

CS1138

Machine Learning

Lecture : Evaluation Metrics, Model Selection, Data Splits
and Bias Variance Tradeoff

(Slide Credits: Sebastian Raschka, Andrew Ng)

Arpan Gupta

Overview

- Confusion Matrix
- Precision, Recall, F1-score
- Balanced Accuracy
- ROC Curve
- Model Selection / Data Splits
- Bias Variance Tradeoff

Find the accuracy

- Suppose there are 100 patients, 50 of them actually have a disease, while remaining 50 are healthy. An ML model (which pretends to classify the patients based on disease), just takes the input features of patients and outputs “healthy” for all the cases. (It actually doesn’t use the features, just outputs the same thing for all the patients). What is the accuracy of the “model”?
- Answer: **50%**

Find the accuracy

- Suppose there are 100 patients, 10 of them actually have a disease (say cancer), while remaining 90 are healthy. What is the accuracy when the same “model” is used?

- Answer: **90%**

Class Imbalance problem arises. Therefore, our accuracy metric does not show the correct evaluation.

2x2 Confusion Matrix

		Predicted class	
		P	N
Actual class	P	True positives (TP)	False negatives (FN)
	N	False positives (FP)	True negatives (TN)

Malignant Tumour:
Positive Class

Benign Tumour:
Negative Class

$$ERR = \frac{FP + FN}{FP + FN + TP + TN} = 1 - ACC$$

$$ACC = \frac{TP + TN}{FP + FN + TP + TN} = 1 - ERR$$

TP, TN, FP, FN

- True Positive (TP): #Samples of positive class that were correctly predicted. Malignant tumours correctly predicted.
- True Negative (TN): Benign tumours correctly predicted.
- False Positive (FP): Predicted malignant, but were actually benign.
- False Negative (FN): Predicted benign, but were actually malignant.

TPR, FPR, TNR, FNR

- TPR: ratio of correct +ve predictions to the total positives in the dataset.
- FPR: ratio of false positives to the total number of negatives in the dataset.
- TNR: ratio of correct -ve predictions to the total negatives in the dataset.
- FNR: ratio of false negatives to the total number of positives in the dataset.

False Positive Rate and False Negative Rate

$$\text{TPR}^* = \frac{\text{TP}}{\text{P}} = \frac{\text{TP}}{\text{TP} + \text{FN}} = 1 - \text{FNR}$$

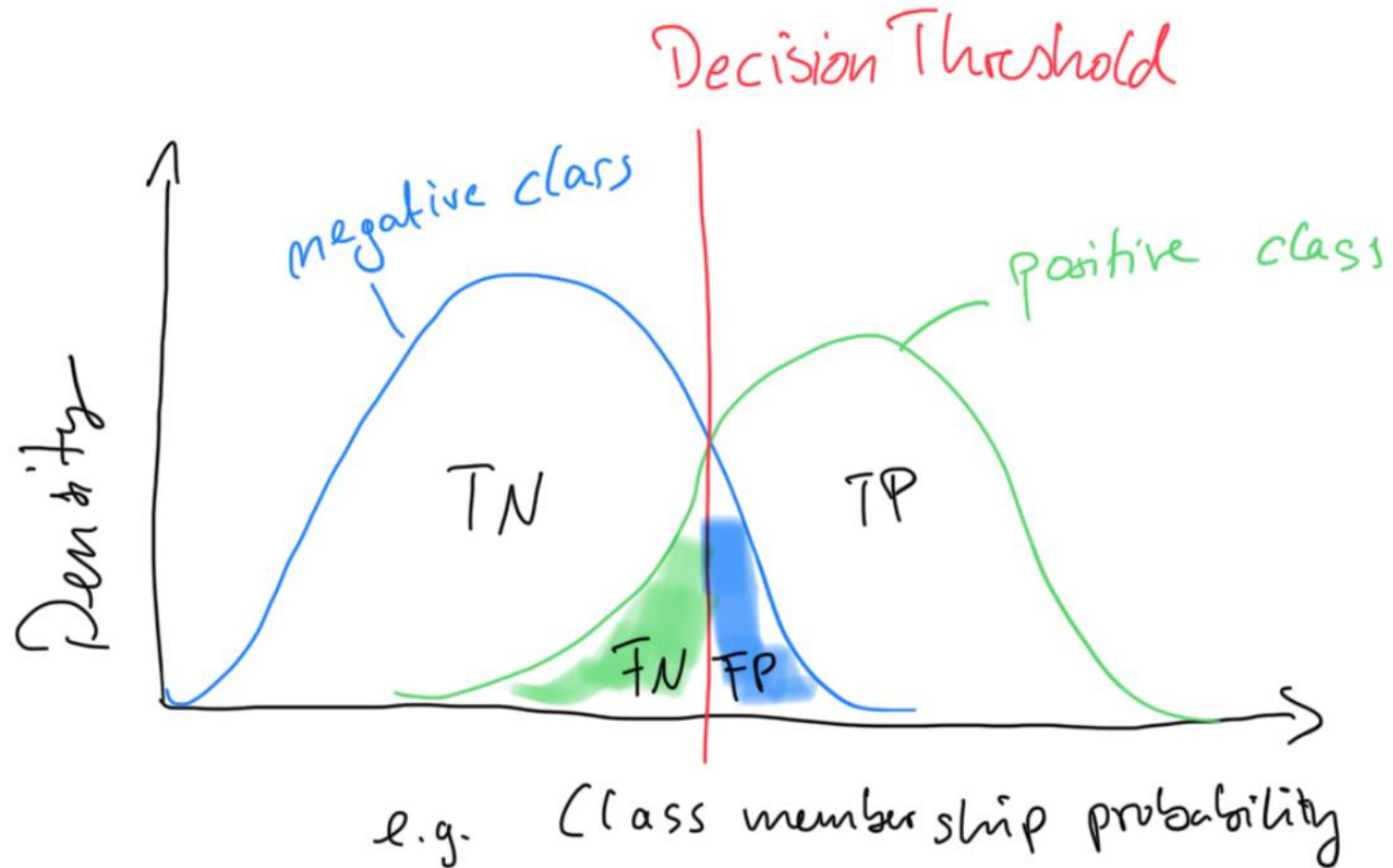
* Relevant later for ROC

$$\text{FPR}^* = \frac{\text{FP}}{\text{N}} = \frac{\text{FP}}{\text{FP} + \text{TN}} = 1 - \text{TNR}$$

$$\text{FNR} = \frac{\text{FN}}{\text{P}} = \frac{\text{FN}}{\text{FN} + \text{TP}} = 1 - \text{TPR}$$

$$\text{TNR} = \frac{\text{TN}}{\text{N}} = \frac{\text{TN}}{\text{TN} + \text{FP}} = 1 - \text{FPR}$$

FPR Vs FNR



Confusion Matrix for multi-class setting

		Predicted Labels		
		Class 0	Class 1	Class 2
True Labels	Class 0	$T(0,0)$		
	Class 1		$T(1,1)$	
	Class 2			$T(2,2)$

Confusions matrices are traditionally for binary class problems but we can be readily generalized it to multi-class settings

Precision and Recall

$$Precision = \frac{TP}{TP + FP}$$

→ #Correct +ve predictions

→ #Total +ve predictions

$$Recall = \frac{TP}{TP + FN}$$

→ #Correct +ve predictions

→ #Total +ve Ground Truth values

Recall is another name for “True Positive Rate (TPR)” or “sensitivity”.

F1-Score

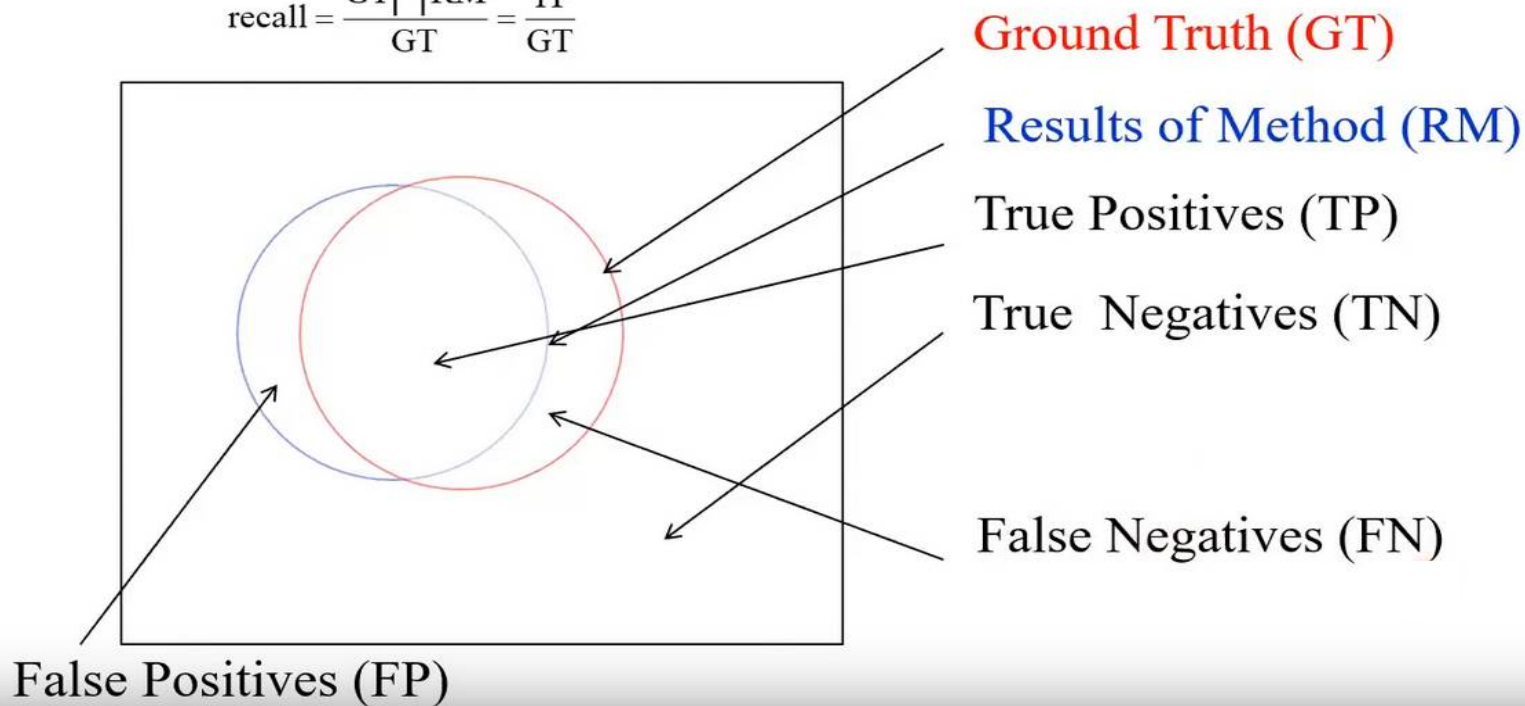
- It is the harmonic mean of Precision and Recall.

$$F1-Score = \frac{2 \times Precision \times Recall}{Precision + Recall}$$

Evaluation Metrics

$$\text{precision} = \frac{GT \cap RM}{RM} = \frac{TP}{RM}$$

$$\text{recall} = \frac{GT \cap RM}{GT} = \frac{TP}{GT}$$



Find the precision recall and F1-score

- Suppose we want to classify 100 patients into having cancer or not. The expert doctors (after taking into account their features/lab reports etc.) have said that 5 have cancer, and rest of them do not have cancer. We build an ML model, which generates the outputs and we find the TP, TN, FP and FN values as follows. What is the precision and recall values of the same “model”?

TP=3, TN = 92, FP=3, FN = 2

- Answer:

Precision: $3 / (3+3) = 0.5$ (or 50%)

Recall: $3 / (3 + 2) = 0.6$ (or 60%)

F1-Score = $2 * 0.5 * 0.6 / (1.1) = 0.5454$

Sensitivity and Specificity

- ***Sensitivity*(SEN)** measures the recovery rate of the Positives and complimentary, ***Specificity*(SPC)** measure the recovery rate of the Negatives.
- $SEN = TPR = Recall = TP/P = TP/(TP+FN)$
- $SPC = TNR = TN/N = TN/(TN+FP)$

Matthew's Correlation Coefficient

- Matthews correlation coefficient (MCC) was first formulated by Brian W. Matthews [1] in 1975 to assess the performance of protein secondary structure predictions
- The MCC can be understood as a specific case of a linear correlation coefficient (Pearson r) for a binary classification setting
- Considered as especially useful in unbalanced class settings
- The previous metrics take values in the range between 0 (worst) and 1 (best)
- The MCC is bounded between the range 1 (perfect correlation between ground truth and predicted outcome) and -1 (inverse or negative correlation) — a value of 0 denotes a random prediction.

$$MCC = \frac{TP \cdot TN - FP \cdot FN}{\sqrt{(TP + FP)(TP + FN)(TN + FP)(TN + FN)}}$$

[1] Brian W Matthews. Comparison of the predicted and observed secondary structure of T4 phage lysozyme. *Biochimica et Biophysica Acta (BBA) – Protein Structure*, 402(2), 1975

Implementation : Confusion matrix

```
from sklearn.preprocessing import StandardScaler
from sklearn.neighbors import KNeighborsClassifier
from sklearn.pipeline import make_pipeline
from mlxtend.evaluate import confusion_matrix

pipe_knn = make_pipeline(StandardScaler(),
                          KNeighborsClassifier(n_neighbors=5))

pipe_knn.fit(X_train, y_train)

y_pred = pipe_knn.predict(X_test)

confmat = confusion_matrix(y_test, y_pred)

print(confmat)

[[71  1]
 [ 3 39]]
```

sklearn.metrics.confusion_matrix

```
>>> from sklearn.metrics import confusion_matrix
>>> y_true = [2, 0, 2, 2, 0, 1]
>>> y_pred = [0, 0, 2, 2, 0, 2]
>>> confusion_matrix(y_true, y_pred)
array([[2, 0, 0],
       [0, 0, 1],
       [1, 0, 2]])
```

```
>>> tn, fp, fn, tp = confusion_matrix([0, 1, 0, 1], [1, 1, 1, 0]).ravel()
>>> (tn, fp, fn, tp)
(0, 2, 1, 1)
```

Precision, Recall, F1, MCC

```
from sklearn.metrics import accuracy_score, precision_score, \
    recall_score, f1_score, matthews_corrcoef
```

```
print('Accuracy: %.3f' % accuracy_score(y_true=y_test, y_pred=y_pred))
print('Precision: %.3f' % precision_score(y_true=y_test, y_pred=y_pred))
print('Recall: %.3f' % recall_score(y_true=y_test, y_pred=y_pred))
print('F1: %.3f' % f1_score(y_true=y_test, y_pred=y_pred))
print('MCC: %.3f' % matthews_corrcoef(y_true=y_test, y_pred=y_pred))
```

Accuracy: 0.965
Precision: 0.975
Recall: 0.929
F1: 0.951
MCC: 0.925

https://github.com/rasbt/stat451-machine-learning-fs20/blob/master/L12/code/12_2_pre-recall-f1.ipynb

Balanced Accuracy / Average Per-Class (APC) Accuracy

		Predicted Labels		
		Class 0	Class 1	Class 2
True Labels	Class 0	T(0,0)		
	Class 1		T(1,1)	
	Class 2			T(2,2)

$$ACC = \frac{T}{n}$$

Correct predictions

Total samples in
dataset

Balanced Accuracy / Average Per-Class (APC) Accuracy

		Predicted Labels		
		Class 0	Class 1	Class 2
True Labels	Class 0	3	0	0
	Class 1	7	50	12
	Class 2	0	0	18

$$ACC = \frac{3 + 50 + 18}{90} \approx 0.79$$

$$APC\ ACC = \frac{83/90 + 71/90 + 78/90}{3} \approx 0.86$$

Balanced Accuracy / Average Per-Class (APC) Accuracy

True Labels	Predicted Labels	
	Class 0	Neg Class
Class 0	3	0
Neg Class	7	80

True Labels	Predicted Labels	
	Class 1	Neg Class
Class 1	50	19
Neg Class	0	21

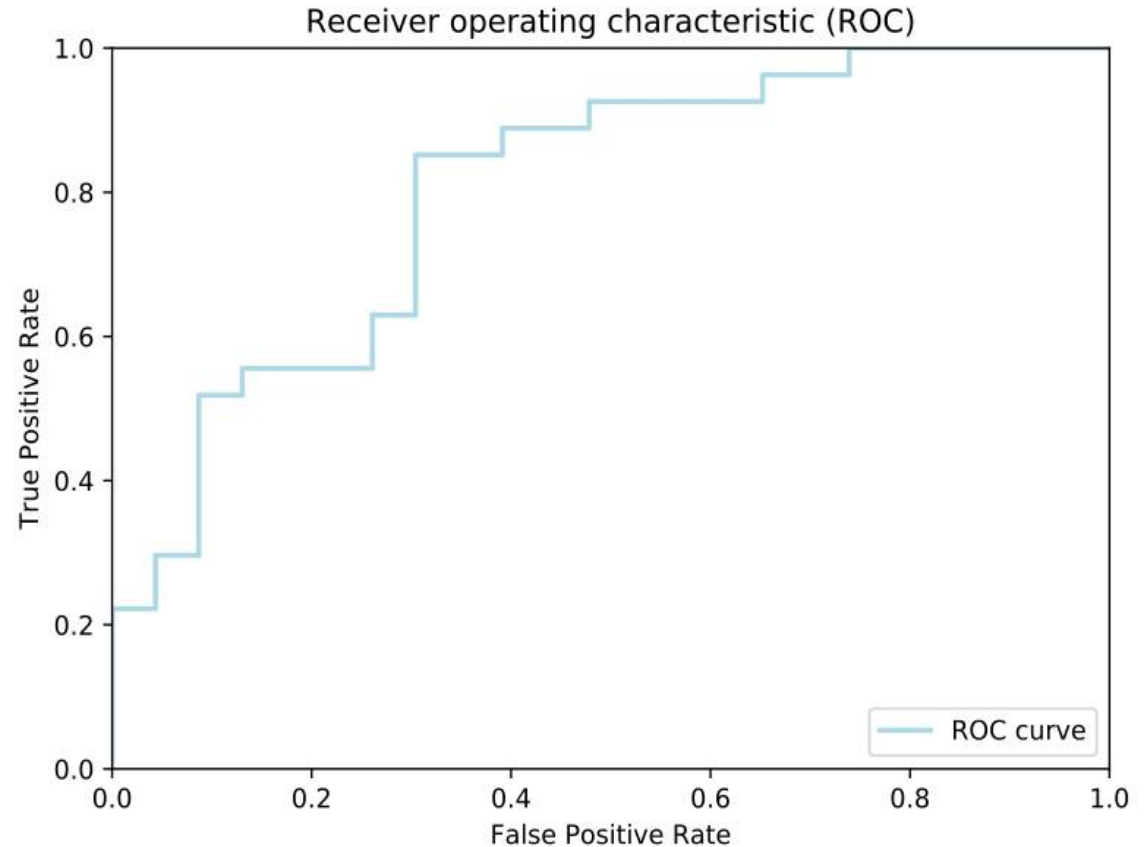
True Labels	Predicted Labels	
	Class 2	Neg Class
Class 2	18	0
Neg Class	12	60

True Labels	Predicted Labels		
	Class 0	Class 1	Class 2
Class 0	3	0	0
Class 1	7	50	12
Class 2	0	0	18

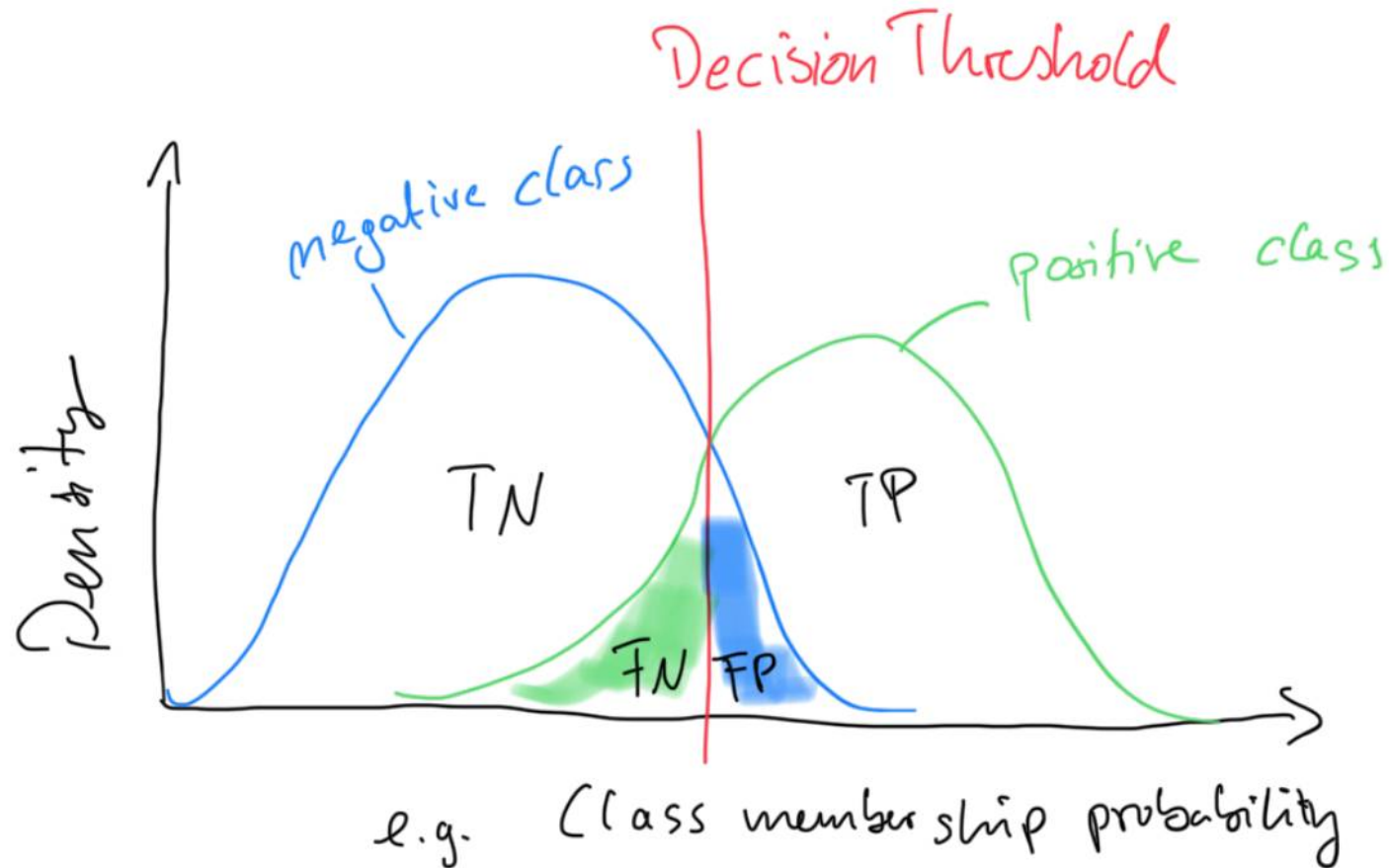
$$APC\ ACC = \frac{83/90 + 71/90 + 78/90}{3} \approx 0.86$$

Receiver Operating Characteristic Curve (ROC Curve)

- Trade-off between True Positive Rate(TPR) and False Positive Rate(FPR)
- ROC can be plotted by changing the prediction threshold

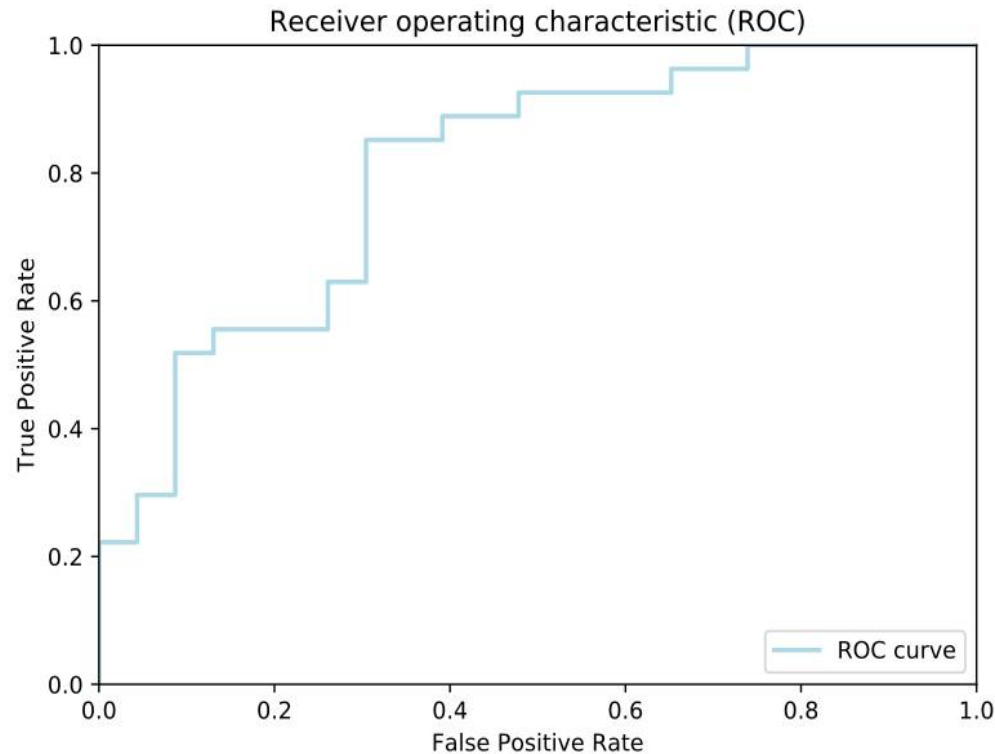


RECAP: False Positive Rate and False Negative Rate

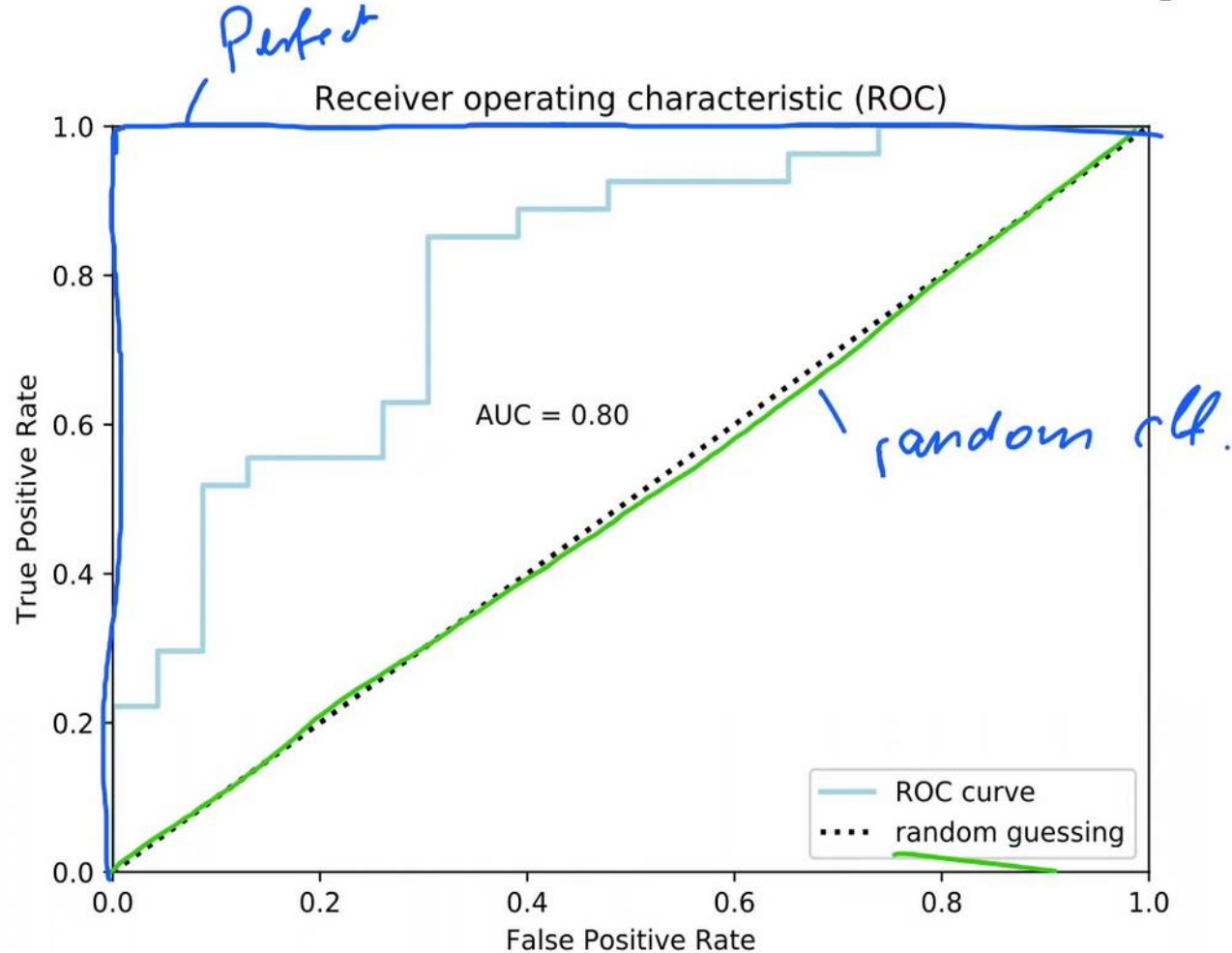


Receiver Operating Characteristic curve (ROC curve)

- 1.0 = Perfect Prediction
- 0.5 = Random Prediction



ROC Area Under the Curve (AUC)



ROC Curve Demo

https://developers.google.com/machine-learning/crash-course/classification/roc-and-auc#auc_and_roc_for_choosing_model_and_threshold

Model Selection and Data Splits

Debugging a learning algorithm:

Suppose you have implemented regularized linear regression to predict housing prices.

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

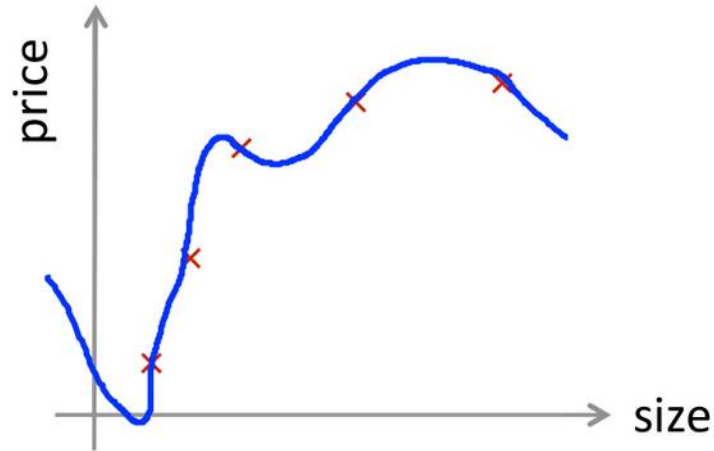
However, when you test your hypothesis on a new set of houses, you find that it makes unacceptably large errors in its predictions. What should you try next?

- Getting more training examples.
- Try smaller set of features.
- Try getting additional features.
- Try getting polynomial features (x_1^2 , x_2^2 , x_1x_2 , etc.)
- Try decreasing λ
- Try increasing λ

Evaluating a hypothesis

- Suppose we train a model and we get a very low training error. But it fails to generalize on the new unseen examples.
- One reason can be that the hypothesis is overfitting
- How do we know that it is overfitting?

Evaluating your hypothesis



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

Fails to generalize to new examples not in training set.

x_1 = size of house

x_2 = no. of bedrooms

x_3 = no. of floors

x_4 = age of house

x_5 = average income in neighborhood

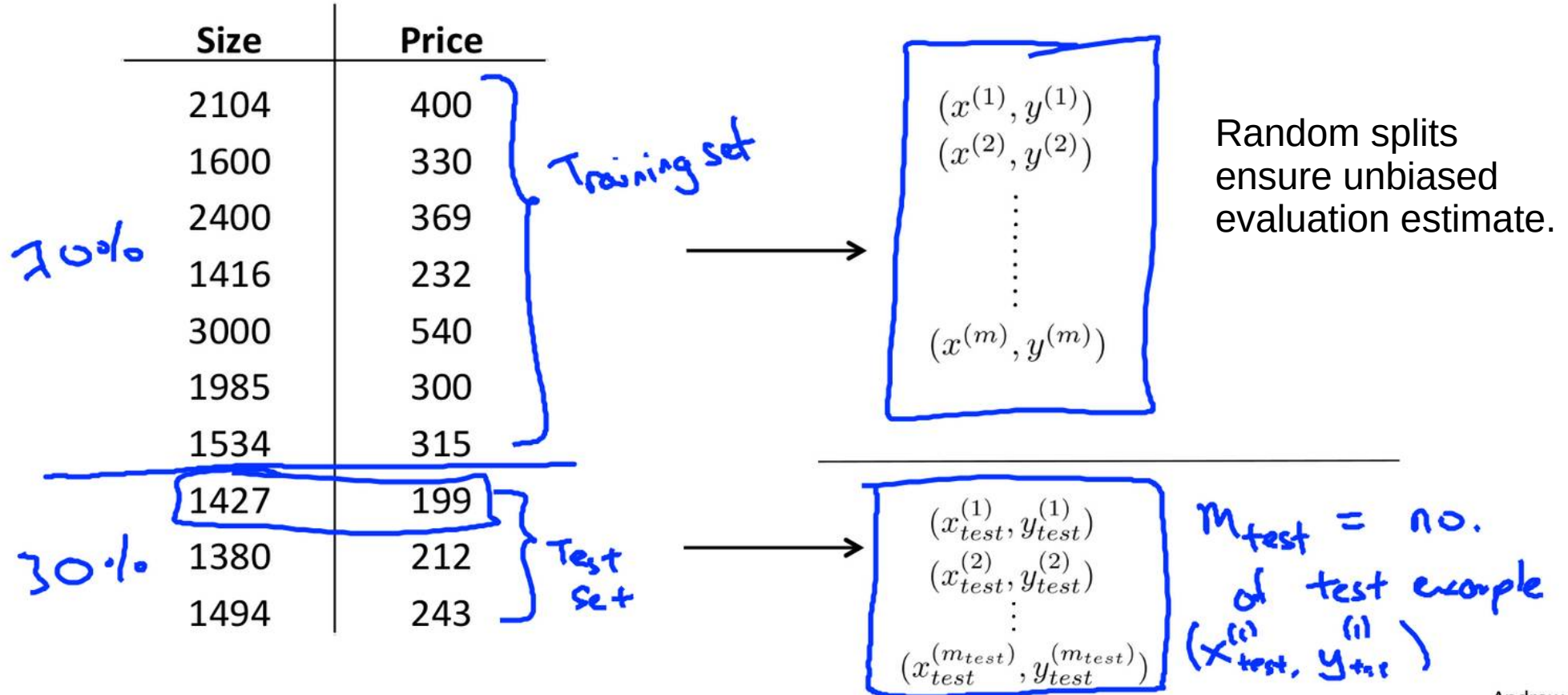
x_6 = kitchen size

\vdots

x_{100}

Evaluating your hypothesis

Dataset:



Evaluation on test set

- Linear Regression:

Compute test set error:

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

- Logistic Regression:

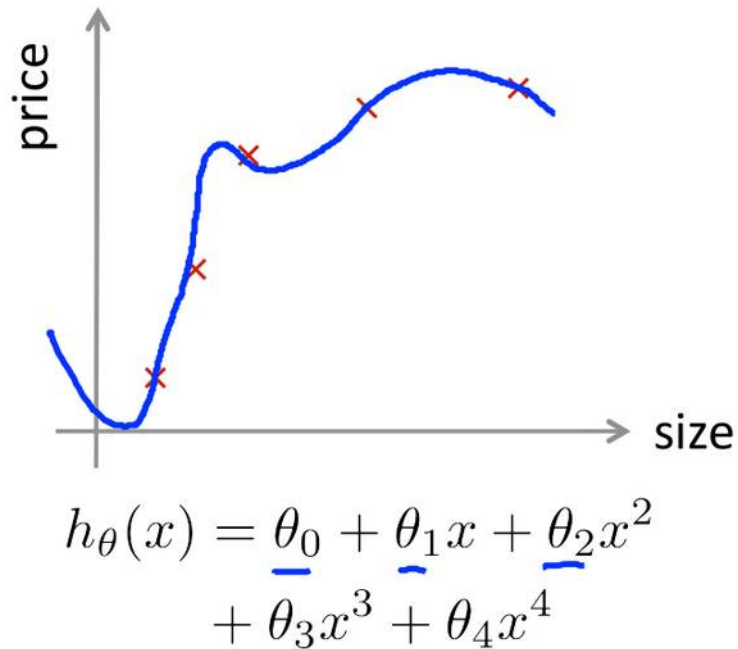
Compute test set error:

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

Model selection and training/validation and test sets

- If we want to decide what degree of polynomial features to use for the hypothesis while using logistic regression, or what should be the value of the regularization parameter λ .
- These are known as model selection problems.

Overfitting example



Once parameters $\theta_0, \theta_1, \dots, \theta_4$ were fit to some set of data (training set), the error of the parameters as measured on that data (the training error $J(\theta)$) is likely to be lower than the actual generalization error.

Model selection

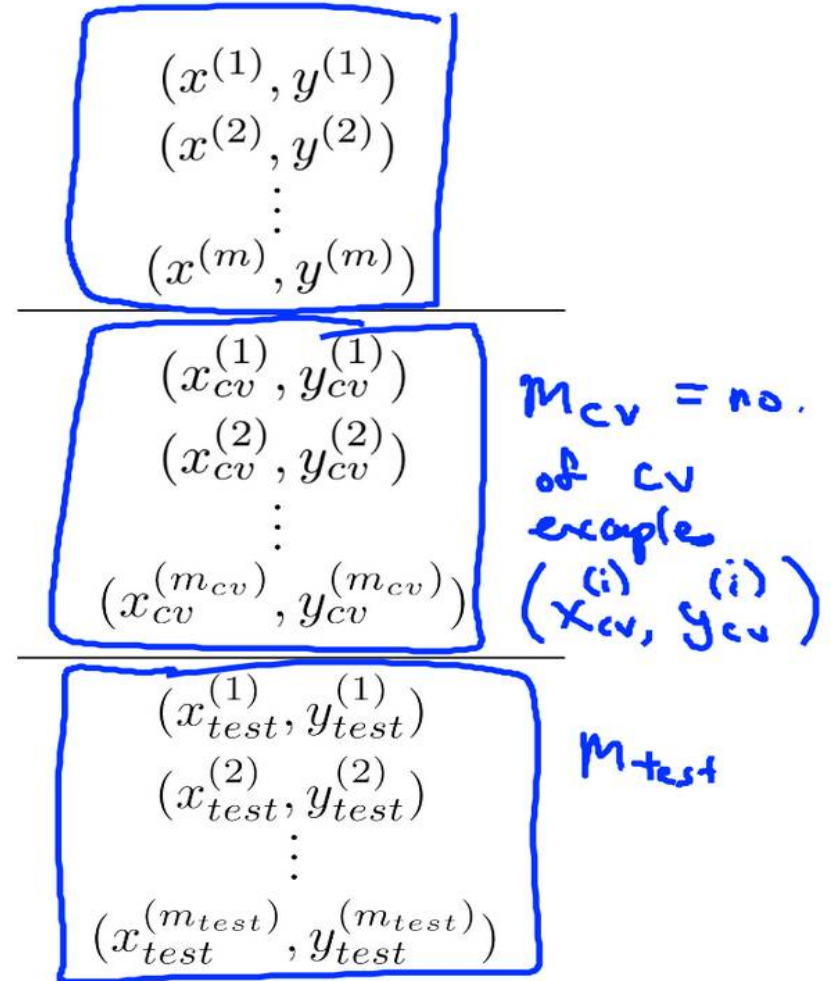
1. $h_{\theta}(x) = \theta_0 + \theta_1 x$
2. $h_{\theta}(x) = \theta_0 + \theta_1 x + \theta_2 x^2$
3. $h_{\theta}(x) = \theta_0 + \theta_1 x + \dots + \theta_3 x^3$
- \vdots
10. $h_{\theta}(x) = \theta_0 + \theta_1 x + \dots + \theta_{10} x^{10}$

- Suppose, we train 10 models with different degrees of polynomial (**d**), and evaluate on the test set. Finally, report the model with the lowest error.
- This will not be a fair estimate of the model. It is likely to be an optimistic estimate of the generalization error i.e., our extra parameter of **d** is fit to the test set.
- Therefore, we need a Cross Validation (or simply validation) set for model selection.

Evaluating your hypothesis

Dataset:

Size	Price	
2104	400	60% Training set
1600	330	
2400	369	
1416	232	
3000	540	
1985	300	
1534	315	20% Cross validation set (CV)
1427	199	
1380	212	20% Test set
1494	243	



Train/validation/test error

Training error:

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

Cross Validation error:

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

Test error:

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

Train / Val / Test splits: Points to note

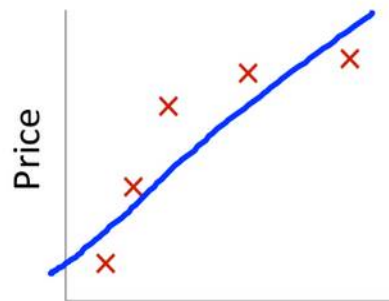
- The validation set helps in choosing model hyperparameters (such as degree of polynomial, or regularization parameter λ), while the training set will determine the parameters (Θ).
- To have an **unbiased estimate** of the performance, use the test set once on the trained model. Do not make any decisions based on the accuracy that you get on the test set.
- E.g., You may test the new models on a weekly basis in production environment, but if you revert to an older model (considering that in the current week the performance was worse), then the estimate will not remain unbiased.

Train / Val / Test splits: Points to note

- Generally the train val test splits are in the ratio 60%:20%:20. But it is quite arbitrary and depends on the dataset size and error improvement that we want to achieve.
- If the dataset is small, so we may want more data to be put in the training set. In this case, we may use **k-fold cross validation**. This divides the training data into k equal parts (k ranging from 1 to 10), train the model on k-1 folds and test on the remaining 1 fold. This is done for all the folds and validation accuracy is obtained by averaging all the k accuracy values.
- **Leave-One-Out CV** is an extreme case of k-fold CV, where we train on m-1 samples and test on one sample (k = m = training set size).

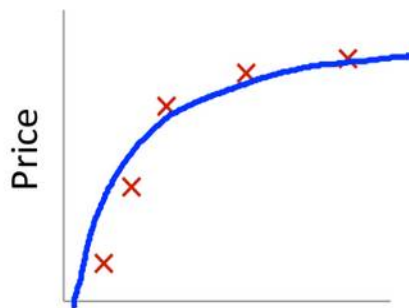
Bias Variance Tradeoff

Bias/variance



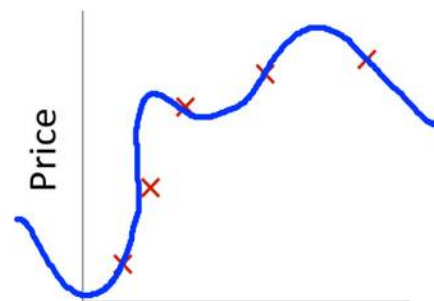
Size
 $\theta_0 + \theta_1 x$

High bias
(underfit)
 $d=1$



Size
 $\theta_0 + \theta_1 x + \theta_2 x^2$

“Just right”
 $d=2$



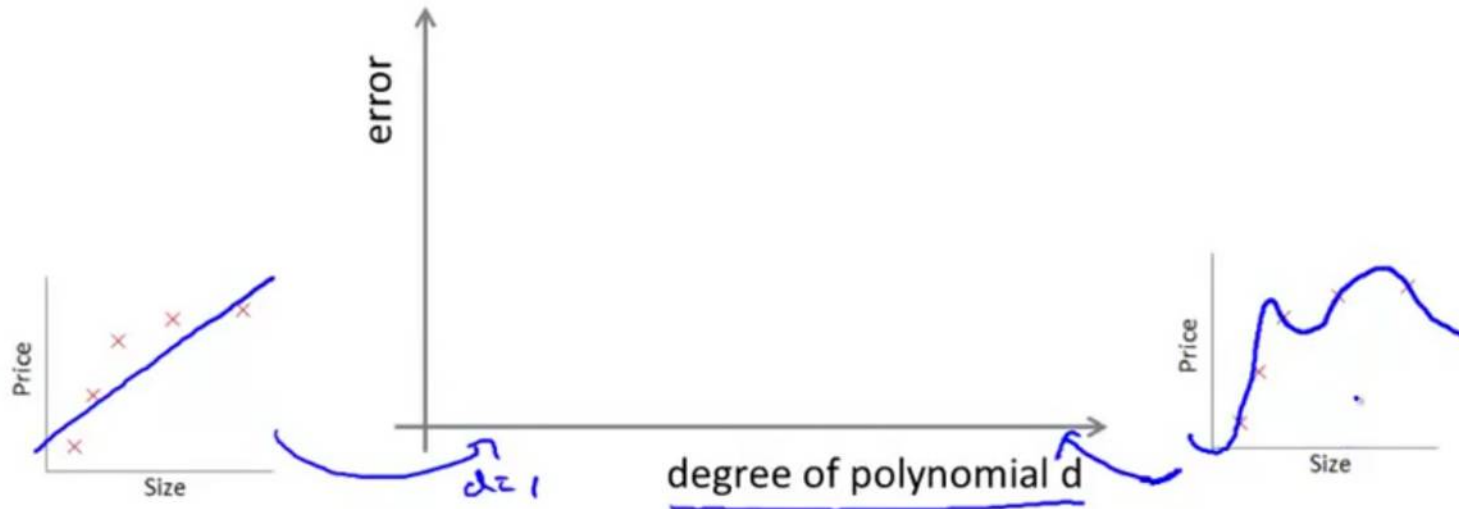
Size
 $\theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3 + \theta_4 x^4$

High variance
(overfit)
 $d=4$

Bias/variance

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

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

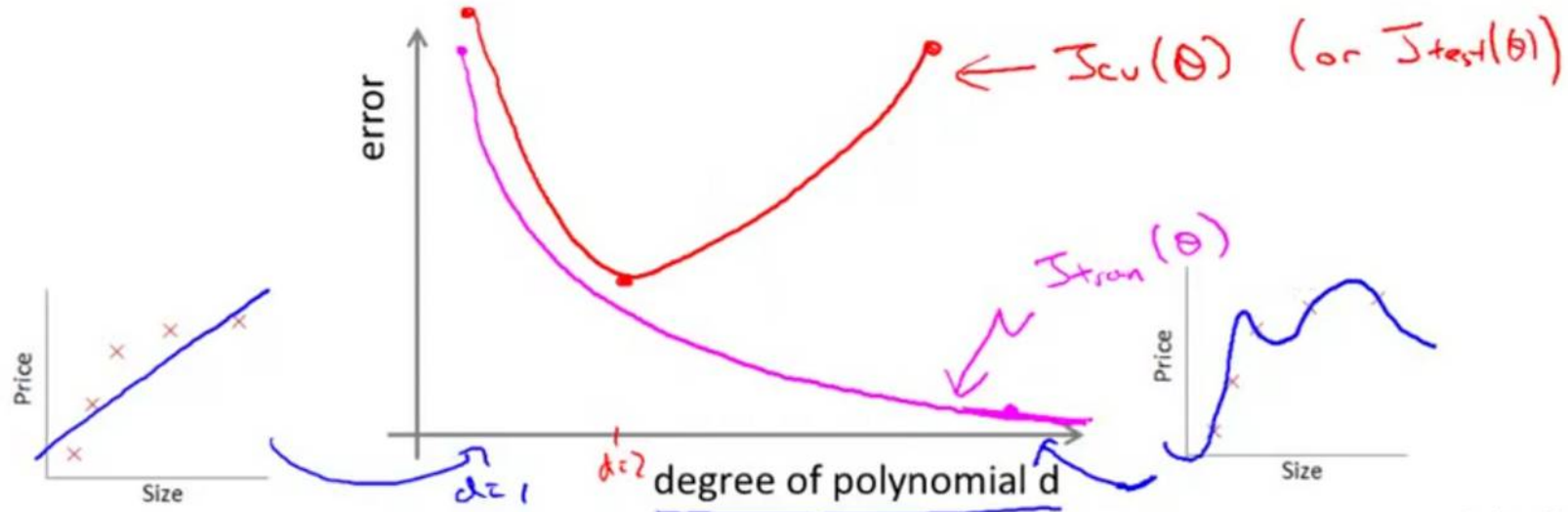


Andrew Ng

Bias/variance

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

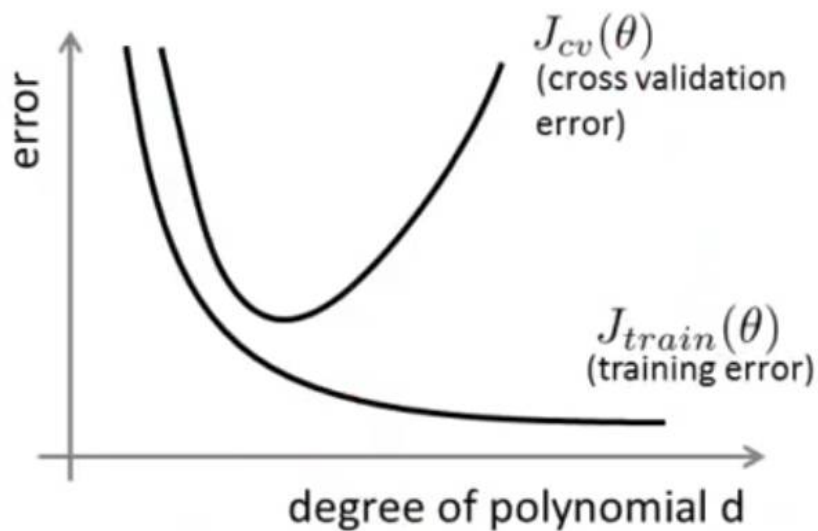
Cross validation error: $J_{cv}(\theta) = \frac{1}{2m_{cv}} \sum_{i=1}^{m_{cv}} (h_{\theta}(x_{cv}^{(i)}) - y_{cv}^{(i)})^2$ (or $J_{test}(\theta)$)



Andrew Ng

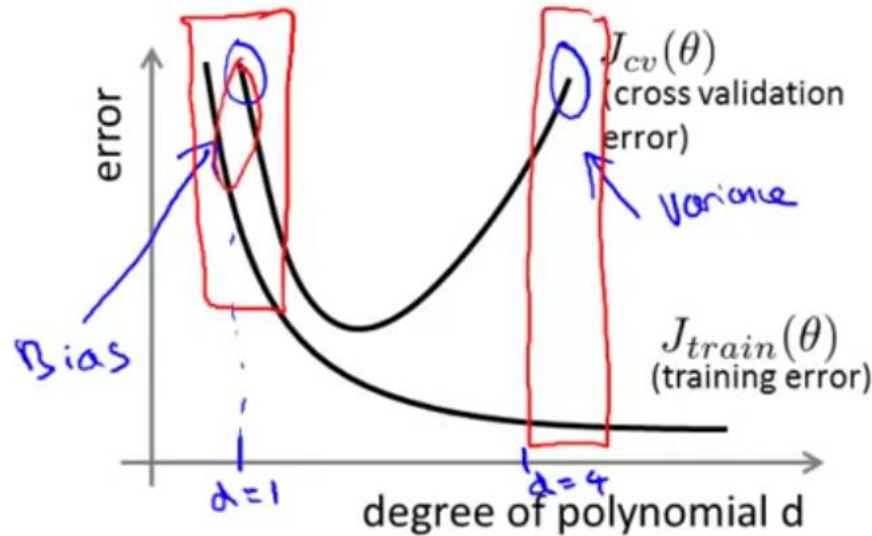
Diagnosing bias vs. variance

Suppose your learning algorithm is performing less well than you were hoping. ($J_{cv}(\theta)$ or $J_{test}(\theta)$ is high.) Is it a bias problem or a variance problem?



Diagnosing bias vs. variance

Suppose your learning algorithm is performing less well than you were hoping. ($J_{cv}(\theta)$ or $J_{test}(\theta)$ is high.) Is it a bias problem or a variance problem?



Bias (underfit):

$$\left. \begin{array}{l} J_{train}(\theta) \text{ will be high} \\ J_{cv}(\theta) \approx J_{train}(\theta) \end{array} \right\}$$

Variance (overfit):

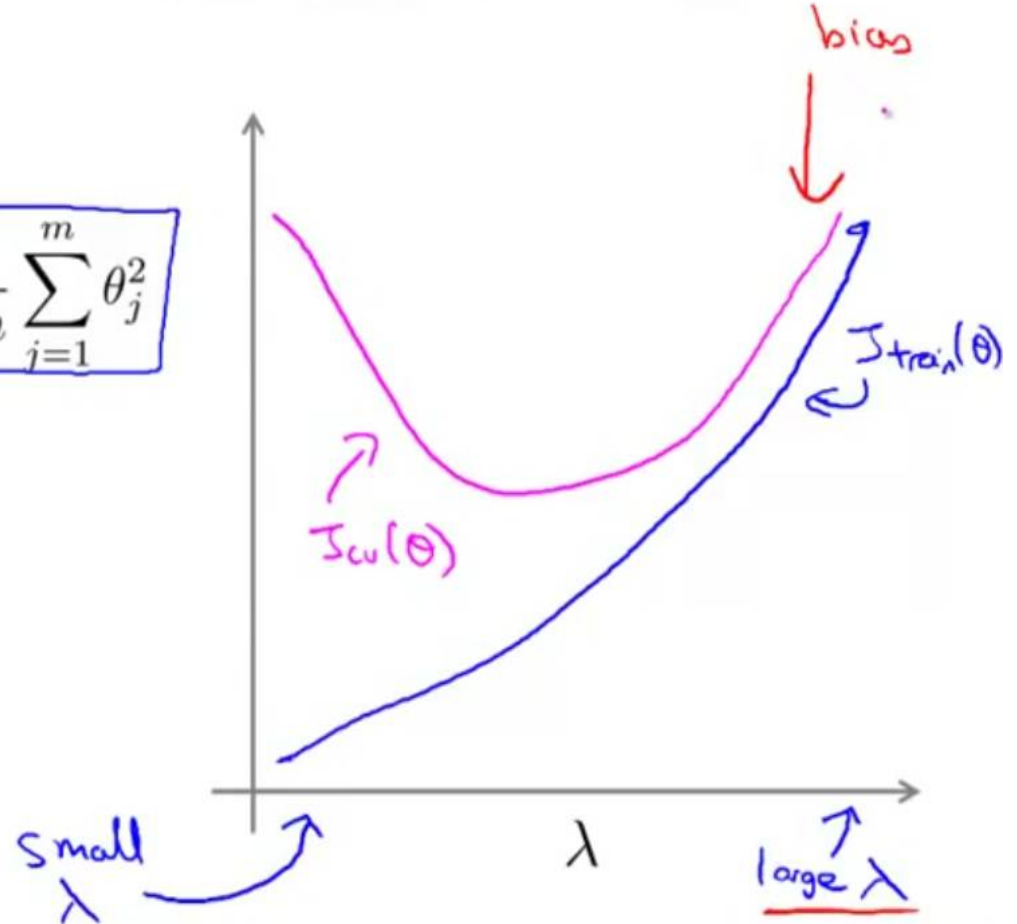
$$\left. \begin{array}{l} J_{train}(\theta) \text{ will be low} \\ J_{cv}(\theta) \gg J_{train}(\theta) \end{array} \right\}$$

Bias/variance as a function of the regularization parameter λ

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

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

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



References

- Sebastian Raschka – Slides + Code
 - <https://github.com/rasbt/stat451-machine-learning-fs20/tree/master>
- Andrew Ng: Coursera course on Introduction to Machine Learning
- Stanford CS229 (Andrew Ng).

End of Lecture