

Evaluating the performance of a model and hypertuning:

Methods to improve the performance of a model on a new set of data input.

- ➔ Get more training examples
  - However this may not help very much in many cases and leads to a lot of time wasted
- ➔ Try a smaller set of features
- ➔ Try getting additional features
- ➔ Try adding polynomial features
- ➔ Increase  $\lambda$
- ➔ Decrease  $\lambda$

Hence how should we decide which of these diagnoses we should consider pursuing first?

This is where Machine Learning Diagnostics come into play, as tests to gain insight into what is/is not working with a learnt model and gain direction on how to improve it's performance.

Evaluating the Hypothesis:

Overfitting the dataset?

Soln: Split the raw dataset into a train-test split

70% train, 30%

Train-Test procedure for Linear Regression/ any other models:

- 1) Learn Params from Train split of the data.
- 2) Compute test set error on cost function.

$$\text{Eg: for linear regression: } J_{test}(\theta) = \frac{1}{2m} \sum_{i=1}^{m_{test}} (h_{\theta}(x_{test}^{(i)}) - y_{test}^{(i)})^2$$

For Classification Problems:

Another option is misclassification error:

$$err(h_{\theta}(x^{(i)}), y^{(i)}) = \begin{cases} 1 & \text{if } h_{\theta}(x^{(i)}) \geq \text{threshold}; y = 0 \text{ or } h_{\theta}(x^{(i)}) < \text{threshold}; y = 1 \\ 0 & \text{otherwise} \end{cases}$$

$$\text{Test error} = \frac{1}{m_{test}} \sum_{i=1}^{m_{test}} err(h_{\theta}(x^{(i)}), y^{(i)})$$

A percentage of misclassified images

Model Selection:

Once params  $\theta_0, \dots, \theta_n$  are fit to a training set the error of params as measured on that set of data is likely to be lower than the actual generalization error, as the params are optimized to that particular set of data.

Hence we further split the data into:

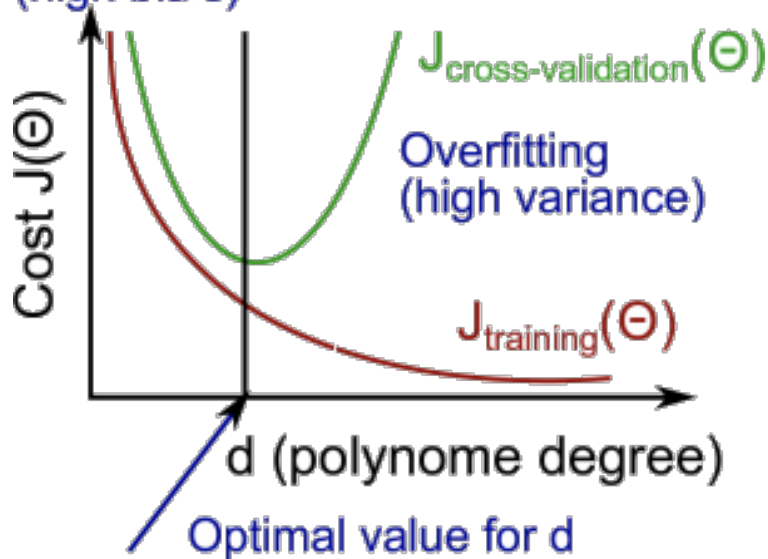
Train, Test, Cross Validation(CV) sets at a proportion of 60%,20%,20% after random permutation of the dataset.

Then we iterate through degrees of polynomials on the CV set in order to select a best fit polynomial for the parameters on the CV set. The data would then have a bias on the minimization of the CV set. However, not on the test set which measures the ability of the model to predict on new data.

This helps resolve the performance bias we may face on the test dataset, and hence obtain a more generalizable metric of error.

Diagnosing Bias vs Variance:

Underfitting  
(high bias)



Bias (Underfit):

$$J_{\text{train}}(\theta) \text{ is high}$$

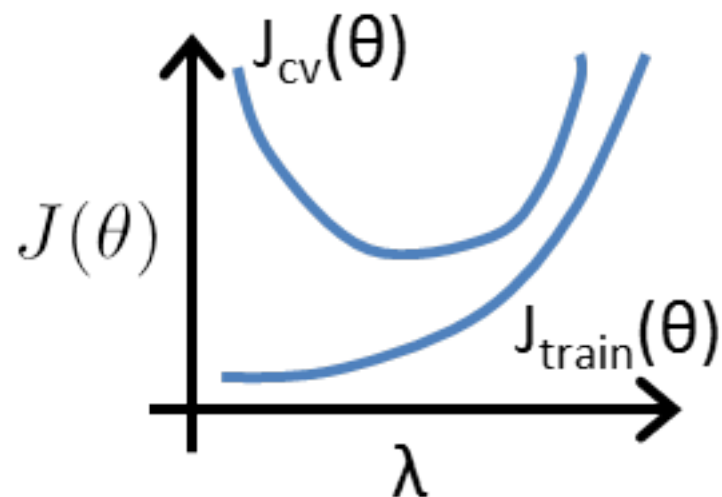
$$J_{\text{cv}}(\theta) \approx J_{\text{train}}(\theta)$$

Variance(Overfit):

$$J_{\text{train}}(\theta) \text{ is low}$$

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

Regularization:



If  $\lambda$  too small  $J_{\text{cv}}(\theta)$  will be large and  $J_{\text{train}}(\theta)$  is small due to overfitting

If  $\lambda$  too large  $J_{\text{cv}}(\theta)$  will be large and  $J_{\text{train}}(\theta)$  is large due to underfitting

Hence how do we solve this?

Workflow:

1. Create a list of lambdas (i.e.  $\lambda \in \{0, 0.01, 0.02, 0.04, 0.08, 0.16, 0.32, 0.64, 1.28, 2.56, 5.12, 10.24\}$ );
2. Create a set of models with different degrees or any other variants.
3. Iterate through the  $\lambda$ s and for each  $\lambda$  go through all the models to learn some  $\Theta$ .
4. Compute the cross validation error using the learned  $\Theta$  (computed with  $\lambda$ ) on the  $J_{cv}(\Theta)$  **without** regularization or  $\lambda = 0$ .
5. Select the best combo that produces the lowest error on the cross validation set.
6. Using the best combo  $\Theta$  and  $\lambda$ , apply it on  $J_{test}(\Theta)$  to see if it has a good generalization of the problem.

Learning Curves:

High Bias:

[More on Bias vs. Variance](#)

Typical learning curve for high bias (at fixed model complexity):



**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  $J_{train}(\Theta) \approx J_{cv}(\Theta)$ .

Getting more training data **will not** (by itself) **help much**.

Possibly requires feature analysis, or polynomial regressions

**Experiencing high variance:**

[More on Bias vs. Variance](#)

Typical learning curve for high variance (at fixed model complexity):



**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.

**More training data is likely to help.**

Applied to Neural Networks:

“Small” Neural Networks (few parameter, more prone to underfitting)

“Large” Neural Networks (more prone to overfitting)

- Use  $\lambda$  to address overfitting

CrossValidation is a good method to select the number of hidden layers in a network.