# Diagnostics of Machine Learning algorithm:

**Summary:**

- **Getting more training examples:** Fixes high variance
- **Trying smaller sets of features:** Fixes high variance
- **Adding features:** Fixes high bias
- **Adding polynomial features:** Fixes high bias
- **Decreasing λ(Lambda):** Fixes high bias
- **Increasing λ(Lambda):** Fixes high variance.

☑ Diagnostics can give guidance as to what might be more fruitful things to try to
improve a learning algorithm.

**Correct**

☑ Diagnostics can be time-consuming to implement and try, but they can still be a very
good use of your time.

**Correct**

☑ A diagnostic can sometimes rule out certain courses of action (changes to your
learning algorithm) as being unlikely to improve its performance significantly.

**Correct**

## Evaluating a Hypothesis

Once we have done some trouble shooting for errors in our predictions by:

- Getting more training examples
- Trying smaller sets of features
- Trying additional features
- Trying polynomial features
- Increasing or decreasing λ

We can move on to evaluate our new hypothesis.

A hypothesis may have a low error for the training examples but still be inaccurate (because of overfitting). Thus, to evaluate a hypothesis, given a dataset of training examples, we can split up the data into two sets: a **training set** and a **test set**. Typically, the training set consists of 70 % of your data and the test set is the remaining 30 %.

The new procedure using these two sets is then:

1. Learn $\Theta$ and minimize $J_{train}(\Theta)$ using the training set
2. Compute the test set error $J_{test}(\Theta)$

### The test set error

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

2. For classification ~ Misclassification error (aka 0/1 misclassification error):

$$err(h_\Theta(x), y) = \begin{array}{ll} 1 & \text{if } h_\Theta(x) \geq 0.5 \text{ and } y = 0 \text{ or } h_\Theta(x) < 0.5 \text{ and } y = 1 \\ 0 & \text{otherwise} \end{array}$$

This gives us a binary 0 or 1 error result based on a misclassification. The average test error for the test set is:

$$\text{Test Error} = \frac{1}{m_{test}} \sum_{i=1}^{m_{test}} err(h_\Theta(x_{test}^{(i)}), y_{test}^{(i)})$$

This gives us the proportion of the test data that was misclassified.

## Model Selection and Train/Validation/Test Sets

Just because a learning algorithm fits a training set well, that does not mean it is a good hypothesis. It could over fit and as a result your predictions on the test set would be poor. The error of your hypothesis as measured on the data set with which you trained the parameters will be lower than the error on any other data set.

Given many models with different polynomial degrees, we can use a systematic approach to identify the 'best' function. In order to choose the model of your hypothesis, you can test each degree of polynomial and look at the error result.
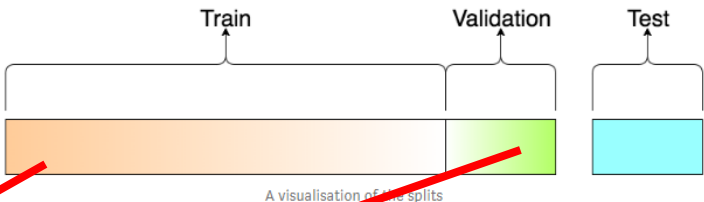
One way to break down our dataset into the three sets is:

- Training set: 60%
- Cross validation set: 20%
- Test set: 20%

We can now calculate three separate error values for the three different sets using the following method:

1. Optimize the parameters in $\Theta$ using the training set for each polynomial degree. ←
2. Find the polynomial degree d with the least error using the cross validation set.
3. Estimate the generalization error using the test set with $J_{test}(\Theta^{(d)})$, (d = theta from polynomial with lower error);

we get the error using the same training data.
ex: svmp.predict(traning_data)
and then compare with training sollution.

This way, the degree of the polynomial d has not been trained using the test set.

| Train | Validation | Test |

A visualisation of the splits

**Training Error:** We get the by calculating the classification error of a model on *the same data the model was trained on.*

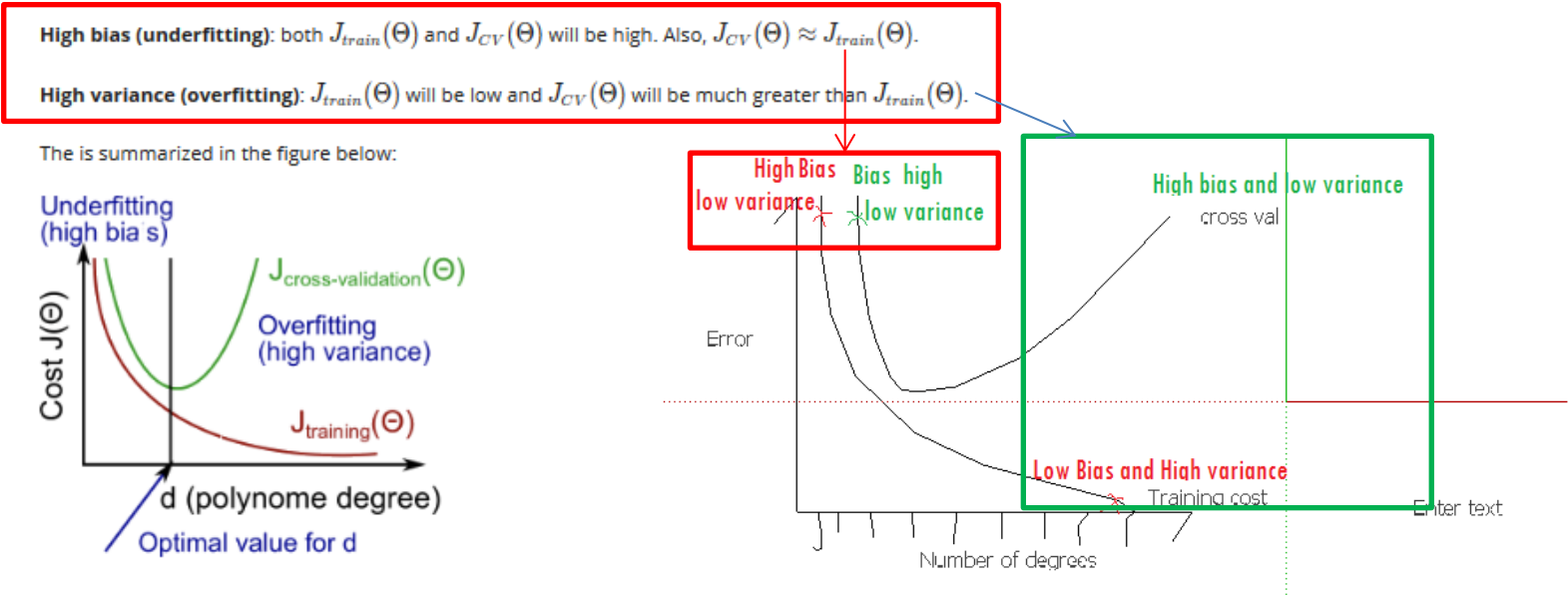Train_x, train_y,**Val_x,Val_y**= train_test_split(trainingdata)

# Diagnosing Bias vs. Variance

In this section we examine the relationship between the degree of the polynomial d and the underfitting or overfitting of our hypothesis.

- We need to distinguish whether **bias** or **variance** is the problem contributing to bad predictions.
- High bias is underfitting and high variance is overfitting. Ideally, we need to find a golden mean between these two.

<mark>The training error will tend to **decrease** as we increase the degree d of the polynomial.</mark>
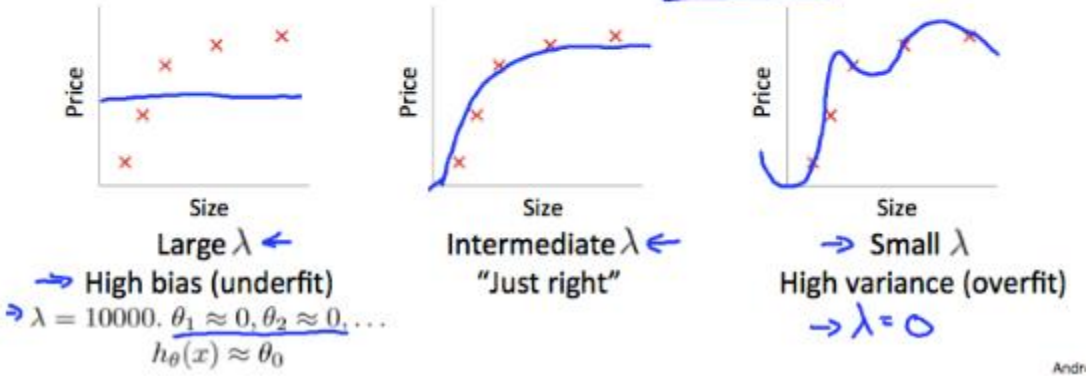
At the same time, the cross validation error will tend to **decrease** as we increase d up to a point, and then it will **increase** as d is increased, forming a convex curve.

**High bias (underfitting)**: both $J_{train}(\Theta)$ and $J_{CV}(\Theta)$ will be high. Also, $J_{CV}(\Theta) \approx J_{train}(\Theta)$.

**High variance (overfitting)**: $J_{train}(\Theta)$ will be low and $J_{CV}(\Theta)$ will be much greater than $J_{train}(\Theta)$.

The is summarized in the figure below:



## Linear regression with regularization

Model: $h_\theta(x) = \theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3 + \theta_4 x^4$ ←

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



| Large $\lambda$ ← | Intermediate $\lambda$ ← | Small $\lambda$ |
| High bias (underfit) | "Just right" | High variance (overfit) |

$\lambda = 10000.\ \theta_1 \approx 0, \theta_2 \approx 0, \ldots$
$h_\theta(x) \approx \theta_0$

$\rightarrow \lambda = 0$

Andn

If we have high lambda value then, Underfit

If we have small lambda value then , Overfit

Then we need to set a sweet inbetween value.

**Experiencing high variance:**

**Low training set size**: $J_{train}(\Theta)$ will be low and $J_{CV}(\Theta)$ will be high.

**Large training set size**: $J_{train}(\Theta)$ increases with training set size and $J_{CV}(\Theta)$ continues to decrease without leveling off. Also, $J_{train}(\Theta) < J_{CV}(\Theta)$ but the difference between them remains significant.

If a learning algorithm is suffering from **high variance**, getting more training data is likely to help.

More on Bias vs. Variance
Typical learning curve for high variance (at fixed model complexity):



# Learning Curves

Training an algorithm on a very few number of data points (such as 1, 2 or 3) will easily have 0 errors because we can always find a quadratic curve that touches exactly those number of points. Hence:

- As the training set gets larger, the error for a quadratic function increases.
- The error value will plateau out after a certain m, or training set size.

**Experiencing high bias:**

**Low training set size**: causes $J_{train}(\Theta)$ to be low and $J_{CV}(\Theta)$ to be high.

**Large training set size**: causes both $J_{train}(\Theta)$ and $J_{CV}(\Theta)$ to be high with $J_{train}(\Theta) \approx J_{CV}(\Theta)$.

If a learning algorithm is suffering from **high bias**, getting more training data will not **(by itself)** help much.

More on Bias vs. Variance
Typical learning curve for high bias (at fixed model complexity):



In which of the following circumstances is getting more training data likely to significantly help a learning algorithm's performance?

☐ Algorithm is suffering from high bias.

**Un-selected is correct**

☑ Algorithm is suffering from high variance.

**Correct**

☑ $J_{CV}(\theta)$ (cross validation error) is much larger than $J_{train}(\theta)$ (training error).
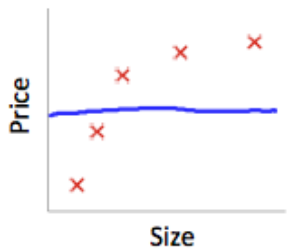
**Correct**

Lambda and its effects:



## Linear regression with regularization

**Model:** $h_\theta(x) = \theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3 + \theta_4 x^4$
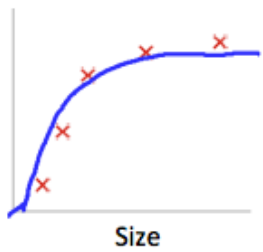
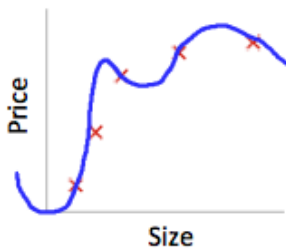$$J(\theta) = \frac{1}{2m}\sum_{i=1}^{m}(h_\theta(x^{(i)}) - y^{(i)})^2 + \frac{\lambda}{2m}\sum_{j=1}^{m}\theta_j^2$$

Large $\lambda$
High bias (underfit)
$\lambda = 10000$. $\theta_1 \approx 0, \theta_2 \approx 0, \dots$
$h_\theta(x) \approx \theta_0$

Intermediate $\lambda$
"Just right"

Small $\lambda$
High variance (overfit)
$\lambda = 0$

Andrew Ng

High Lambda, will penalize each theta and will make them close to zero

Underfit means, High Bias and reducing lambda will help us in Fixing Lambda

## Building Spam Classifier:

Given a data set of emails, we could construct a vector for each email. Each entry in this vector represents a word. The vector normally contains 10,000 to 50,000 entries gathered by finding the most frequently used words in our data set. If a word is to be found in the email, we would assign its respective entry a 1, else if it is not found, that entry would be a 0. Once we have all our x vectors ready, we train our algorithm and finally, we could use it to classify if an email is a spam or not.

**Building a spam classifier**

Supervised learning. $x$ = features of email. $y$ = spam (1) or not spam (0).
Features $x$: Choose 100 words indicative of spam/not spam.



So how could you spend your time to improve the accuracy of this classifier?

- Collect lots of data (for example "honeypot" project but doesn't always work)
- Develop sophisticated features (for example: using email header data in spam emails)
- Develop algorithms to process your input in different ways (recognizing misspellings in spam).

It is difficult to tell which of the options will be most helpful.

Q: Why it is recommended to perform error analysis using the cross validation data rather than the test data ?

Answer: If we develop new features by examining the test set, then we may end up choosing features that work well specifically for the test set, so Jtest(θ) is no longer a good estimate of how well we generalize to new examples.

## How to Evaluate Metrics?

## Using:

## Precision/Recall:



$$\text{Precision} = \frac{\text{True positives}}{\text{\# predicted as positive}} = \frac{\text{True positives}}{\text{True positives} + \text{False positives}}$$

$$\text{Recall} = \frac{\text{True positives}}{\text{\# actual positives}} = \frac{\text{True positives}}{\text{True positives} + \text{False negatives}}$$

Your algorithm's performance on the test set is given to the right. What is the algorithm's precision?

|  | Actual class 1 | Actual class 0 |
|---|---|---|
| Predicted class 1 | 80 | 20 |
| Predicted class 0 | 80 | 820 |

Precision: 0.8

Recall : 0.8