

# Diagnosing ML for Applications

## Debugging a learning algorithm:

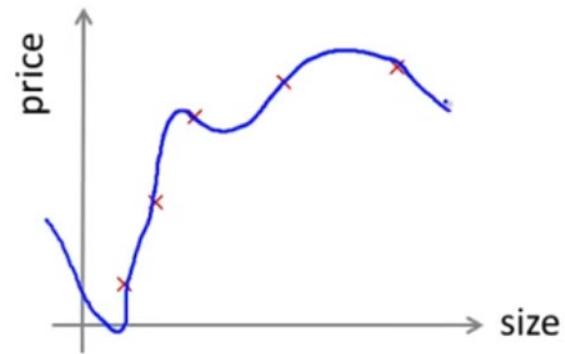
- A regularized linear regression model has been implemented for predicting house re-sale price. The parameters in the model are obtained from minimizing the following cost:

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

- However, the model does not generalize well (i.e., when the realtor uses your model to list the house on the market, the price is way off.)
- How to improve this model?
  - Collecting more data?
  - Using smaller sets of features?
  - Including additional features?
  - Using polynomial features ( $x_1^2, x_2^2, x_3^2, x_1x_2, x_1x_3, \dots$ )?
  - Increasing/decreasing  $\lambda$ ?

## How to address the problem?

- Intuitive idea – plot price against features. Not feasible.
- Run tests to evaluate a hypothesis



$$h_{\theta}(x) = \theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3 + \theta_4 x^4$$

(just one physical feature)

$x_1$  = size of house  
 $x_2$  = number of bedrooms  
 $x_3$  = area safety  
 $x_4$  = school rating  
 $x_5$  = shopping convenience  
 $x_6$  = average income of neighborhood  
 $x_7$  = newly upgraded kitchen  
.....

Then how to evaluate a hypothesis?  
And then how to select a good hypothesis?

# Model Selection – prepare datasets

Dataset:

Size	Price
2104	400
1600	330
2400	369
1416	232
3000	540
1985	300
1534	315
1427	199
1380	212
1494	243



$$(x^{(1)}, y^{(1)})$$

$$(x^{(2)}, y^{(2)})$$

$\vdots$

$$(x^{(m)}, y^{(m)})$$

$$(x_{cv}^{(1)}, y_{cv}^{(1)})$$

$$(x_{cv}^{(2)}, y_{cv}^{(2)})$$

$\vdots$

$$(x_{cv}^{(m_{cv})}, y_{cv}^{(m_{cv})})$$

$$(x_{test}^{(1)}, y_{test}^{(1)})$$

$$(x_{test}^{(2)}, y_{test}^{(2)})$$

$\vdots$

$$(x_{test}^{(m_{test})}, y_{test}^{(m_{test})})$$

Training set

~60% (random shuffle if possible)

$m$ : # of training examples

Cross validation (cv) set

~20% (random shuffle if possible)

$m_{cv}$ : # of cv examples

Testing set

~20% (random shuffle if possible)

$m_{test}$ : # of testing examples

# Model Selection (through **model evaluation**)

Training error:

$$J_{train}(\theta) = \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2 \rightarrow \text{to obtain parameters in a given hypothesis}$$

Cross Validation error:

$$J_{cv}(\theta) = \frac{1}{2m_{cv}} \sum_{i=1}^{m_{cv}} (h_{\theta}(x_{cv}^{(i)}) - y_{cv}^{(i)})^2 \rightarrow \text{to evaluate a hypothesis for model selection}$$

Test error:

$$J_{test}(\theta) = \frac{1}{2m_{test}} \sum_{i=1}^{m_{test}} (h_{\theta}(x_{test}^{(i)}) - y_{test}^{(i)})^2 \rightarrow \text{to test how well the hypothesis generalizes}$$

## Model Selection (choose the complexity of a hypothesis)

- First obtain model parameters given a structure using training set
- Then evaluate the hypothesis
- A set of models with different complexity then become available for choice

$$1. h_{\theta}(x) = \theta_0 + \theta_1 x \rightarrow \theta^{(1)} \rightarrow J_{\text{CV}}(\theta^{(1)})$$

$$2. h_{\theta}(x) = \theta_0 + \theta_1 x + \theta_2 x^2 \rightarrow \theta^{(2)} \rightarrow J_{\text{CV}}(\theta^{(2)})$$

$$3. h_{\theta}(x) = \theta_0 + \theta_1 x + \dots + \theta_3 x^3 \rightarrow \theta^{(3)} \rightarrow J_{\text{CV}}(\theta^{(3)})$$

:

$$10. h_{\theta}(x) = \theta_0 + \theta_1 x + \dots + \theta_{10} x^{10} \rightarrow \theta^{(10)} \rightarrow J_{\text{CV}}(\theta^{(10)})$$

Pick the 3<sup>rd</sup> order hypothesis if  $J_{\text{CV}}(\theta^{(3)})$  has small value

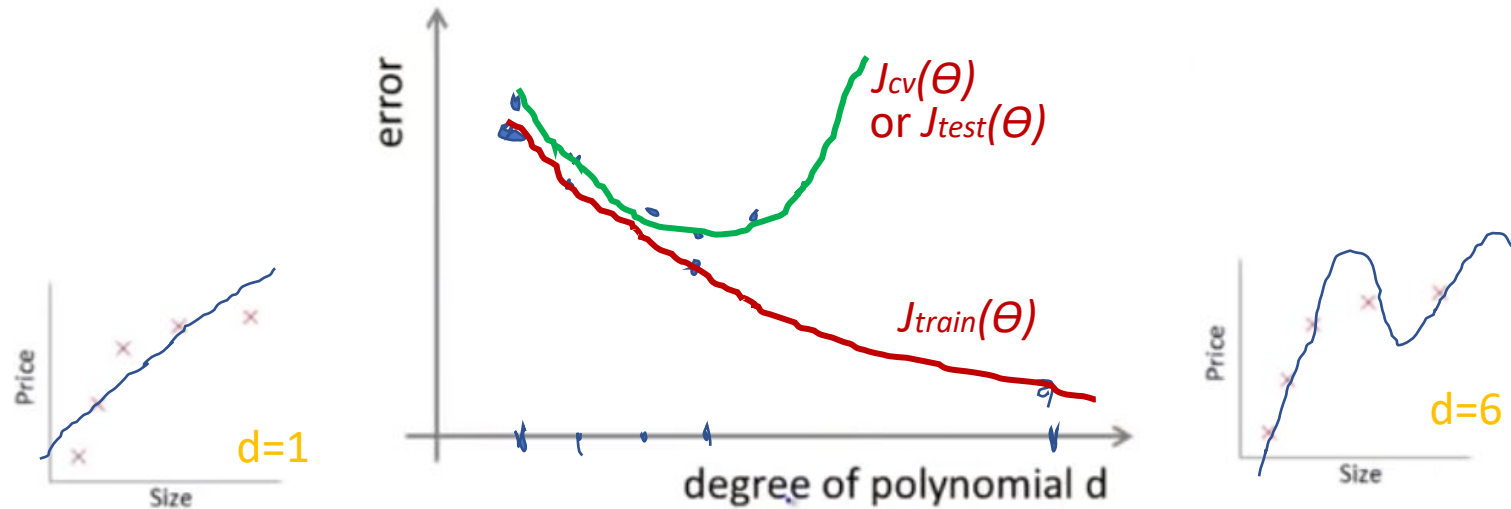
## **Diagnosing Bias vs Variance problem (classical view)**



## Bias or Variance (let's make an observation)

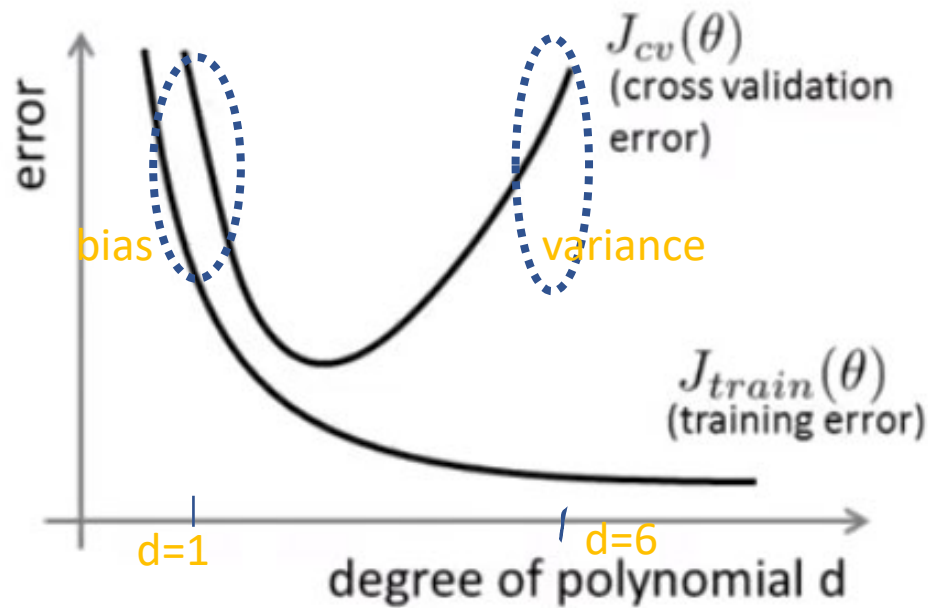
$$J_{train}(\theta) = \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2$$

$$J_{cv}(\theta) = \frac{1}{2m_{cv}} \sum_{i=1}^{m_{cv}} (h_{\theta}(x_{cv}^{(i)}) - y_{cv}^{(i)})^2$$



# Diagnosis of a bias or variance problem

- Realizing that there is a problem with the hypothesis/model you generated
- Either  $J_{cv}(\theta)$  or  $J_{test}(\theta)$  is high
- Two possibilities: model too simple (underfit) or model too complex (overfit)



Bias problem (underfit)

$J_{train}(\theta)$  - high

$J_{cv}(\theta) \sim J_{train}(\theta)$

Variance problem (overfit)

$J_{train}(\theta)$  - low

$J_{cv}(\theta) \gg J_{train}(\theta)$

**With the bias vs variance diagnosis tool, we can now  
apply it to the diagnosis of a learning problem  
(use regularization problem as an example)**

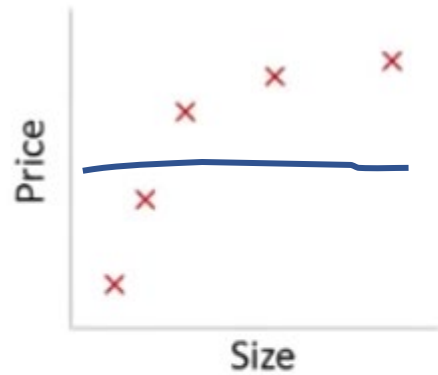
Recall the regularized linear regression problem

$$h_{\theta}(x) = \theta^T x = \theta_0 x_0 + \theta_1 x_1 + \theta_2 x_2 + \cdots + \theta_n x_n$$

The parameters  $\theta_0 \theta_1 \dots \theta_n$  are determined from minimizing the following cost function,

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

# How to choose $\lambda$

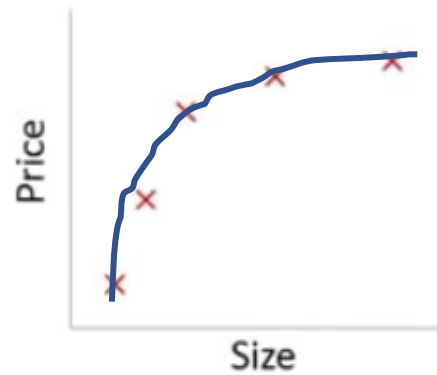


Large  $\lambda$

High bias (underfit)

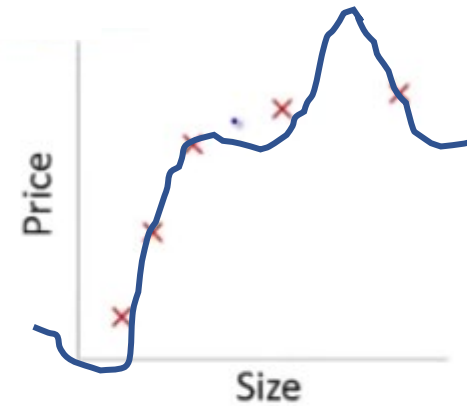
$= 10000. \theta_1 \approx 0, \theta_2 \approx 0, \dots$

$h_{\theta}(x) \approx \theta_0$



Intermediate  $\lambda$

"Just right"



Small  $\lambda$

High variance (overfit)

Given a hypothesis with parameters  $\theta_0 \theta_1 \dots \theta_n$

$$h_{\theta}(x) = \theta^T x = \theta_0 x_0 + \theta_1 x_1 + \theta_2 x_2 + \dots + \theta_n x_n$$

Obtain parameters  $\theta_0 \theta_1 \dots \theta_n$  from minimizing the following cost function,

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

### Evaluation of a Hypothesis

For selecting an appropriate regularization parameter  $\lambda$ , use the following costs:

$$\text{Training: } J_{\text{train}}(\theta) = \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2$$

$$\text{Cross Validation: } J_{\text{cv}}(\theta) = \frac{1}{2m_{\text{cv}}} \sum_{i=1}^{m_{\text{cv}}} (h_{\theta}(x_{\text{cv}}^{(i)}) - y_{\text{cv}}^{(i)})^2$$

$$\text{Test: } J_{\text{test}}(\theta) = \frac{1}{2m_{\text{test}}} \sum_{i=1}^{m_{\text{test}}} (h_{\theta}(x_{\text{test}}^{(i)}) - y_{\text{test}}^{(i)})^2$$

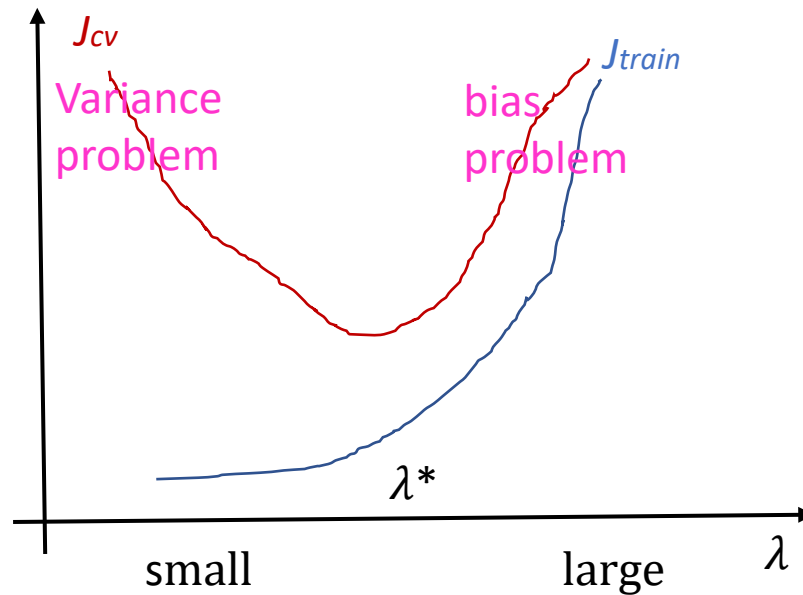
## Choosing the regularization parameter $\lambda$

1. Try  $\lambda = 0$   $\rightarrow \theta^{(1)} \rightarrow J_{\text{cv}}(\theta^{(1)})$
2. Try  $\lambda = 0.01$   $\rightarrow \theta^{(2)} \rightarrow J_{\text{cv}}(\theta^{(2)})$
3. Try  $\lambda = 0.02$   $\rightarrow \theta^{(3)} \rightarrow J_{\text{cv}}(\theta^{(3)})$
4. Try  $\lambda = 0.04$   $\rightarrow \theta^{(4)} \rightarrow J_{\text{cv}}(\theta^{(4)})$
5. Try  $\lambda = 0.08$   $\rightarrow \theta^{(5)} \rightarrow J_{\text{cv}}(\theta^{(5)})$
- ...
20. Try  $\lambda = 10$   $\rightarrow \theta^{(20)} \rightarrow J_{\text{cv}}(\theta^{(20)})$

Pick  $\lambda = 0.08$  say if  $J_{\text{cv}}(\theta^{(5)})$  has small value

# Diagnosis of a bias or variance problem

- Realizing that there is a problem with the hypothesis/model, which is related to choice of  $\lambda$
- Either  $J_{cv}(\theta)$  or  $J_{test}(\theta)$  is high
- Two possibilities: model too simple/very large  $\lambda$  (underfit) or model too complex/very small  $\lambda$  (overfit)



Bias problem (underfit)

$J_{train}(\theta)$  - high

$J_{cv}(\theta) \sim J_{train}(\theta)$

Variance problem (overfit)

$J_{train}(\theta)$  - low

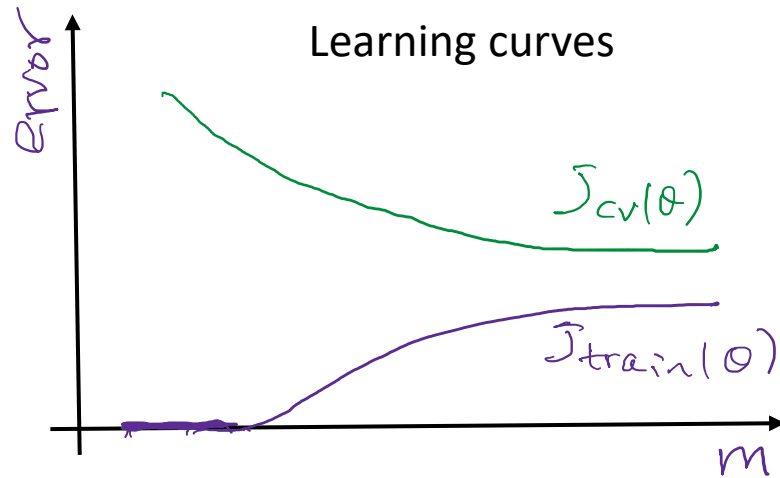
$J_{cv}(\theta) \gg J_{train}(\theta)$



# Impact of # of samples on bias-variance relationship

Consider a hypothesis with given structure/complexity:

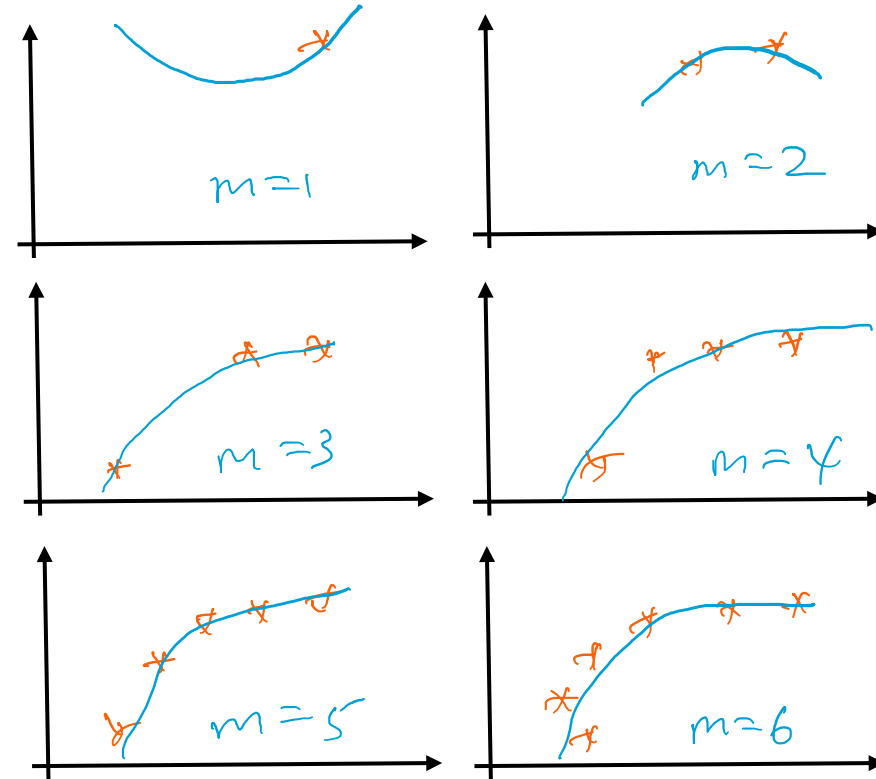
$$h_{\theta}(x) = \theta_0 + \theta_1 x + \theta_2 x^2$$



$$J_{train}(\theta) = \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2$$

$$J_{cv}(\theta) = \frac{1}{2m_{cv}} \sum_{i=1}^{m_{cv}} (h_{\theta}(x_{cv}^{(i)}) - y_{cv}^{(i)})^2$$

Illustration of hypothesis vs # of samples ( $m$ )



Learning curves associated with high bias / high variance

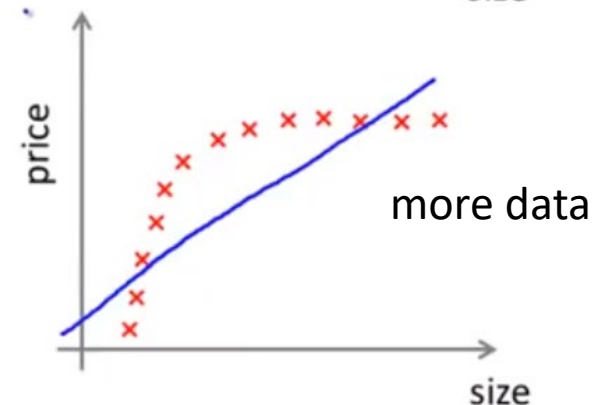
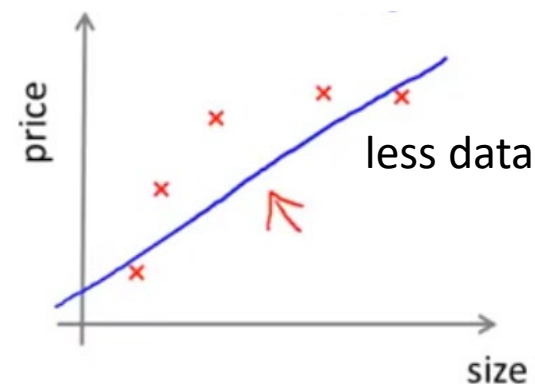
# High bias problems



If learning suffers from a high bias problem, getting more training data by itself is not likely helpful

Given a hypothesis

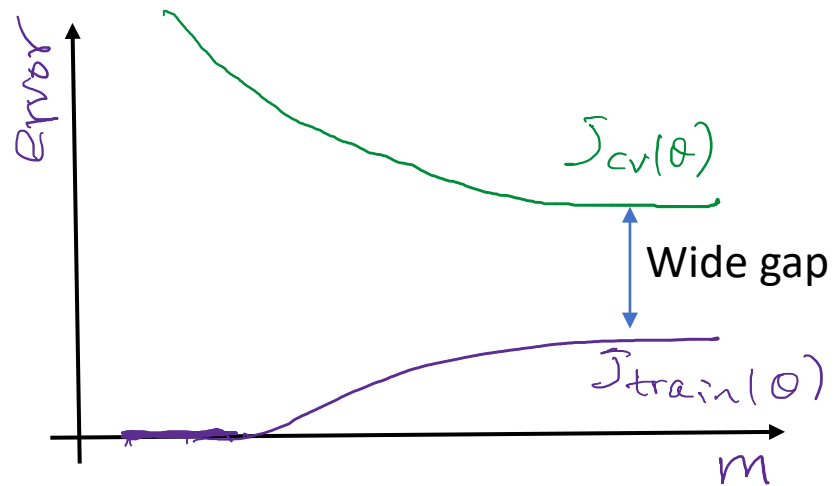
$$h_{\theta}(x) = \theta_0 + \theta_1 x$$



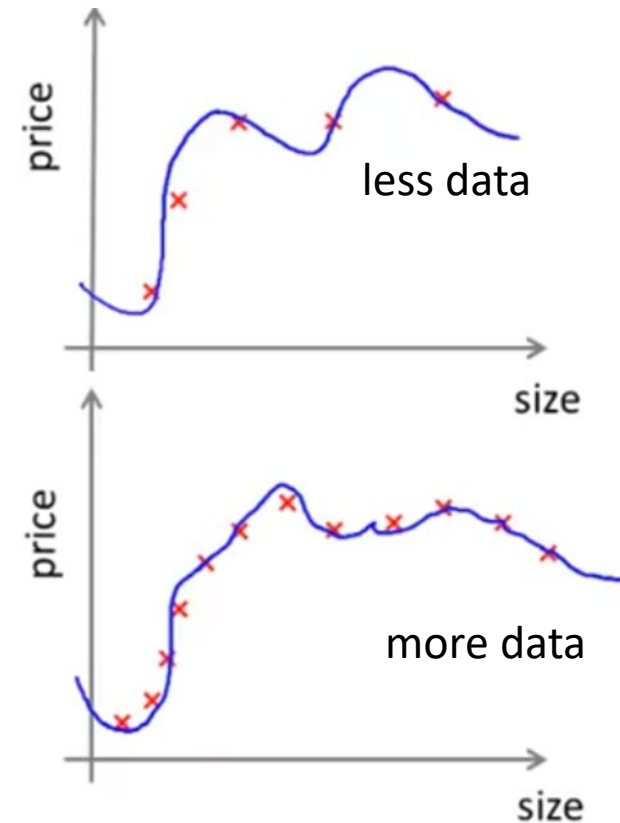
# High variance problems

Given a hypothesis (the parameters of which are obtained using a small  $\lambda$  for regulation)

$$h_{\theta}(x) = \theta_0 x_0 + \theta_1 x_1 + \theta_2 x_2 + \cdots + \theta_{100} x_{100}$$



If learning suffers from a high variance problem, getting more training data is likely helpful



## Debugging a learning algorithm:

- A regularized linear regression model has been implemented for predicting house re-sale price. The parameters in the model are obtained from minimizing the following cost:

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

- However, the model does not generalize well (i.e., when the realtor uses your model to list the house on the market, the price is way off.)
- How to improve this model?
  - Collecting more data? addresses high variance problem
  - Using smaller sets of features? addresses high variance problem
  - Including additional features? addresses high bias
  - Using polynomial features ( $x_1^2, x_2^2, x_3^2, x_1x_2, x_1x_3, \dots$ )? addresses high bias
  - Increasing  $\lambda$ ? addresses high variance
  - Decreasing  $\lambda$ ? addresses high bias

# Modern development

- Rich models such as deep neural networks can fit or interpolate the data really well
- Current evidence shows that they also are accurate on test data
- The “double-descent” performance curve instead of the U-shaped bias–variance trade-off curve, for beyond the point of interpolation
- On-going research to understand the phenomenon

