

→ Diyelim ki predicting housing prices örneğini yapıyoruz. Regularized linear regression uyguluyoruz.

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

→ Diyelim ki hipotezini new set of houses üzerinde test ediyor sun, ama bir bakıyonsun ki çok büyük hatalar var. Ne yapacağız? Learning algorithmi nasıl geliştireceğiz? Bazı yollar şunlar:

- Get more training examples (her zaman işe yaramaz)
- Try smaller sets of features ( $x_1, x_2, \dots, x_n$  çok fazla ise overfitting olabilir. Arolarından seçilebilir.)
- Try getting additional features (belki eldeki featurelar yetmez)
- Try adding polynomial features ( $x_1^2, x_2^2, x_1 \cdot x_2$ , etc.)
- Try decreasing  $\lambda$
- Try increasing  $\lambda$

• Peki hangisini yapmamız gerektiğini nasıl bileceğiz? Çoğu insan bunların arasından rastgele seçim yapıyor

- Fortunately, there is a technique (simple one) that can let you very quickly rule out half of the things on this list as being potentially promising things to pursue (alemek)
- Bu tekniği uygulayarak bu seçeneklerin bir çoğunu eleyebiliriz.

İlenleyen videolarda bu tekniklerden bahsedilecek, they are called as:

### Machine Learning Diagnostics

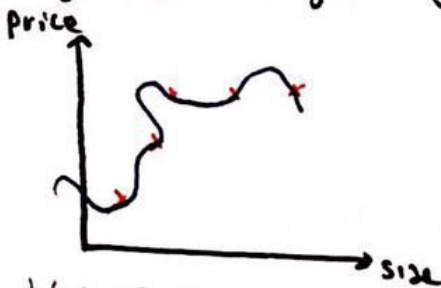
Diagnostic: A test that you can run to gain insight what is/isn't working with a learning algorithm, and gain guidance as to how best to improve its performance.

Diagnostics can take time to implement, but doing so can be a very good use of your time.



## W6 - Evaluating a Hypothesis

→ Bu videoda eğitilmiş hipotezimizi değerlendirmeyi göreceğiz. İstenilen videolarda underfitting ve overfitting'i nasıl anlayabileceğimizi de göreceğiz.



$$h(x) = \theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3 + \theta_4 x^4$$

• Bizim şimdiye kadar odaklandığımız şey cost function'u minimize eden  $\theta$ 'lar ile hipotezi oluşturmak.

• Şunu da biliyoruz ki Low Training Error demek kesinlikle iyi bir hipotez oluşturdu demek değil. Hipotez overfitting yapabilir.

• Fails to generalize to new examples not in training set.

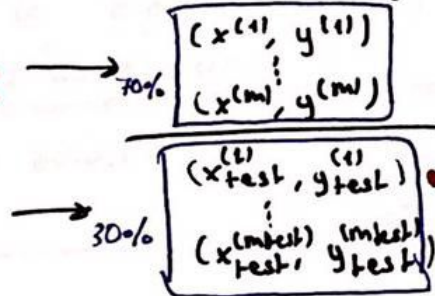
Hipotezimin overfit edip etmediğini nasıl bileceğim? Yukarıdaki örnekte 2 feature olduğu için plot ile görebiliriz ama daha karmaşık problemlerde bunu visualize etmek mümkün değil.

→ O zaman hipotezi değerlendirmek için bise başka bir yol gerek.

Standard way to Evaluate a Learned Hypothesis is as follows:

Size	Price
2104	400
1600	330
2400	369
1416	232
3000	540
1785	300
1534	315
1427	199
1380	212
1494	243

• Yandaki gibi bir training set olsun. Evaluation of hypothesis için bunu ikiye bölüyoruz.



•  $m_{test} \neq$  test examples

✗ Eğer data ordered ise 70% - 30% split'i random yapmamız faydalıdır olur!

### Training/Testing procedure for Linear Regression

• Learn parameter  $\theta$  from training data (minimizing training error  $J(\theta)$ ) →  $J(\theta)$  70% lik data kullanılarak oluşturuluyor!

Eğitilmiş parametre kullanılıp test yapılıyor.

• Compute test set error:

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

• Logistic Regression ile Classification yapılıyor olsaydı:

• Learn parameter  $\theta$  from training data

• Compute Test Set Error:

$$J_{test}(\theta) = \frac{1}{m_{test}} \sum_{i=1}^{m_{test}} y_{test}^{(i)} \log h(x_{test}^{(i)}) + (1 - y_{test}^{(i)}) \log h(x_{test}^{(i)})$$

**Burada yapılan:**  
Training set "en" minimize edilen  $J(\theta)$  y. "en" bir set "en" ile elde edilen en iyi  $\theta$  ile deneyip bir test  $J_{test}(\theta)$  bulunur.



→ Eğitilen parametrelerin hatasını anlayabilmek için  $J_{test}(\theta)$  ya test training ex.s kullanarak hesaplamak mantıklı, bunun yanında alternatif yöntemler de kullanılabilir:

Misclassification Error (0/1 misclassification error):

$$\text{err}(h(x), y) = \begin{cases} 1 & \text{if } h(x) \geq 0.5, y=0 \\ & \text{or if } h(x) < 0.5, y=1 \\ 0 & \text{otherwise} \end{cases}$$

$\text{err}(h(x), y)$  bir tahminin hatasını temsil ediyor eğer tahmin ( $h(x) \geq 0.5$ ) 1 ise ama gerçek sonuç 0 ise error 1 olur veya tahmin 0 ama gerçek sonuç 1 ise.

$$\text{Test error} = \frac{1}{n_{\text{test}}} \sum_{i=1}^{n_{\text{test}}} \text{err}(h(x_{\text{test}}^{(i)}), y^{(i)})$$

Yani sonuçta hipotezin yanlış bildiklerinin sayısının tüm test examples'a oranı!

→ Bu sistem  $w_4, w_5$  ile predict deinde kullanılmıştır.

ilerleyen videolarda bu öğrenilenleri de kullanarak nasıl feature seçileceği, degree & polynomial to use with learning algorithm, nasıl regularization parametre  $\lambda$  seçileceğini göreceğiz...



How to decide what degree of polynomial to fit a data set? What features to include? How to choose regularization parameter  $\lambda$ ? These are called model selection problems.

→ Daha önce de gördük ki, just because the trained hypothesis fits the training set very well doesn't mean it is a good hypothesis çünkü overfitting olabilir.

• More generally, this is why training set's error ( $J(\theta)$  for trained  $\theta$ ) is not a good predictor for how well the hypothesis do in new examples.

• Yani hipotezin eğitildiği setteki hatası genelde actual generalization error'dan daha düşüktür. Yeni tahminlerde daha başarılıdır demek.

### Let's consider the Model Selection Problem

Let's say you're trying to choose what degree polynomial to fit to data.

$$1.) h(x) = \theta_0 + \theta_1 x \quad d=1 \quad 2.) h(x) = \theta_0 + \theta_1 x + \theta_2 x^2 \quad d=2 \quad \dots \quad 10.) h(x) = \dots \quad d=10$$

•  $\theta$  parametresinin yanında bir parametreyi daha bulmaya çalışıyoruz gibi düşün:  $d = \text{degree of polynomial}$ .

Let's say: Yapmak istediğim model seçip, modeli eğitmek ve eğitilen modelin yeni örneklerle ne kadar iyi oturduğunu tahmin etmek istiyorum. Bunun için ne yapabilirim?

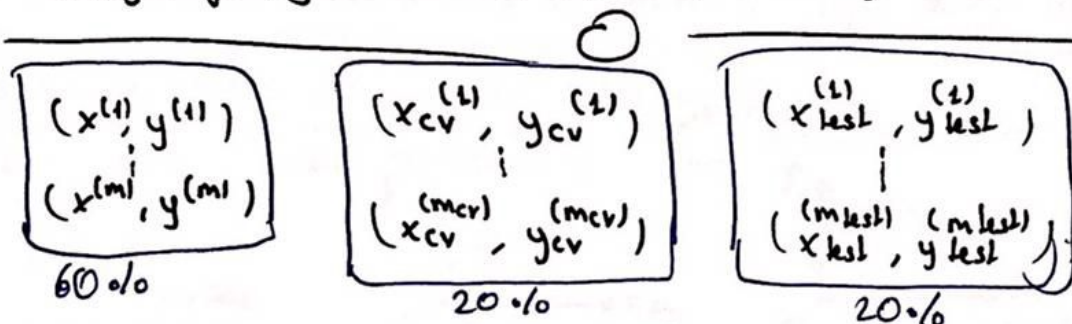
$d=1, d=2, \dots, d=10$  olan hipotezler için ayrı ayrı training set ile eğitimi yaparım hen biri için ayrı bir  $\theta$  elde ederim.  $\theta^{(1)}, \theta^{(2)}, \dots$  şeklinde. Daha sonra hangisinin daha iyi olduğunu anlamak için Test Set'i kullanarak hepsinin performansını test ederim  $J_{\text{test}}(\theta^{(1)}), \dots, J_{\text{test}}(\theta^{(10)})$  hesaplanır. En iyi olanı seçerim!

Mesela  $J_{\text{test}}(\theta^{(5)})$  en düşük düzeyim o zaman  $d=5$  seçtim. Şimdi bu son hipotezin generalized error'unu tahmin etmek istiyorum, nasıl yapacağım? Eğer test set'i kullanırsak sağma olacak çünkü  $d=5$  parametresini zaten test set'e göre seçtik. Yeni training ex.s için iyi bir tahmin yapamayız. Bu yüzden training set'i 2 yerine 3'e bölmeliyiz

Training Set / Cross Validation Set / Test Set -



- ⊗ Sonuçta Training Set 3 parçaya bölünecek %60 - %20 - %20
1. part: Training Set,  $\Theta$  ları eğitmek için kullanılacak
  2. part: Cross Validation Set, bu set ile  $\Theta$  lar test edilecek ve sonucunda  $d$  parametresine karar verilecek.
  3. part: Test Set, bu set ile karar verilen sonuç hipotezinin ( $\Theta$  ve  $d$  'den oluşuyor) yeni training ex.s için nasıl bir performans sergileyeceği simüle edilmiş oluyor



### Train / Test / Validation Errors

Training Error:

$$J(\Theta)_{\text{train}} = \frac{1}{2m} \sum_{i=1}^m (h(x^{(i)}) - y^{(i)})^2$$

Cross Validation Error:

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

Test Error

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

Sonuçta şu yol izlenecek:

- ⊙ Farklı hipotezler için Training Set kullanılarak  $\Theta$  lar bulunacak!
  - 1)  $h(x) = \Theta_0 + \Theta_1 x \rightarrow \min_{\Theta} J(\Theta) \rightarrow \Theta^{(1)} \rightarrow J_{cv}(\Theta^{(1)})$
  - 2)  $h(x) = \Theta_0 + \Theta_1 x + \Theta_2 x^2 \rightarrow \Theta^{(2)} \rightarrow J_{cv}(\Theta^{(2)})$
  - 10)  $h(x) = \Theta_0 + \dots + \Theta_{10} x^{10} \rightarrow \Theta^{(10)} \rightarrow J_{cv}(\Theta^{(10)})$
- ⊙ Her hipotez için  $J_{cv}$  hesaplanacak (cv set kullanılarak) min  $J_{cv}$  seç!
- ⊙ Sonuçta seçilen hipotezin  $\Theta$  ve  $d$  belli, generalized error est.  $J_{test}(\Theta)$  hesaplanır. (Test Set kullanılarak)

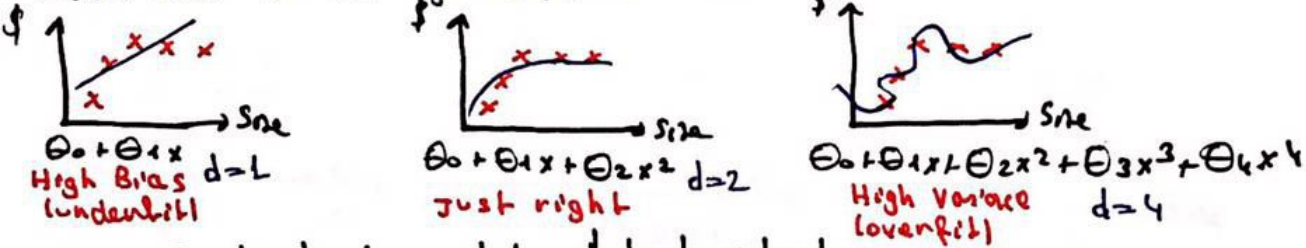


→ If you run a learning algorithm and it doesn't do as well as you are hoping, almost **EVERYTIME**, it will be because you have either high bias problem ~~or~~ high variance problem. In other words either overfitting problem ~~or~~ underfitting problem

○ Bu durumda hangisinin olduğunu veya ikisinden de as az mı olduğunu bilmek çok işe yaran çünkü bunları algoritmayı nasıl geliştirebileceğimizi hakkında güçlü işaretlerdir.

○ Bu videoda overfitting ve underfitting'e daha detaylı değinerek ve bir learning algoritmasının ikisinden birine sahip olup olmadığını nasıl değlendirebileceğimizi göreceğiz.

Daha önce şunları görmüştük:



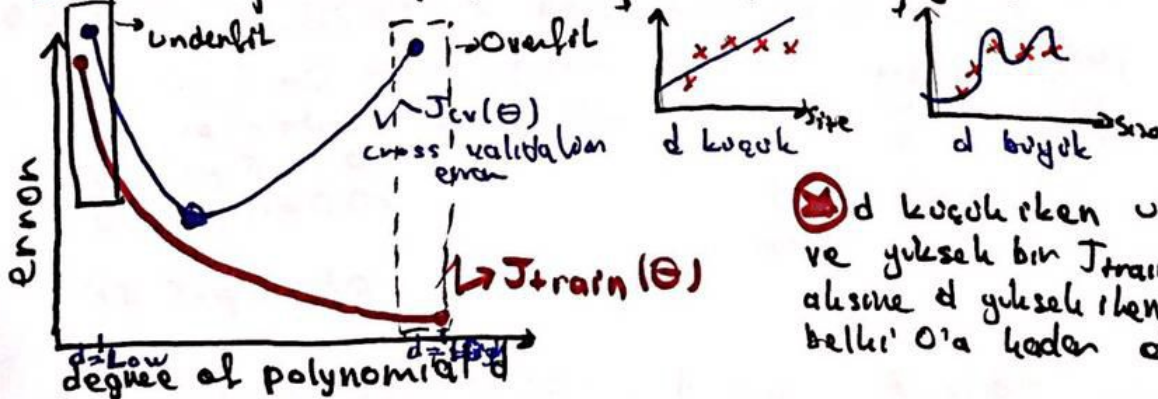
Şimdi bunları daha detaylı anlatalım.

Bias / Variance : Let's assume Training Error and Validation Error is defined as:

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

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

○ Şimdi degree of polynomial  $d$  karşı error grafiklerini çizelim:



○  $d$  küçükken underfitting olur ve yüksek bir  $J_{\text{train}}(\theta)$  bulursak  $d$  yüksekken ise  $J_{\text{train}}(\theta)$  belli  $0$ 'a kadar düşebilir.

○ Ayrıca yeni dataset üzerinde yapılan testlerde  $d$ 'nin küçük olması, kadar büyük olması da yüksek  $J_{\text{cv}}(\theta)$  hatası doğurur!

So we can assume your LA performing less well than you were hoping. ( $J_{\text{cv}}(\theta)$  or  $J_{\text{test}}(\theta)$  is high) Is it a bias problem or variance problem?

Bias (Underfit)

$J_{\text{train}}(\theta)$  will be high

$J_{\text{cv}}(\theta) \approx J_{\text{train}}(\theta)$

Clue of that the problem MAY BE underfitting

High Variance (Overfit)

$J_{\text{train}}(\theta)$  will be low

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

Clue of that problem MAY BE Overfitting



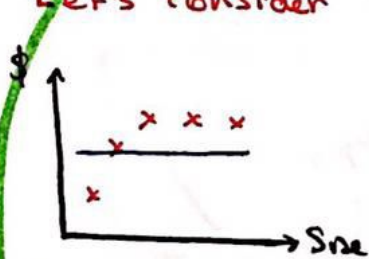
→ Regularization in overfitting: önleyebileceğini görmüşük. How does it affect the bias and variances of a learning algorithm?

→ Suppose we're fitting a high order polynomial like below but to prevent overfitting we're gonna use regularization. By regularization we are trying to keep values of parameters small.

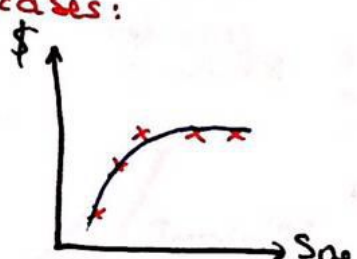
$$\text{Model: } h(x) = \theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3 + \theta_4 x^4$$

$$J(\theta) = \frac{1}{2m} \sum_{i=1}^m (h(x^{(i)}) - y^{(i)})^2 + \frac{\lambda}{2m} \sum_{j=1}^4 \theta_j^2$$

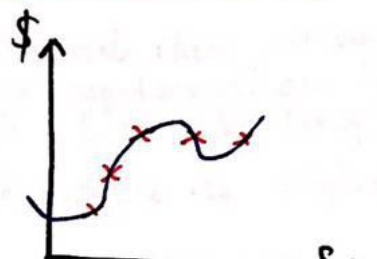
Let's consider 3 cases:



- Large  $\lambda$  (10000)
- High Bias (Underfit)
- $\theta_1 \approx 0, \theta_2 \approx 0, \dots$
- $h(x) \approx \theta_0$



- Intermediate  $\lambda$
- Just Right



- Small  $\lambda$
- High Variance (overfit)
- or  $\lambda = 0$

So how can we choose a good value for the  $\lambda$ ?

Assume this is our model and  $J(\theta)$  is defined as above:

• Deneyerek işledikimiz  $\lambda$  değerlerini alalım mesela:

- |                        |                           |       |                 |                                     |
|------------------------|---------------------------|-------|-----------------|-------------------------------------|
| ① Try $\lambda = 0$    | $\min_{\theta} J(\theta)$ | gives | $\theta^{(1)}$  | $\rightarrow J_{cv}(\theta^{(1)})$  |
| ② Try $\lambda = 0.01$ | $\min_{\theta} J(\theta)$ | gives | $\theta^{(2)}$  | $\rightarrow J_{cv}(\theta^{(2)})$  |
| ③ Try $\lambda = 0.02$ | "                         | "     | $\theta^{(3)}$  | $\rightarrow J_{cv}(\theta^{(3)})$  |
| ④ Try $\lambda = 0.04$ | "                         | "     | $\theta^{(4)}$  | $\rightarrow J_{cv}(\theta^{(4)})$  |
| ⋮                      | ⋮                         | ⋮     | ⋮               | ⋮                                   |
| ⑫ Try $\lambda = 10$   | "                         | "     | $\theta^{(12)}$ | $\rightarrow J_{cv}(\theta^{(12)})$ |

En düşük  $J_{cv}$ 'si olanı seçeriz burada  $\lambda = 0.08$  olsun yani 5. sıradakisi.

→ Son olarak  $J_{test}(\theta^{(5)})$ 'i bularak  $\lambda = 0.08$  için  $h(x)$  hipotezinin generalised emonunu bulabiliriz.

⊛ Aslında bu 2 sayfa önce yaptığımız işlemlerle aynı sayılır  $J(\theta)$  belki!

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

? Bunu burada neden kullanmadık anlamadım yavaşça  $J(\theta)$  kullandık.

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

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



Learning Curves are often a very useful thing to plot if either you wanted to sanity check that your algorithm is working correctly or if you want to improve the performance of the algorithm.

Learning Curves is a very often use tool to diagnose if a particular learning algorithm may be suffering from bias, or a variance problem or a bit of both.

Learning Curves dediğimiz  $J_{train}(\theta)$  ve  $J_{cv}(\theta)$ 'nin m'e yani # of training examples'a göre çizilmesi ile oluşur.

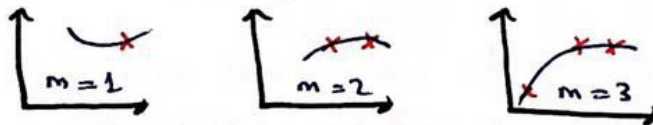
Bu grafiği çizmek için m'i artificially düşürmek yani normalde Training Set size 100 ise bunu 10'ıan 20'ıan vb. çizebiliriz. Training Error ve Validation Error bu azaltılmış yeni set için çizilir

$$J_{train}(\theta) = \frac{1}{2m} \sum_{i=1}^m (h(x^{(i)}) - y^{(i)})^2 \quad J_{cv}(\theta) = \frac{1}{2m_{cv}} \sum_{i=L}^{m_{cv}} (h(x^{(i)}) - y^{(i)})^2$$

Let's see what these plots may look like; **Plotting  $J_{train}(\theta)$**

1) Let's say I have only 1 training ex. and I am fitting a quadratic function  $h(x) = \theta_0 + \theta_1 x + \theta_2 x^2$  (Diyenleri sen de bu  $h(x)$ !)

Training Ex sayısı 1, 2 ve 3 iken quadratic hipotes iyi olur mu reg. yok ise tam olur mu varsa da çok küçük hata.



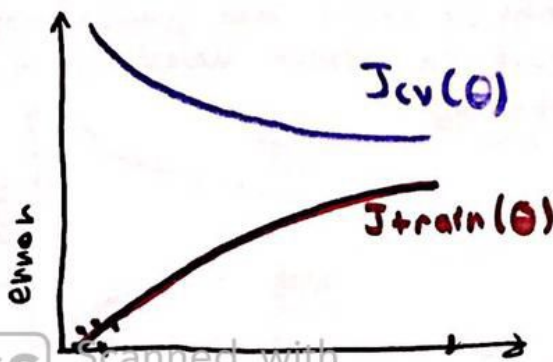
$J_{train}(\theta) = 0$  if no reg.  
 $J_{train} \approx 0$  if reg.

2)  $m \geq 4$  için quadratic function data set'e artık perfectly oturur. m arttıkça tüm data'lara ~~oturma~~ oturması zorlaşır!

Kısaca, training set boyutu arttıkça average training error  $J_{train}(\theta)$  da azalır



**Intuition:** when m's small it's easy to fit every single training ex-s perfectly. When m is getting larger it becomes harder



How about the  $J_{cv}(\theta)$ ?  $J_{cv}$  is my error on cross validation set that I haven't seen. ~~if~~ when I have very small training set I am not gonna generalize well.

when I get a larger training set I am starting to get a hypothesis maybe fit to the data better

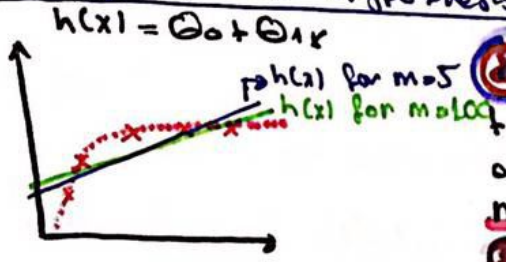
The more data you have, the better you do at generalizing to new examples.



Let's look at what the learning curves may look like if we have either **HIGH BIAS** or **HIGH VARIANCE** problems:

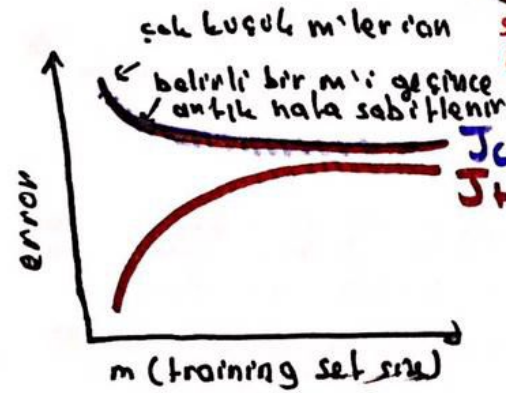
1) Suppose our Hypothesis have **HIGH BIAS**:

Plot  $J_{cv}(\theta)$



Let's think what would happen if we were to increase the training set size. So instead of five ex. let's say we have **100**. (Pembe notulen)

● Büyük bir değişim olmuyor zaten belli bir data sayısında sonra best hypothesis'e çok yaklaşıyor.  $m=1, 2$  vb. gibi küçük m'ler için hata yüksek olabilir.



How about  $J_{train}(\theta)$ : It will start small and for **HIGH BIAS** case it will end up close to the  $J_{cv}(\theta)$

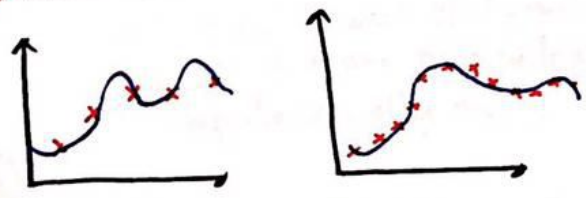
● Because you have so few parameters and so much data, at least when  $m$  is large the performance on the training set and the cross validation set will be very similar.

The problem of high bias is reflected in the fact that both the  $J_{cv}$  and  $J_{train}$  are high (relatively)

**Also Note That:** If a learning algorithm is suffering from high bias, getting more training data will not (by itself) help much. Navi ile yeşil line'ler aşağı yukarı aynı!

2) Suppose our Hypothesis have **HIGH VARIANCE**:

Assume  $h(x) = \theta_0 + \theta_1 x + \dots + \theta_{100} x^{100}$  and  $\lambda = \text{small value}$

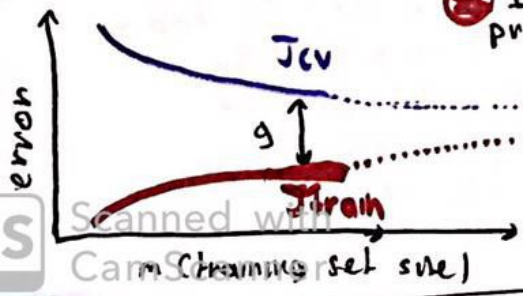


●  $m$  küçük iken hipotez gayet rahat oluşturulabilir zaten  $\lambda$  çok büyük o zaman  $J_{train}(\theta)$  küçük olur.

●  $m$  büyük iken hala hipotez güzel bir şekilde oluşturulabilir çünkü  $\lambda$  büyük ama  $m$  küçük olduğu kadar rahat olmaz. Sonuçta biraz daha büyük bir  $J_{train}(\theta)$ !

How about the  $J_{cv}(\theta)$ : In high variance setting, a hypothesis is overfitting and cross validation error will remain high even if we get moderate number of training examples.

Indicative diagnostic that we have a high variance problem is the large gap,  $g$ , between  $J_{cv}$  and  $J_{train}$



Addig more training data is likely to help for high variance problem.

Curveler her zaman böyle net olmaz ama yine de high bias mi high variance mi problem olduğunu anlayabiliriz.



1. Sayısal Örneğe Geri Dönelim: Suppose you have implemented regularized linear regression to predict housing prices. However, when you test your hypothesis in a new set of houses, you find that it makes unacceptably large errors in its prediction. What should you try next?

• Get more training examples → It fixes high variance

Eğer high bias problem varsa bir işe yaramaz.

• Try smaller sets of features → It fixes high variance

High bias can be useless

• Try getting additional features → Usually a solution for fixing a high bias problem.

Current hypothesis is too simple so we want to get additional features to make our hypothesis better able to fit the training set

• Try adding polynomial features → Similarly fix a high bias problem, High var. problem isn't solved almost.

• Try decreasing  $\lambda$  → Fixes high bias

• Try increasing  $\lambda$  → Fixes high variance

⊗ Let's take everything we have learned and relate it back to Neural Networks:

1) "Small" Neural Networks: Relatively few hidden units, maybe only one hidden layer. • This kind of networks have relatively few parameters • More prone to underfitting • ~~Small~~ advantage is that they are computationally cheaper.

2) "Large" Neural Networks: Either more hidden units or more hidden layers • More parameters • More prone to overfitting • Computationally more expensive.

⊗ Genelde using a larger NN by using regularization is more effective than using a small NN.

⊗ Fehi number of hidden layers kaç tane olmalı?

• Using a single hidden layer is reasonable default.

• But if you want to choose the # of hidden layers we can try → Find a training, cross-validation and test sets and try training NNs for different number of hidden layers and see which of those performs best on the cross validation set.



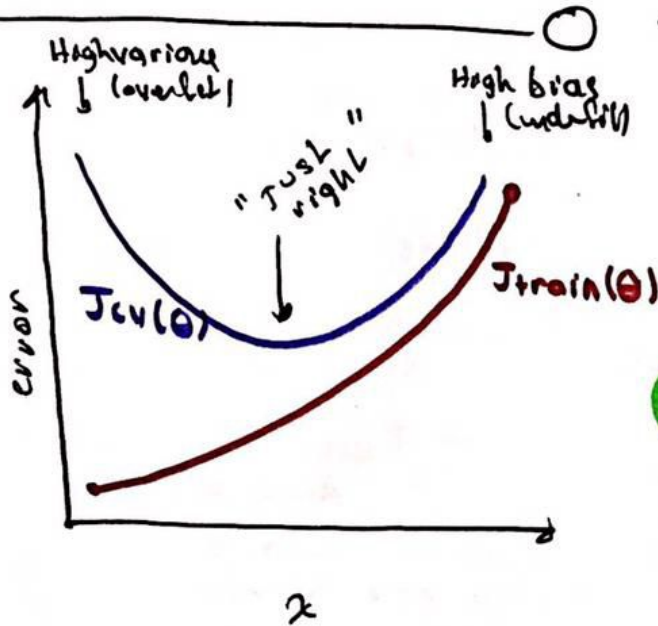
Let's look at how Cross Validation Error  $J_{cv}$ , and Training Error  $J_{train}$  vary as we vary the  $\lambda$ .

$$J(\theta) = \frac{1}{2m} \sum_{i=1}^m (h(x^{(i)}) - y^{(i)})^2 + \frac{\lambda}{2m} \sum_{j=1}^n \theta_j^2 \rightarrow \text{This is our original cost function}$$

Plot them against  $\lambda$

$$J_{train}(\theta) = \frac{1}{2m} \sum_{i=1}^m (h(x^{(i)}) - y^{(i)})^2 \rightarrow \text{Bunu yeni kurduğumuz grafik ve karşılaştırma için kullanacağız.}$$

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



- $\lambda$  küçük iken we are not using much regularization then larger risk of overfitting
- $\lambda$  büyük ise higher risk of underfit

✳ Burada sanırım şöyle kabul ediyoruz  $\lambda$  için  $\theta$  bulunuyor sonra bulunan  $\theta$  için  $J_{train}(\theta)$  ve  $J_{cv}(\theta)$  hesaplanır bu grafik o şekilde çizilir

- Gerçekle elde edilen curves daha complex olur biraz daha gürültülü olur. Bazi data setler için gerçekten böyle grafikler çizebiliriz
- Bu grafiği elde edip  $\lambda$ 'yı ona göre seçmek mantıklı olacaktır.

$\lambda$  is only used during training.  $\lambda$  is used to train an algorithm to obtain a good  $\theta$  which will balance the desire for low bias and low variance (no underfit and no overfit). Once training is complete, the training error and cross validation errors can be calculated to understand the quality of  $\theta$ . These error metrics are calculated using actual outputs and predicted outputs, and not the