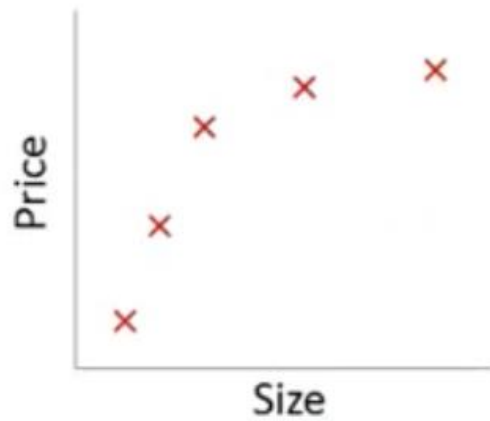
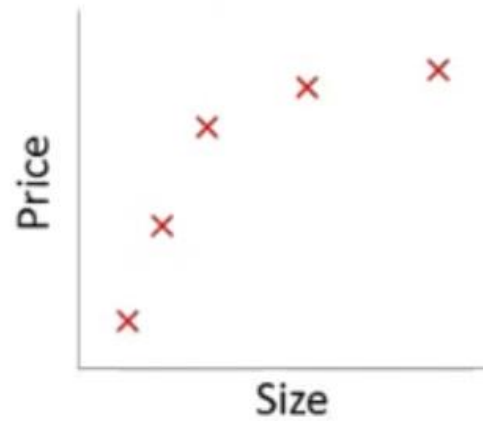


Regularization

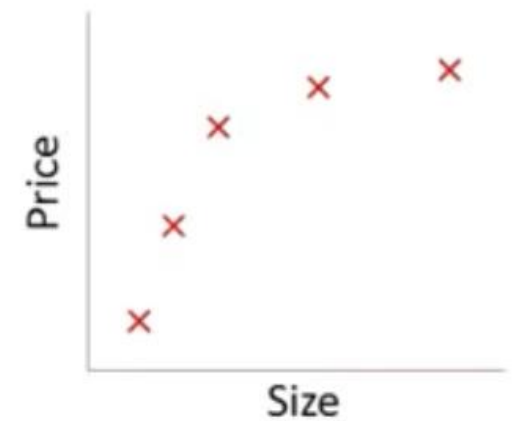
Some data about house prices



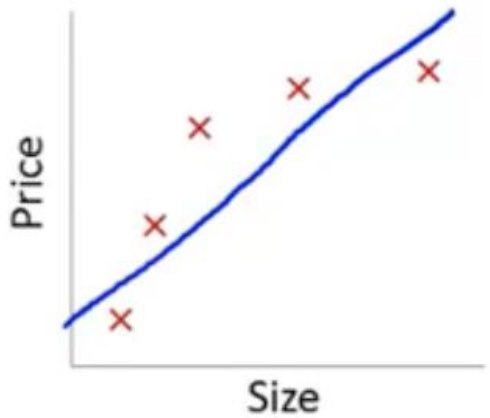
$$w_0 + w_1x$$



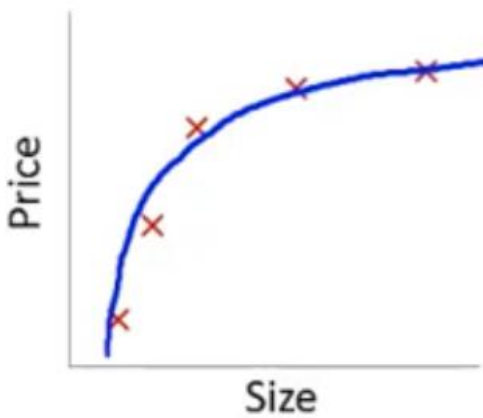
$$w_0 + w_1x + w_2x^2$$



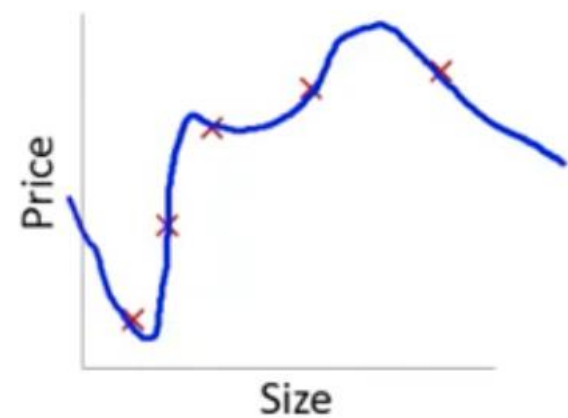
$$w_0 + w_1x + w_2x^2 + w_3x^3 + w_4x^4$$



"Underfit" "High bias"



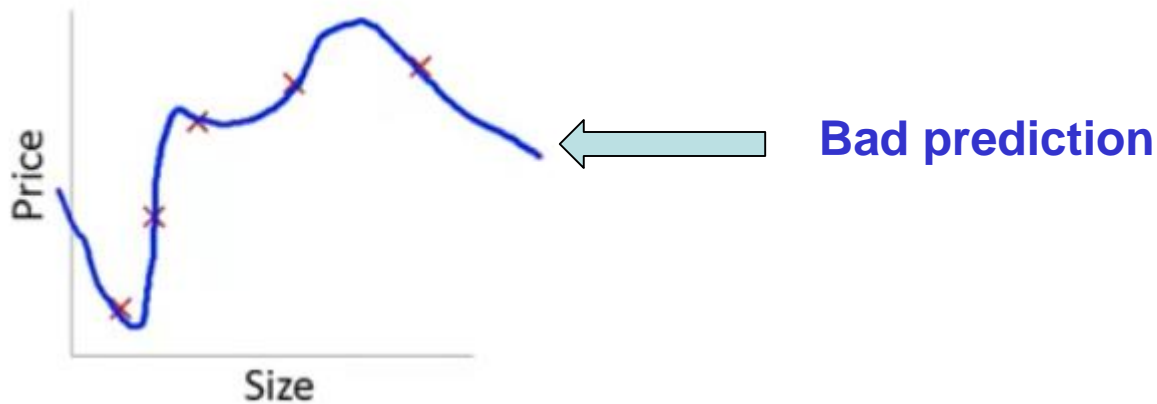
"Just right"



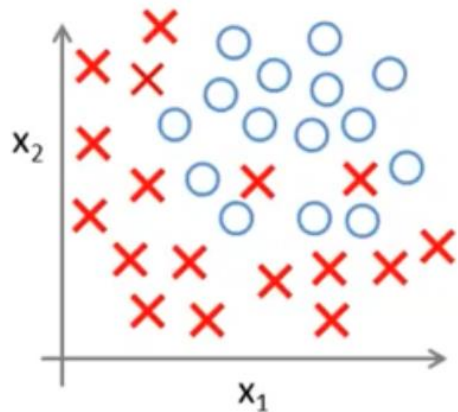
"Overfit" "High variance"

Overfitting

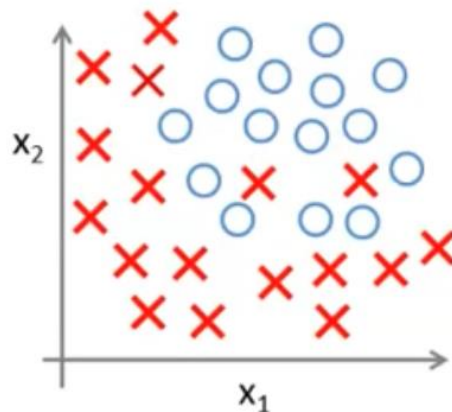
- If we have too many features, the learned hypothesis may fit the training set very well, but fails to generalize to new examples (predict prices on new examples).



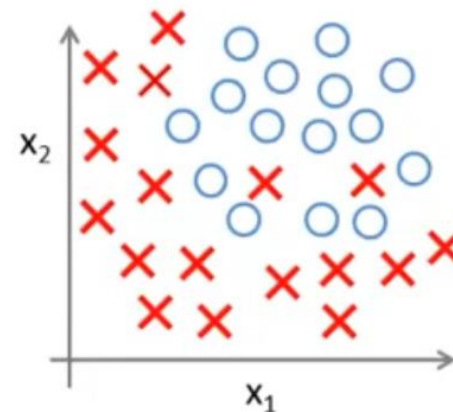
Example: Logistic Regression



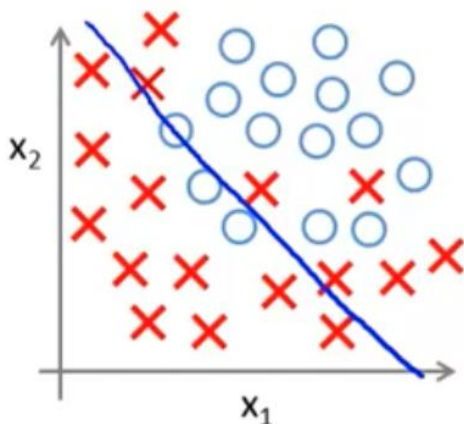
$$\sigma(w_0 + w_1x_1 + w_2x_2)$$



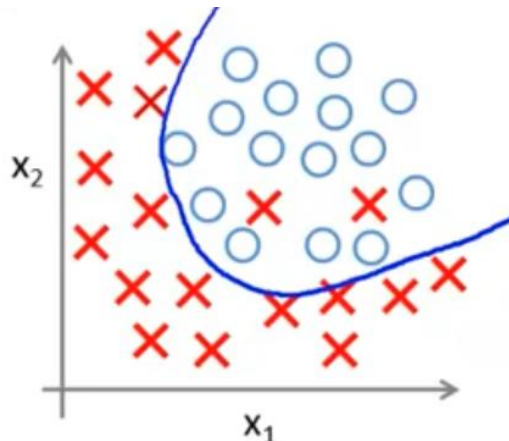
$$\sigma(w_0 + w_1x_1 + w_2x_2 + w_3x_1^2 + w_4x_2^2 + w_5x_1x_2)$$



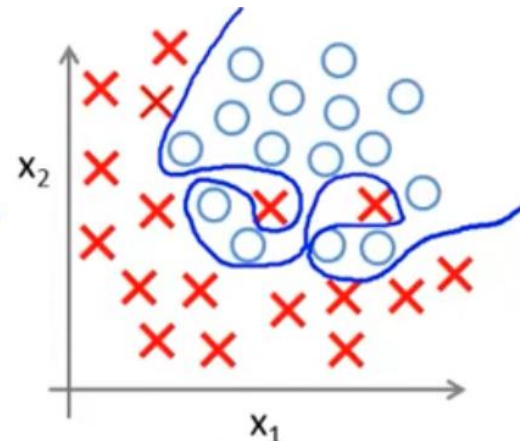
$$\sigma(w_0 + w_1x_1 + w_2x_2 + w_3x_1^2 + w_4x_2^2 + w_5x_1x_2 + w_6x_1^3x_2 + \dots)$$



"Underfit" "High bias"



"Just right"

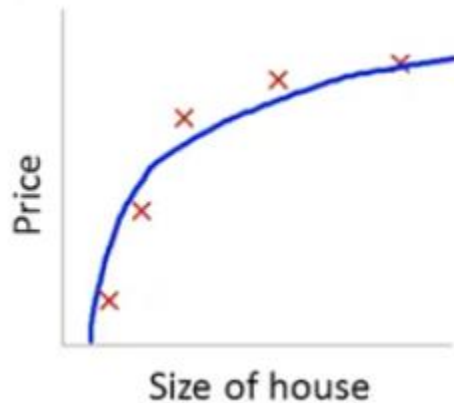


"Overfit" "High variance"

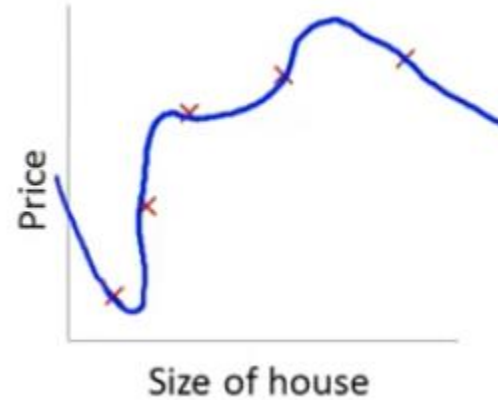
Address Overfitting

- Reduce number of features
 - Not always easy to do
- Regularization
 - Keep all the features, but reduce magnitude/values of parameters w .
 - Works well when we have a of features, each of which contributes a bit to predicting y .

Regularization - Intuition



$$w_0 + w_1x + w_2x^2$$

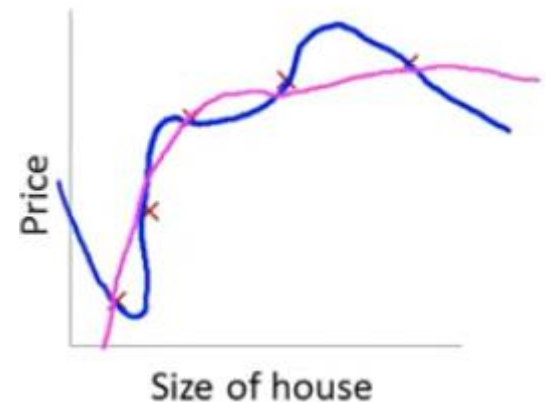


$$w_0 + w_1x + w_2x^2 + w_3x^3 + w_4x^4$$

Suppose we penalize and make w_3, w_4 small.

$$\min \frac{1}{2n} \sum_{k=1}^n (y^k - h_w(x^k))^2 + 1000w_3^2 + 1000w_4^2$$

We will end up with small w_3, w_4



Regularization

- Small values for parameters w_1, w_2, \dots, w_m
 - Simpler hypothesis
 - Less prone to overfitting
- Housing:
 - Features: x_1, x_2, \dots, x_{100}
 - Parameters: w_1, w_2, \dots, w_{100}

$$E(\mathbf{w}) = \frac{1}{2n} \left[\sum_{k=1}^n (y^k - h_{\mathbf{w}}(\mathbf{x}^k))^2 + \lambda \sum_{i=1}^m w_i^2 \right]$$

By convention, we don't include w_0 in the second sum. However, we can include it, without causing much of a difference.

Regularization

$$E(\mathbf{w}) = \frac{1}{2n} \left[\sum_{k=1}^n (y^k - h_{\mathbf{w}}(\mathbf{x}^k))^2 + \lambda \sum_{i=1}^m w_i^2 \right]$$

- λ is the **regularization parameter** controls trade off between two goals
 - 1) Want to fit the training set well
 - 2) Want to keep parameters small
- If λ is large we penalize ALL the parameters so all the parameters end up being close to 0
 - If this happens, it's like we got rid of all the terms in the hypothesis
 - This results in **underfitting**
- So, λ should be chosen carefully - not too big...
 - We look at some automatic ways to select λ later in the course



If λ is large, say 10^{10}

$$w_0 + \cancel{w_1 x} + \cancel{w_2 x^2} + \cancel{w_3 x^3} + \cancel{w_4 x^4}$$

Linear Regression – Gradient Descent


Repeat {

$$w_0 = w_0 - \alpha \frac{1}{n} \sum_{k=1}^n (h_w(\mathbf{x}^k) - y^k) x_0^k$$

$$w_j = w_j - \alpha \left(\frac{1}{n} \sum_{k=1}^n (h_w(\mathbf{x}^k) - y^k) x_j^k + \frac{\lambda}{n} w_j \right)$$

$$j = 1, 2, \dots, m$$

}


$$w_j = w_j \underbrace{\left(1 - \alpha \frac{\lambda}{n} \right)}_{\text{Less than 1, e.g. 0.99}} - \underbrace{\alpha \frac{1}{n} \sum_{k=1}^n (h_w(\mathbf{x}^k) - y^k) x_j^k}_{\text{Same as before (without regularization)}}$$

We are slightly shrinking w_j , then performing an update as before.

Less than 1,
e.g. 0.99

Same as before (without regularization)

Linear Regression – Canonical Eq.

We had:

$$\mathbf{w} = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{y}$$

Now, with regularization it becomes:

$$\mathbf{w} = \left(\mathbf{X}'\mathbf{X} + \lambda \begin{pmatrix} 0 & 0 & \cdots & 0 & 0 \\ 0 & 1 & & 0 & 0 \\ \vdots & & \ddots & \vdots & \\ 0 & 0 & & 1 & 0 \\ 0 & 0 & \cdots & 0 & 1 \end{pmatrix} \right)^{-1} \mathbf{X}'\mathbf{y}$$

- \mathbf{X} is the $n \times (m+1)$ data matrix
 - one row of $m+1$ elements for each data instance
 - without the y attribute
- \mathbf{y} is the n -vector of class values
- $\mathbf{X}'\mathbf{X}$ is $(m+1) \times (m+1)$ matrix
 - Good if number m of attributes is not too big.
- \mathbf{w} is m -vector, i.e. $(m+1) \times 1$

Regularization – Logistic Regression

- Before:

$$E(\mathbf{w}) = \frac{1}{n} \sum_{k=1}^k \ln(1 + e^{-y^k \mathbf{w}^T \mathbf{x}^k})$$

- Now:

$$E(\mathbf{w}) = \frac{1}{n} \sum_{k=1}^k \ln(1 + e^{-y^k \mathbf{w}^T \mathbf{x}^k}) + \frac{\lambda}{2n} \sum_{j=1}^m w_j^2$$

Gradient Descent Algorithm

Initialize $\mathbf{w}=\mathbf{0}$

For $t=0,1,2,\dots$ do

 Compute the gradient $\nabla_E(\mathbf{w}) = -\frac{1}{n} \sum_{k=1}^n \frac{y^k \mathbf{x}^k}{1 + e^{y^k \mathbf{w}^T \mathbf{x}^k}} + \frac{\lambda}{n} \mathbf{w}$

 Update the weights $\mathbf{w} \leftarrow \mathbf{w} - \kappa \nabla_E(\mathbf{w})$

 Iterate with the next step until \mathbf{w} doesn't change too much
 (or for a fixed number of iterations)

Return final \mathbf{w} .