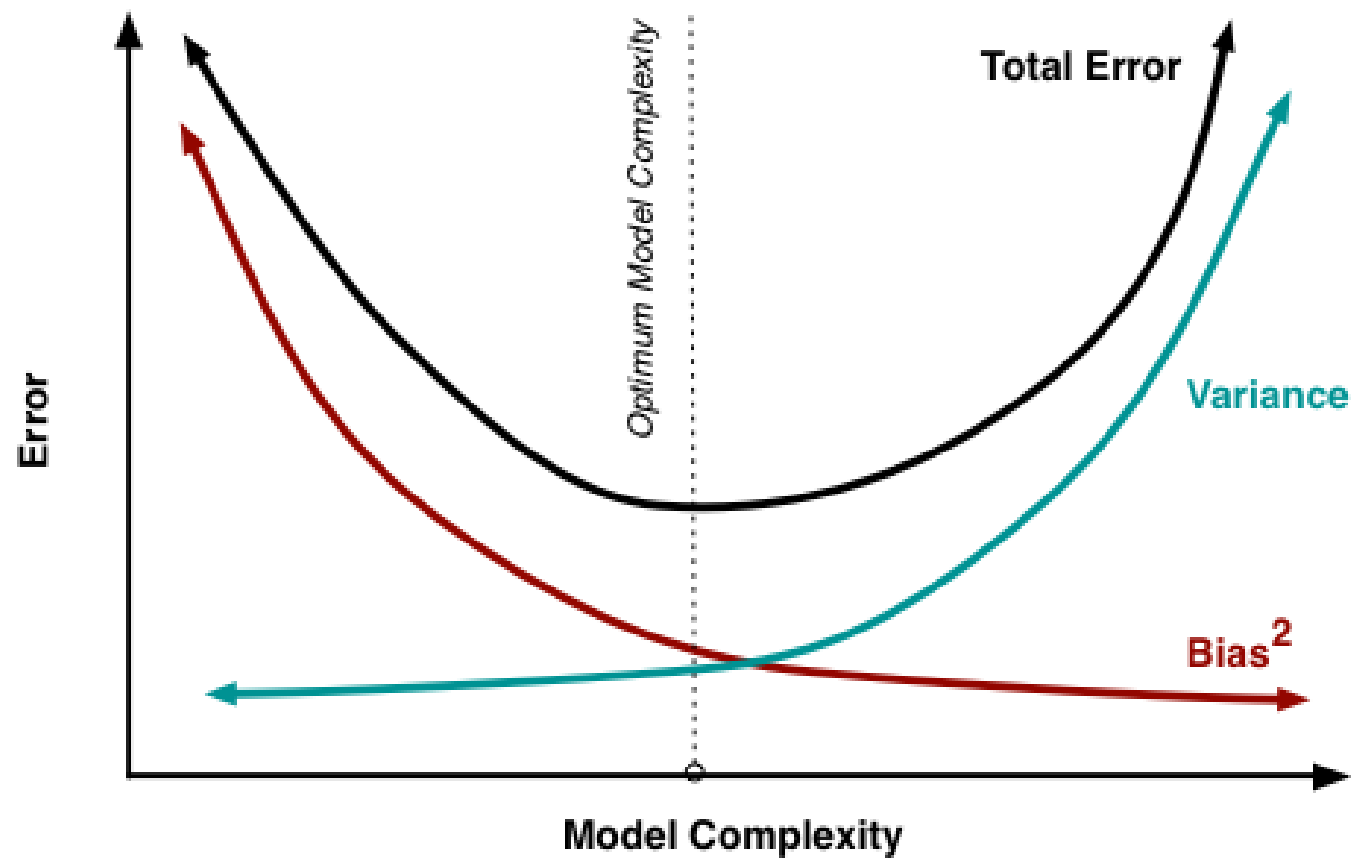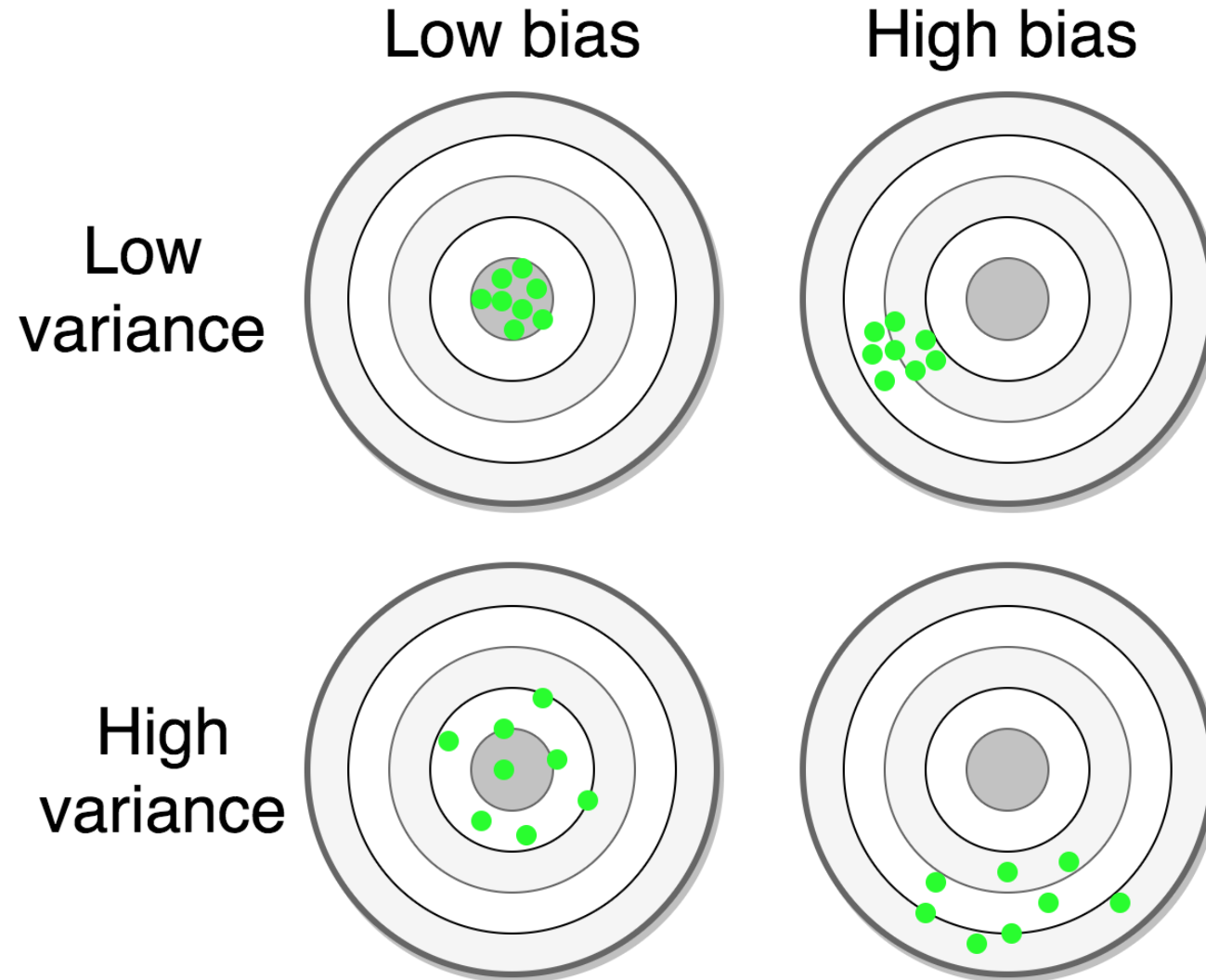
# Model Complexity

# Model Complexity

# Impact of training points

# Bias Variance Trade-off

# Bias Variance Trade-off

# Regularization / Shrinkage

Ridge / L2

$$\sum_{i=1}^{n} \left( y_i - \beta_0 - \sum_{j=1}^{p} \beta_j x_{ij} \right)^2 + \lambda \sum_{j=1}^{p} \beta_j^2 = \text{RSS} + \lambda \sum_{j=1}^{p} \beta_j^2$$
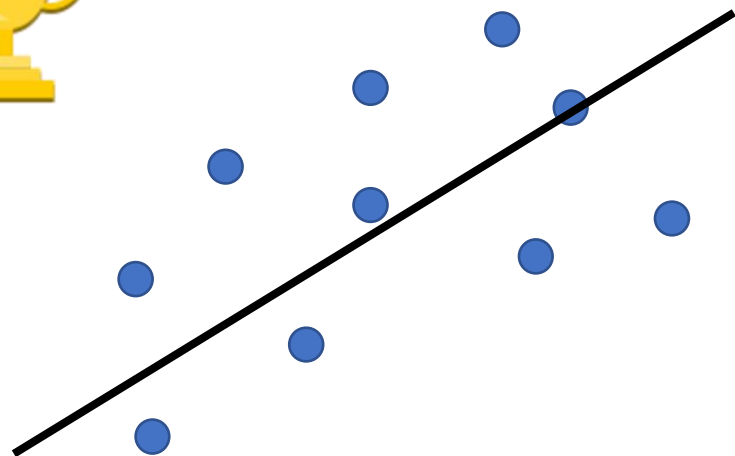
Lasso / L1

$$\sum_{i=1}^{n} \left( y_i - \beta_0 - \sum_{j=1}^{p} \beta_j x_{ij} \right)^2 + \lambda \sum_{j=1}^{p} |\beta_j| = \text{RSS} + \lambda \sum_{j=1}^{p} |\beta_j|$$
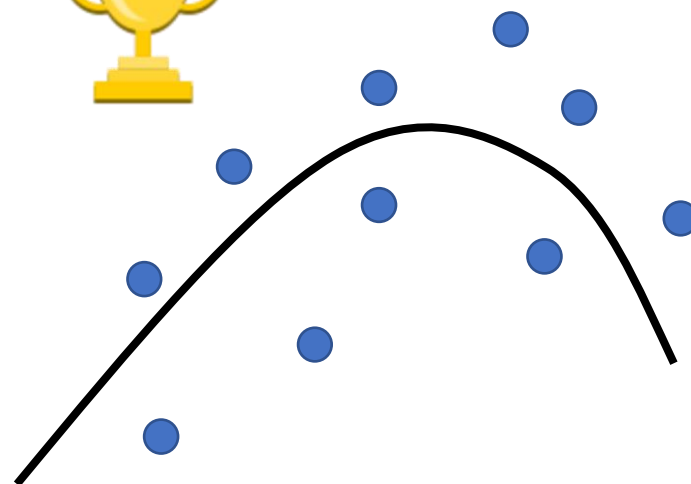
SIMPLE MODEL

COMPLEX MODEL

ERROR:

$3x_1 + 4x_2 + 5$
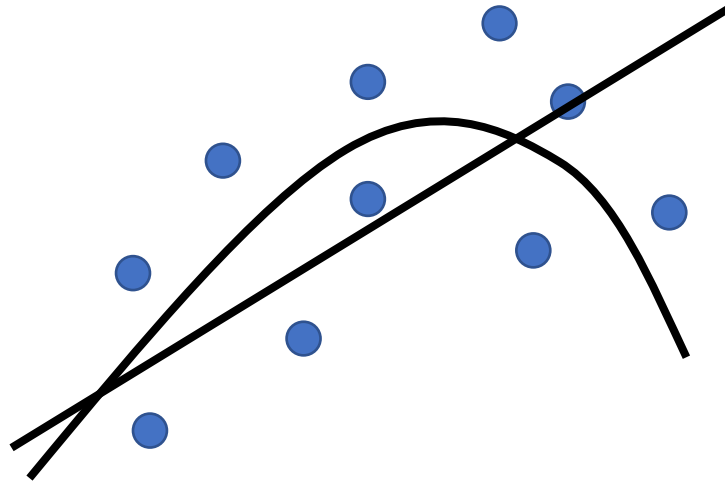
ERROR:

$2x_1^3 - 2x_1^2x_2 - 4x_2^3 + 3x_1^2 + 6x_1x_2 + 4x_2^2 + 5$

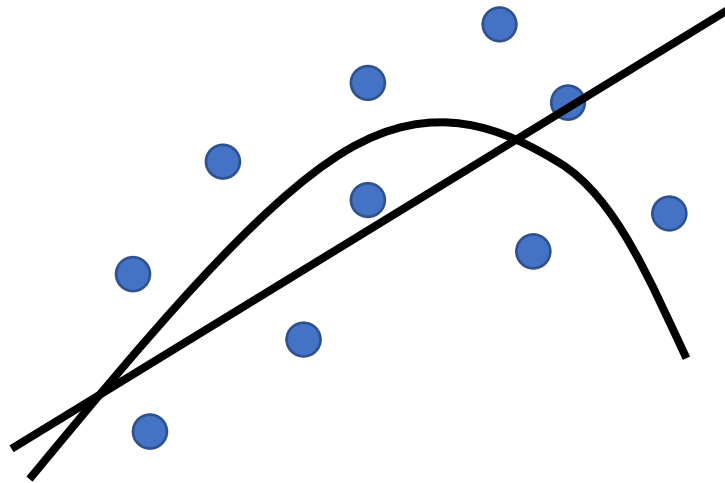# L1 (Lasso) Regularization



$2x_1^3 - 2x_1^2x_2 - 3x_1^3 + 3x_1^2 + 6x_1x_2 + 4x_2^2 + 5$

$|2| + |-2| + |3| + |4| + |6| + |4| = 21$

# L2 (Ridge) Regularization



$2x_1^3 - 2x_1^2x_2 - 4x_1^3 + 3x_1^2 5 + 6x_1x_2 + 4x_2^2 + 5$

$2^2 + (-2)^2 + (-4)^2 + 3^2 + 6^2 + 4^2 = 85$

# Simple vs Complex Models
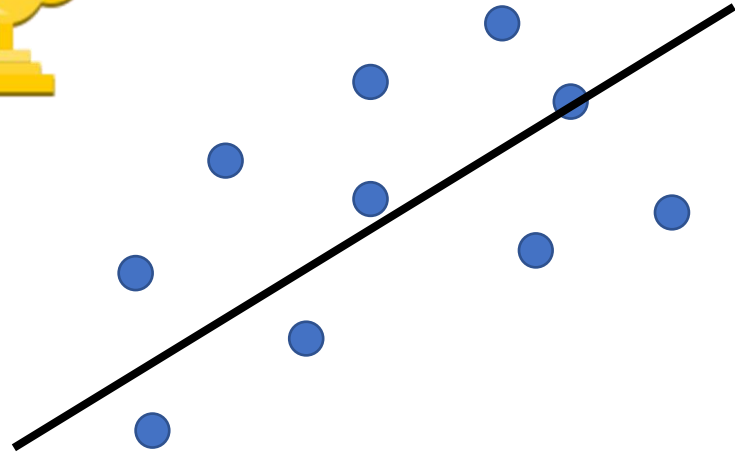
Requires LOW ERROR
OK if it's COMPLEX

PUNISHMENT on COMPLEXITY should be SMALL

Requires SIMPLICITY
OK with ERRORs

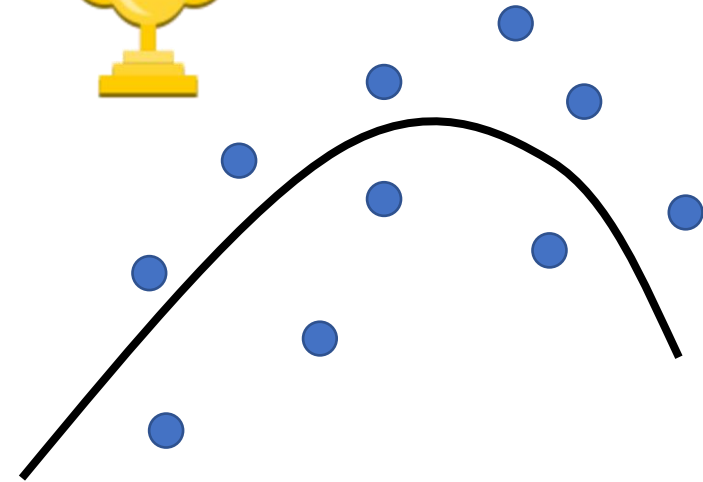PUNISHMENT on COMPLEXITY should be BIG

# The λ Parameter



λ

ERROR:

$3x_1 + 4x_2 + 5$

λ

ERROR:

$2x_1^3 - 2x_1^2x_2 - 4x_2^3 + 3x_1^2 + 6x_1x_2 + 4x_2^2 + 5$