# Regularized Regression Models

# Linear regression fitting example
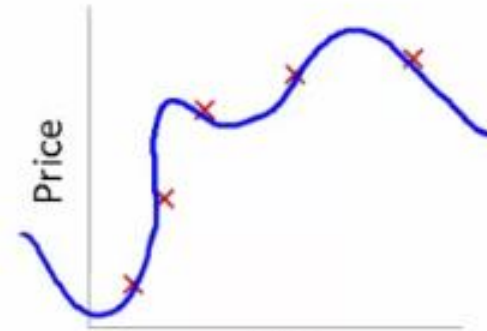


| High bias (underfit) | "Just right" | High variance (overfit) |

$$\theta_0 + \theta_1 x$$

$$\theta_0 + \theta_1 x + \theta_2 x^2$$

$$\theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3 + \theta_4 x^4$$

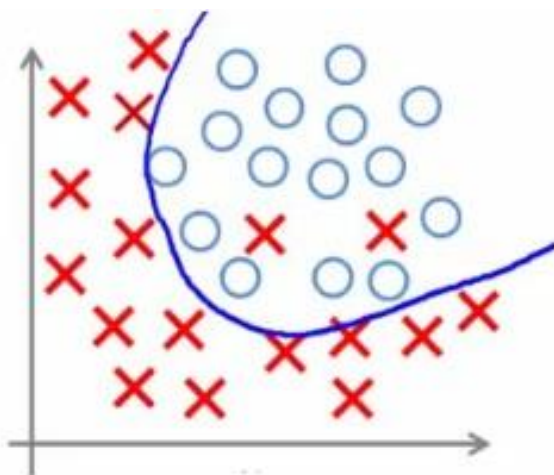**High bias (underfit)**     **"Just right"**     **High variance (overfit)**

# Logistic regression fitting example
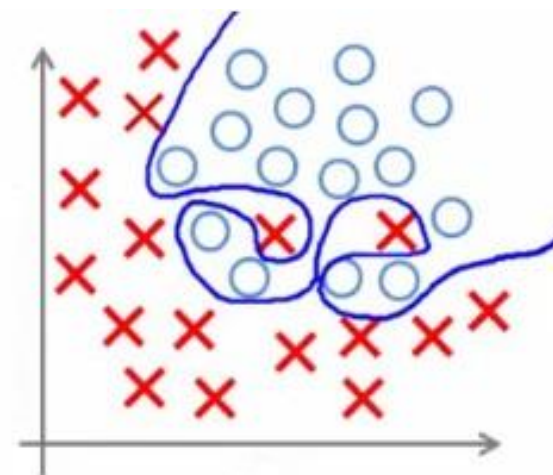


**Under-fitting**

(too simple to explain the variance)

**Appropriate-fitting**

**Over-fitting**

(forcefitting -- too good to be true)

datasciencedojo

# Overfitting

**Overfitting when......**

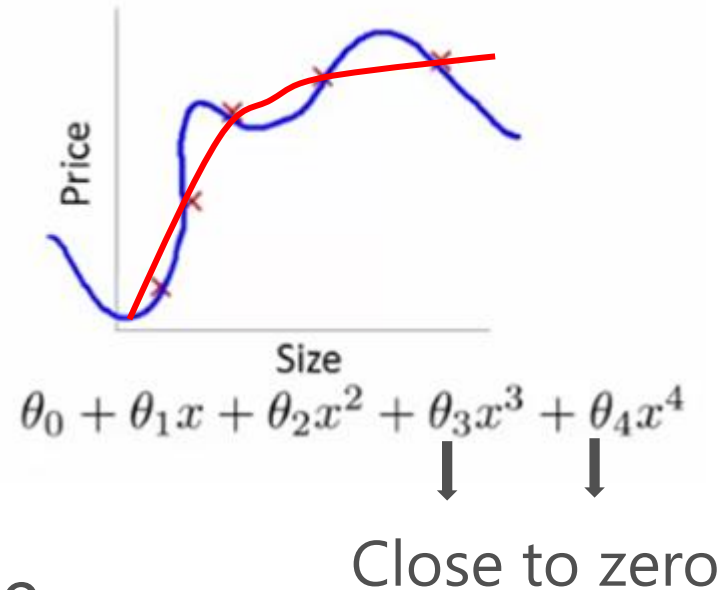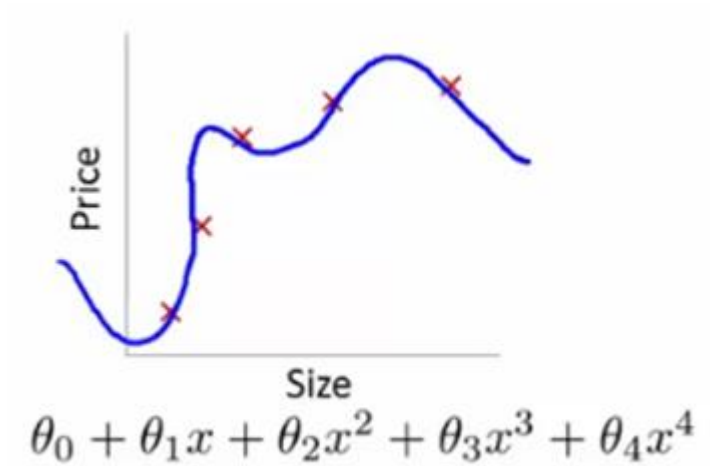- Complex model, too many features, not enough training samples.

**How to address overfitting??**

- Go through each features to decide which to keep.

- Use model selection algorithm to automatically choose features.

# Idea of Regularization

- Keep all the features, but reducing their magnitude of parameter effects in model.

- Shrink $\theta_j$ parameters

# Regularized regression intuition



$$\theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3 + \theta_4 x^4$$

$$\theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3 + \theta_4 x^4$$

Close to zero

- Goal: To minimize  cost function $\theta_j$

$$\frac{1}{2m} \sum_{i=1}^{m} (h_\theta(x^{(i)}) - y^{(i)})^2 \quad +1000\ \theta_3 + 1000\ \theta_4$$

- Suppose we penalize and make $\theta_3$ and $\theta_4$ very small

datasciencedojo

# Regularization

- Two common types of regularization in linear regression

- L2 regularization (a.k.a ridge regression)

$$\sum_{j=1}^{N}(y_j - \sum_{i=0}^{d}\theta_i \cdot x_i)^2 + \lambda\sum_{i=1}^{d}\theta_i^2$$

- L1 regularization (a.k.a lasso regression)

$$\sum_{j=1}^{N}(y_j - \sum_{i=0}^{d}\theta_i \cdot x_i)^2 + \lambda\sum_{i=1}^{d}|\theta_i|$$

datasciencedojo

# Regularized-Ridge regression

Regularization by shrink $\theta_j$ smaller values, as a result

- "less complex" hypothesis function without eliminating features
- More protection from overfitting.

L2: Ridge regression

$$J(\theta) = \frac{1}{2m}\left[\sum_{i=1}^{m}(h_\theta(x^{(i)}) - y^{(i)})^2 + \lambda\sum_{j=1}^{n}\theta_j^2\right]$$

$$\min_\theta J(\theta)$$

datasciencedojo

# Regularized-Ridge regression

$$J(\theta) = \frac{1}{2m} \left[ \sum_{i=1}^{m} (h_\theta(x^{(i)}) - y^{(i)})^2 + \lambda \sum_{j=1}^{n} \theta_j^2 \right]$$

$$\min_\theta J(\theta)$$

- Goal 1: find the best fit

- Goal 2: keep parameter $\theta_j$ small

- $\lambda$ is regularization parameter to controls a trade off

# Regularized-Ridge regression

$$J(\theta) = \frac{1}{2m} \left[ \sum_{i=1}^{m} (h_\theta(x^{(i)}) - y^{(i)})^2 + \lambda \sum_{j=1}^{n} \theta_j^2 \right]$$

$$\min_\theta J(\theta)$$

- If $\lambda$ is too large, $\theta_j$ become too small, as if features have no effect in predicting response.

- If $\lambda$ is too small, $\theta_j$ are not regularized.

datasciencedojo

# Questions?