

Use-case



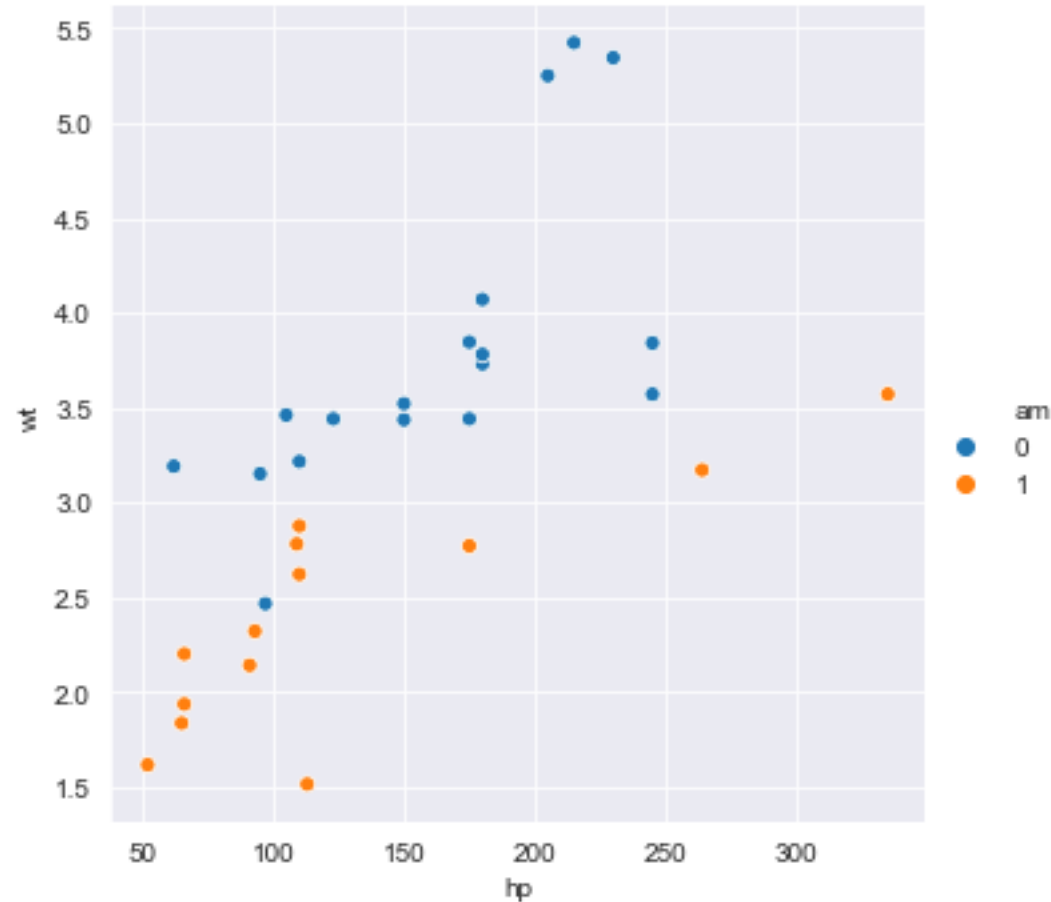
Interpreting Output- Deviance

model	mpg	cyl	disp	hp	drat	wt	qsec	vs	am	gear	carb	
Mazda RX4	21	6	160	110	3.9	2.62	16.46	0	1	4	4	
Mazda RX4 Wag	21	6	160	110	3.9	2.875	17.02	0	1	4	4	
Datsun 710	22.8	4	108	93	3.85	2.32	18.61	1	1	4	1	
Hornet 4 Drive	21.4	6	258	110	3.08	3.215	19.44	1	0	3	1	
Hornet Sportabout	18.7	8	360	175	3.15	3.44	17.02	0	0	3	2	
Valiant	18.1	6	225	105	2.76	3.46	20.22	1	0	3	1	
Duster 360	14.3	8	360	245	3.21	3.57	15.84	0	0	3	4	
Merc 240D	24.4	4	146.7	62	3.69	3.19	20	1	0	4	2	
Merc 230	22.8	4	140.8	95	3.92	3.15	22.9	1	0	4	2	
Merc 280	19.2	6	167.6	123	3.92	3.44	18.3	1	0	4	4	
Merc 280C	17.8	6	167.6	123	3.92	3.44	18.9	1	0	4	4	
Merc 450SE	16.4	8	275.8	180	3.07	4.07	17.4	0	0	3	3	
Merc 450SL	17.3	8	275.8	180	3.07	3.73	17.6	0	0	3	3	
Merc 450SLC	15.2	8	275.8	180	3.07	3.78	18	0	0	3	3	
Cadillac Fleetwood	10.4	8	472	205	2.93	5.25	17.98	0	0	3	4	
Lincoln Continental	10.4	8	460	215	3	5.424	17.82	0	0	3	4	
Chrysler Imperial	14.7	8	440	230	3.23	5.345	17.42	0	0	3	4	
Fiat 128	32.4	4	78.7	66	4.08	2.2	19.47	1	1	4	1	
Honda Civic	30.4	4	75.7	52	4.93	1.615	18.52	1				
Toyota Corolla	33.9	4	71.1	65	4.22	1.835	19.9	1				
Toyota Corona	21.5	4	120.1	97	3.7	2.465	20.01	1				
Dodge Challenger	15.5	8	318	150	2.76	3.52	16.87	0				

mpg	Miles/(US) gallon
cyl	Number of cylinders
disp	Displacement (cu.in.)
hp	Gross horsepower
drat	Rear axle ratio
wt	Weight (1000 lbs)
qsec	1/4 mile time
vs	V/S
am	Transmission (0 = automatic, 1 =
gear	Number of forward gears
carb	Number of carburetors

Example: Automatic or Manual Transmission

Using the MTcars dataset, estimate the probability of a vehicle being fitted with a manual transmission if it has a 120hp engine and weights 2800 lbs.



Example- Will the client subscribe a term deposit or not?

A Portuguese banking institution conducted a direct marketing campaign based on phone calls. Often, more than one contact to the same client was required, in order to assess if the product (bank term deposit) would be subscribed ('yes') or not ('no').



Example – Will the client subscribe a term deposit or not ?

bank client data

- *age (numeric)*
- *job: type of job (categorical: 'admin.', 'blue-collar', 'entrepreneur', 'housemaid', 'management', 'retired', 'self-employed', 'services', 'student', 'technician', 'unemployed', 'unknown')*
- *marital: marital status (categorical: 'divorced', 'married', 'single', 'unknown'; note: 'divorced' means divorced or widowed)*



Example – Will the client subscribe a term deposit or not

bank client data

- *education (categorical:*
'basic.4y', 'basic.6y', 'basic.9y', 'high.school', 'illiterate', 'professional.course', 'university.degree', 'unknown')
- *default: has credit in default? (categorical: 'no', 'yes', 'unknown')*
- *balance: money in account at the end of the year (numeric)*
- *housing: has housing loan? (categorical: 'no', 'yes', 'unknown')*
- *loan: has personal loan? (categorical: 'no', 'yes', 'unknown')*



Example – Will the client subscribe a term deposit or not

related with the last contact of the current campaign

- *contact*: contact communication type (categorical: 'cellular', 'telephone')
- *month*: last contact month of year (categorical: 'jan', 'feb',..., 'nov', 'dec')
- *day_of_week*: last contact day of the week (categorical: 'mon', 'tue', 'wed', 'thu', 'fri')
- *duration*: last contact duration, in seconds (numeric). Important note: this attribute highly affects the output target (e.g., if duration=0 then y='no'). Yet, the duration is not known before a call is performed. Also, after the end of the call, *y* is obviously known. Thus, this input should only be included for benchmark purposes and should be discarded if the intention is to have a realistic predictive model.



Example – Will the client subscribe a term deposit or not

other attributes

- *campaign*: number of contacts performed during this campaign and for this client (numeric, includes last contact)
- *pdays*: number of days that passed by after the client was last contacted from a previous campaign (numeric; 999 means client was not previously contacted)
- *previous*: number of contacts performed before this campaign and for this client (numeric)
- *poutcome*: outcome of the previous marketing campaign (categorical: 'failure', 'nonexistent', 'success')




```
Call: glm(formula = y ~ job + marital + education + balance + housing +
  loan + contact + day + month + duration + campaign + previous +
  poutcome, family = "binomial", data = subscribetermdeposit)
```

Coefficients:

(Intercept)	jobblue-collar	jobentrepreneur	jobhousemaid
-2.555e+00	-3.103e-01	-3.573e-01	-5.028e-01
jobmanagement	jobretired	jobself-employed	jobservices
-1.652e-01	2.552e-01	-2.981e-01	-2.241e-01
jobstudent	jobtechnician	jobunemployed	jobunknown
3.819e-01	-1.758e-01	-1.771e-01	-3.124e-01
maritalmarried	maritalsingle	educationsecondary	educationtertiary
-1.792e-01	9.171e-02	1.832e-01	3.790e-01
educationunknown	balance	housingyes	loanyes
2.506e-01	1.289e-05	-6.767e-01	-4.259e-01
contacttelephone	contactunknown	day	monthaug
-1.629e-01	-1.622e+00	9.976e-03	-6.931e-01
monthdec	monthfeb	monthjan	monthjul
6.920e-01	-1.458e-01	-1.260e+00	-8.305e-01
monthjun	monthmar	monthmay	monthnov
4.544e-01	1.591e+00	-4.001e-01	-8.706e-01
monthoct	monthsep	duration	campaign
8.828e-01	8.741e-01	4.194e-03	-9.082e-02
previous	poutcomeother	poutcomesuccess	poutcomeunknown
1.022e-02	2.049e-01	2.298e+00	-6.803e-02

Degrees of Freedom: 45210 Total (i.e. Null); 45171 Residual
 Null Deviance: 32630
 Residual Deviance: 21560 AIC: 21640



Applications

- Predicting stock price movement (up/down)
- Predict whether a patient has diabetes or not
- Predict whether a customer will buy or not
- Predict the likelihood of loan default



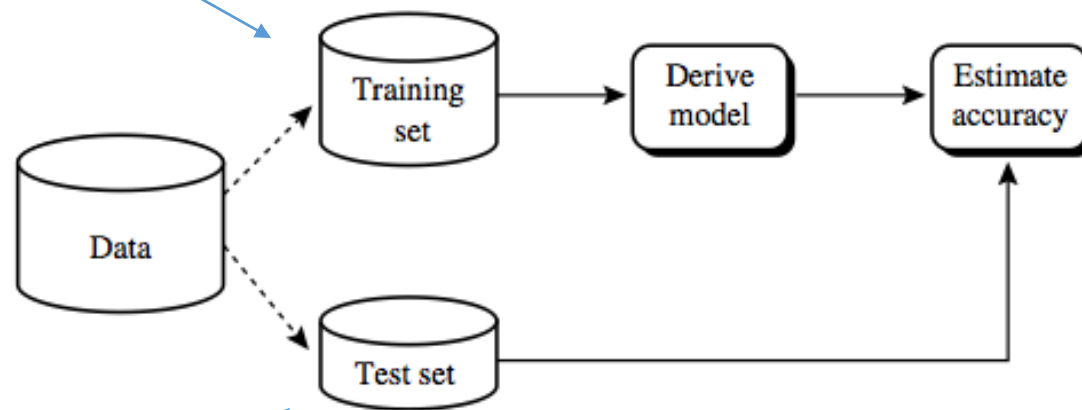
Diagnostic Hints

- Overly large coefficient magnitudes, overly large error bars on the coefficient estimates, and the wrong sign on coefficient could be indications of correlated inputs.
- VIF can be used to check for multicollinearity. Python outputs a Generalized Variance Inflation Factor, which is obtained by correcting VIF to the degrees of freedom for categorical predictors. $GVIF = VIF^{(\frac{1}{2*df})}$



B	C	D	E	F
mpg	cyl	disp	hp	drat
21	6	160	110	3.9
21	6	160	110	3.9
22.8	4	108	93	3.85
21.4	6	258	110	3.08
18.7	8	360	175	3.15
18.1	6	225	105	2.76
14.3	8	360	245	3.21
24.4	4	146.7	62	3.69
22.8	4	140.8	95	3.92
19.2	6	167.6	123	3.92
17.8	6	167.6	123	3.92
16.4	8	275.8	180	3.07
17.3	8	275.8	180	3.07
15.2	8	275.8	180	3.07
10.4	8	472	205	2.93
10.4	8	460	215	3
14.7	8	440	230	3.23
32.4	4	78.7	66	4.08
30.4	4	75.7	52	4.93
33.9	4	71.1	65	4.22
21.5	4	120.1	97	3.7
15.5	8	318	150	2.76
15.2	8	304	150	3.15
13.3	8	350	245	3.73
19.2	8	400	175	3.08
27.3	4	79	66	4.08
26	4	120.3	91	4.43
30.4	4	95.1	113	3.77
15.8	8	351	264	4.22
19.7	6	145	175	3.62
15	8	301	335	3.54
21.4	4	121	109	4.11

70 %



30 %



Case – Framingham Heart Study



Framingham Heart Study

A Project of the National Heart, Lung, and Blood Institute and Boston University

- Committed to identifying common factors contributing to cardiovascular disease (CVD).
- Setup in the town of Framingham, MA in 1948.
- Random sample consisting of 2/3rds of adult population in the town.

AGE-SEX DISTRIBUTION AT ENTRY (1948)				
Age	29-39	40-49	50-62	Totals
Men	835	779	722	2,336
Women	1,042	962	869	2,873
Totals	1,877	1,741	1,591	5,209



Case Study – Data(framinghamheartstudy.org and MITx)

- 5209 men and women participated.
- Age range: 30-62
- People who had not yet developed overt symptoms of CVD or suffered a heart attack or stroke.
- Careful monitoring of Framingham Study population has led to identification of major CVD risk factors.
- Led to development of Framingham Risk Score, a gender specific algorithm used to estimate the 10-year cardiovascular risk of an individual:

<http://cvdrisk.nhlbi.nih.gov/>



Case Study – Predicting Coronary Heart Disease (CHD)

Data description

4240 observations; 15 predictor and 1 predicted variables

- *TenYearCHD* – To be predicted. Risk of having a heart attack or stroke in the next 10 years.

Predictors

- Demographic Risk Factors
 - *male*: Gender of subject – Yes or No
 - *age*: Age of subject at first examination
 - *education*: some high school (1), high school (2), some college/vocational college (3), college (4)



Case Study – Predicting Coronary Heart Disease (CHD)

- Behavioural Risk Factors
 - *currentSmoker*: Yes or No
 - *cigsPerDay*: No. of cigarettes smoked per day if smoker
- Medical History Risk Factors
 - *BPmeds*: On BP medication at the time of first examination – Yes or No
 - *prevalentStroke*: Did the subject have a previous stroke – Yes or No
 - *prevalentHyp*: Is the subject currently hypertensive – Yes or No
 - *diabetes*: Does the subject currently have diabetes – Yes or No



Case Study – Predicting Coronary Heart Disease (CHD)

- Risk Factors from First Examination
 - *totChol*: Total cholesterol (mg/dL)
 - *sysBP*: Systolic blood pressure (the higher number in BP result)
 - *diaBP*: Diastolic blood pressure (the lower number in BP result)
 - *BMI*: Body Mass Index (kg/m²)
 - *heartRate*: # of beats per minute
 - *glucose*: Blood glucose level (mg/dL)



Case Study – Predicting Coronary Heart Disease (CHD)

Approach

- “Randomly” split data into training and test in 70:30 ratio.
- Measure prediction accuracies on training and test data
- Although , the split is random, we need to make sure the frequency of the categories are roughly the same in both training and test set.



Test / Train split



Case Study – Predicting Coronary Heart Disease (CHD)

Results

- Significant variables that cannot be controlled
 - Gender
 - Age
 - Medical history
- Significant variables that can be controlled
 - Smoking habits
 - Cholesterol
 - Systolic BP
 - Blood glucose

```
Call:
glm(formula = TenYearCHD ~ ., family = binomial, data = train)
```

Deviance Residuals:

Min	1Q	Median	3Q	Max
-1.9392	-0.5998	-0.4211	-0.2771	2.8632

Coefficients:

	Estimate	Std. Error	z value	Pr(> z)	
(Intercept)	-8.360272	0.864696	-9.668	< 2e-16	***
male	0.524080	0.130836	4.006	6.19e-05	***
age	0.065429	0.008049	8.129	4.34e-16	***
education	-0.041105	0.059185	-0.695	0.487366	
currentSmoker	0.120498	0.187629	0.642	0.520735	
cigsPerDay	0.016471	0.007488	2.200	0.027825	*
BPMeds	0.169118	0.282140	0.599	0.548898	
prevalentStroke	1.156666	0.560179	2.065	0.038940	*
prevalentHyp	0.307077	0.166034	1.849	0.064389	.
diabetes	-0.319937	0.392574	-0.815	0.415087	
totchol	0.003799	0.001330	2.856	0.004290	**
sysBP	0.011144	0.004446	2.507	0.012188	*
diaBP	-0.001861	0.007760	-0.240	0.810517	
BMI	0.008812	0.015662	0.563	0.573702	
heartRate	-0.007273	0.005131	-1.418	0.156296	
glucose	0.009227	0.002752	3.353	0.000798	***

Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

(Dispersion parameter for binomial family taken to be 1)

Null deviance: 2176.6 on 2565 degrees of freedom
Residual deviance: 1919.9 on 2550 degrees of freedom
(402 observations deleted due to missingness)
AIC: 1951.9



Missing Values

There are several ways of dealing with missing values. If large percentage of data for a given variable is missing, then we don't use that variable for building the model.

If the percentage of missing values is small (5 to 10%)

- Naïve method: Replace the missing values with either mean, median or mode
- Intelligent method: Impute the missing values from the relationship between the variables.

See for eg: <https://www.r-bloggers.com/imputing-missingdata-with-r-mice-package/>



Case Study – Predicting Coronary Heart Disease (CHD)

- Results
- Accuracy in training set = $2200/2566 = 85.7\%$
- Accuracy in testing set = $927/1092 = 84.9\%$
- Accuracy is affected by imbalance between positives and negatives.
- There is a trade-off between sensitivity and specificity.



Some More Performance Measures for Regression and Classification Models



Kappa Metric

- Accuracy can often be a misleading metric, when one category occurs more often than other in the given data-set
 - For eg: Occurrence of cancer in general population is 0.4%
 - If a prediction system blindly marks everyone as “No cancer”, it will 99.6% accurate



Kappa Metric

- Kappa metric quantifies how accurate the prediction algorithm is when compared to a random prediction

$$\text{kappa} = \frac{\text{totalAccuracy} - \text{randomAccuracy}}{1 - \text{randomAccuracy}}$$

$$\text{total Accuracy} = \frac{\text{correctPredictions}}{\text{Total}}$$

$$\text{randomAccuracy} = \frac{\text{ActualFalse}}{\text{Total}} * \frac{\text{predictedFalse}}{\text{Total}} + \frac{\text{ActualTrue}}{\text{Total}} * \frac{\text{PredictedTrue}}{\text{Total}}$$



Kappa Value	
< 0	No agreement
0.0 to 0.2	Slight
0.2 to 0.4	Fair
0.4 to 0.6	Moderate
0.6 to 0.8	Substantial
0.8 to 1.0	Almost Perfect



Kappa Metric

- Total= 30+357+9+2170=2566
- TotalAccuracy=(30+2170)/2566=0.857
- PercTrue=(30+357)/2566 = 0.15 ; PercFalse=(9+2170)/2566 = 0.85
- PredTrue=(30+9)/2566=0.015 ;PredFalse=(357+2170)/2566 = 0.985
- randomAccuracy= 0.15*0.015 + 0.85*0.985 = 0.84

$$Kappa = \frac{TotalAcc - randomAcc}{1 - randomAcc} = \frac{0.857 - 0.84}{1 - 0.84} = 0.10$$

Slightly better than Random!



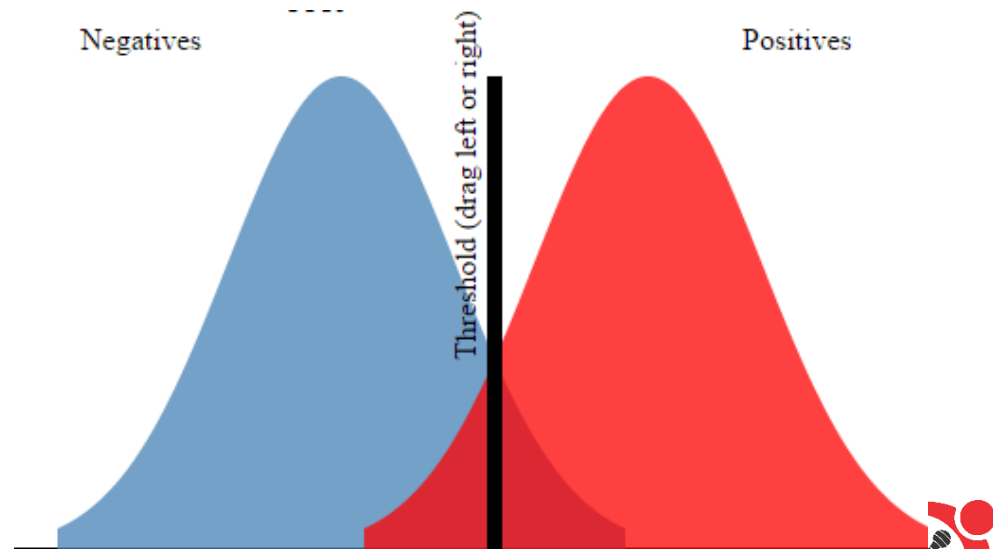
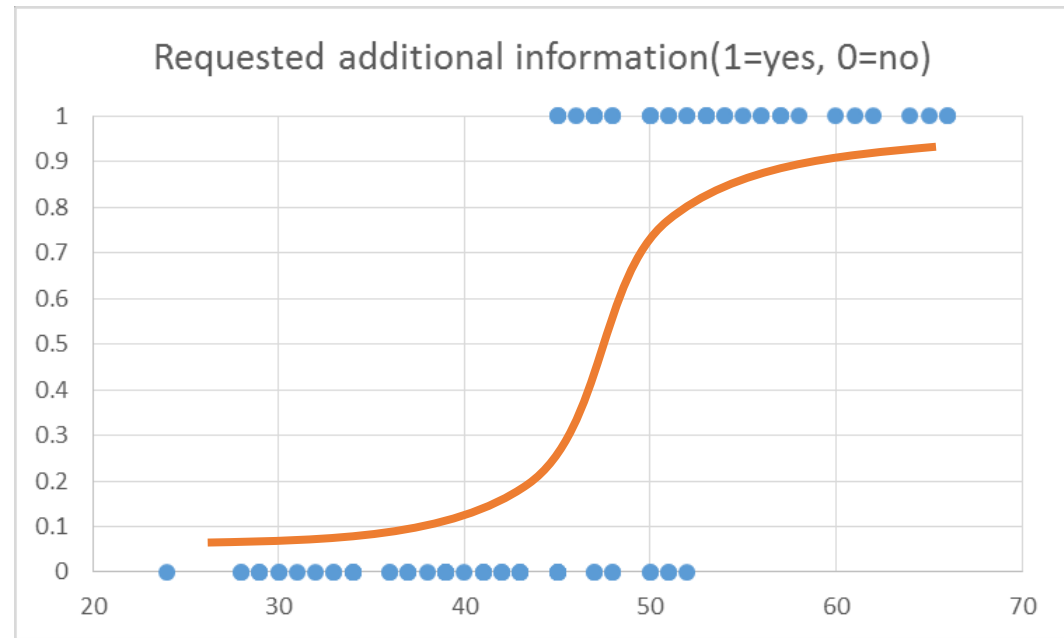
ROC Curves and AUC

- ROC – Receiver Operating Characteristics
- AUC – Area Under the ROC Curve



Logistic regression gives Probability forecasts for the given data point to be in a given bucket.

A threshold needs to be chosen to finally translate this probability to a bucket allocation



- At a given threshold, we can evaluate the classification accuracy (accuracy, sensitivity, recall, kappa etc)
- ROC curve tries to evaluate how well the regression has achieved the separation between the classes at all threshold values

ROC Curve Demo

- <http://www.navan.name/roc/>
- See: <https://youtu.be/OAl6eAyP-yo>



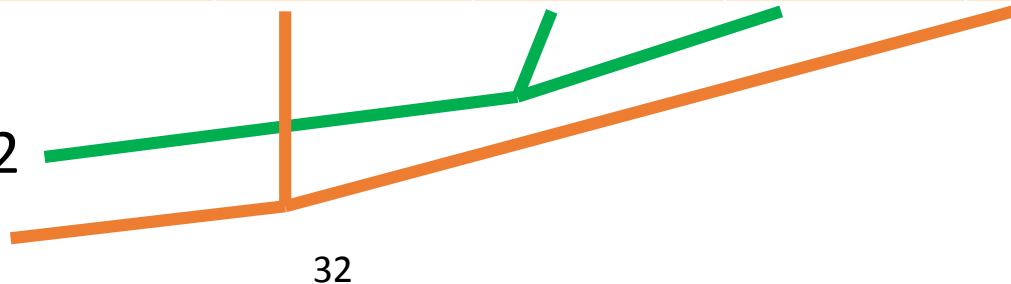
ROC Curves and AUC

- ROC – Plot of True Positive Rate vs False Positive Rate, i.e., Sensitivity vs 1-Specificity

Probability Threshold for Discriminating Between High Risk and Low Risk of Having Ten Year CHD	True Positives	False Positives	True Negatives	False Negative
0.9	0	0	922	170
0.7	1	1	921	169
0.5	12	7	915	158
0.3	46	76	846	124
0.1	140	468	454	30

Actual Counts

- Without CHD: 922
- With CHD: 170



32



ROC Curves and AUC

- ROC – Plot of True Positive Rate vs False Positive Rate, i.e., Sensitivity vs 1-Specificity

Probability Threshold for Discriminating Between High Risk and Low Risk of Having Ten Year CHD	Sensitivity		Specificity	
	True Positives	False Positives	True Negatives	False Negative
0.9	0	0	922	170
0.7	1	1	921	169
0.5	12	7	915	158
0.3	46	76	846	124
0.1	140	468	454	30

Actual Counts

- Without CHD: 922
- With CHD: 170

ROC Curve



ROC Curves and AUC

- ROC – Plot of True Positive Rate vs False Positive Rate, i.e., Sensitivity vs 1-Specificity

Probability Threshold for Discriminating Between High Risk and Low Risk of Having Ten Year CHD	Sensitivity	
	True Positives	False Positives
0.9	0	0
0.7	1	1
0.5	12	7
0.3	46	76
0.1	140	468

ROC Curve

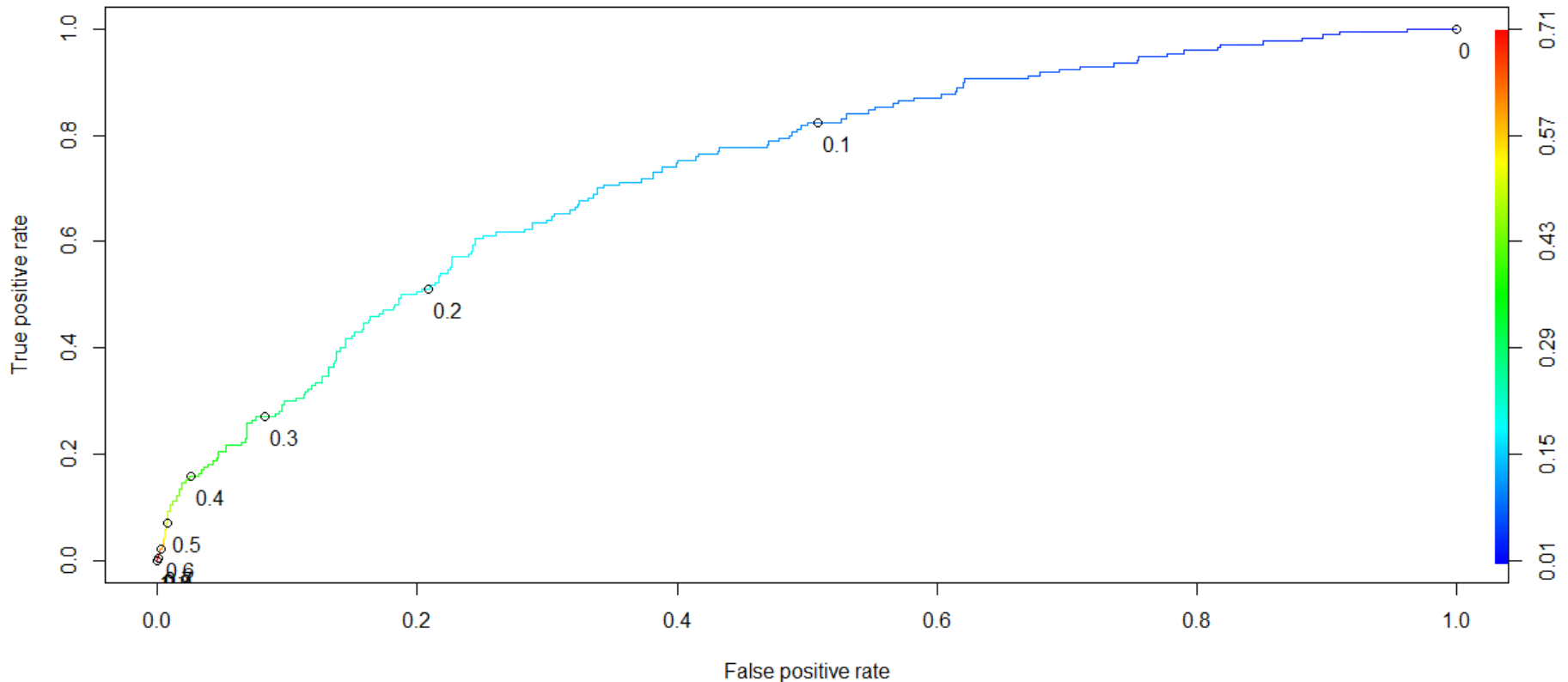
$P(\text{Predicting CHD} \mid \text{Have CHD})$

$P(\text{Predicting CHD} \mid \text{Do Not Have CHD})$



ROC Curves and AUC

- ROC – Plot of True Positive Rate vs False Positive Rate, i.e., Sensitivity vs 1-Specificity



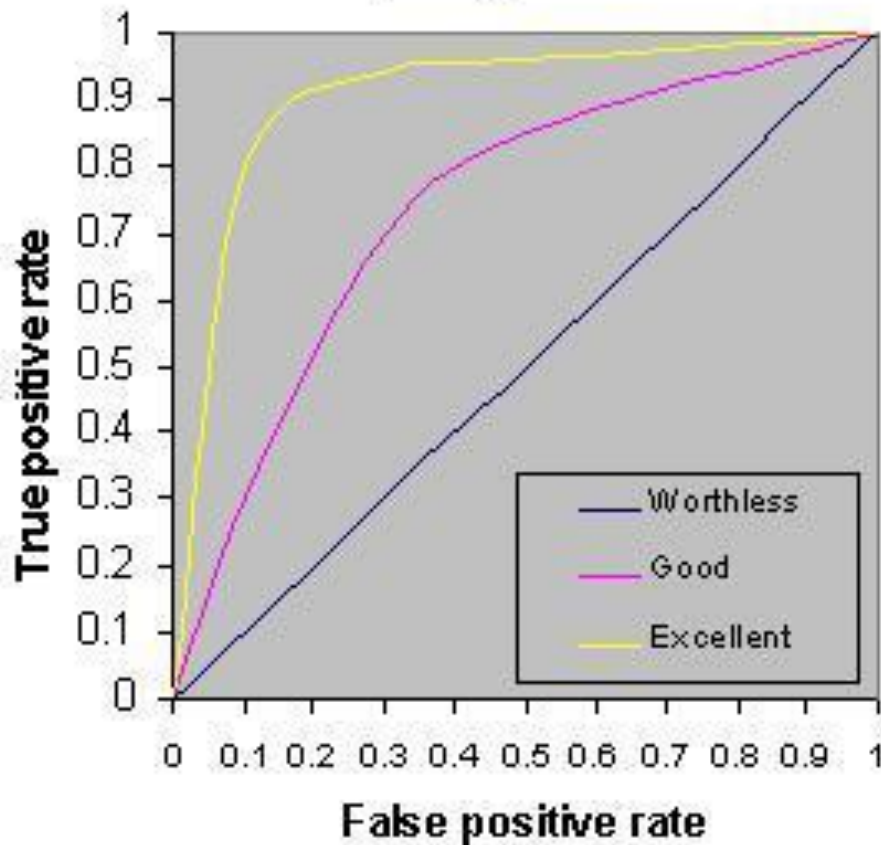
ROC Curves and AUC

- AUC – Measures discrimination, i.e., ability to correctly classify those with and without CHD.
- If you randomly pick one person who HAS CHD and one who DOESN'T and run the model, the one with the higher probability should be from the high risk group.
- AUC is the percentage of randomly drawn such pairs for which the classification is done correctly.



ROC Curves and AUC

Comparing ROC Curves



Rough rule of thumb:

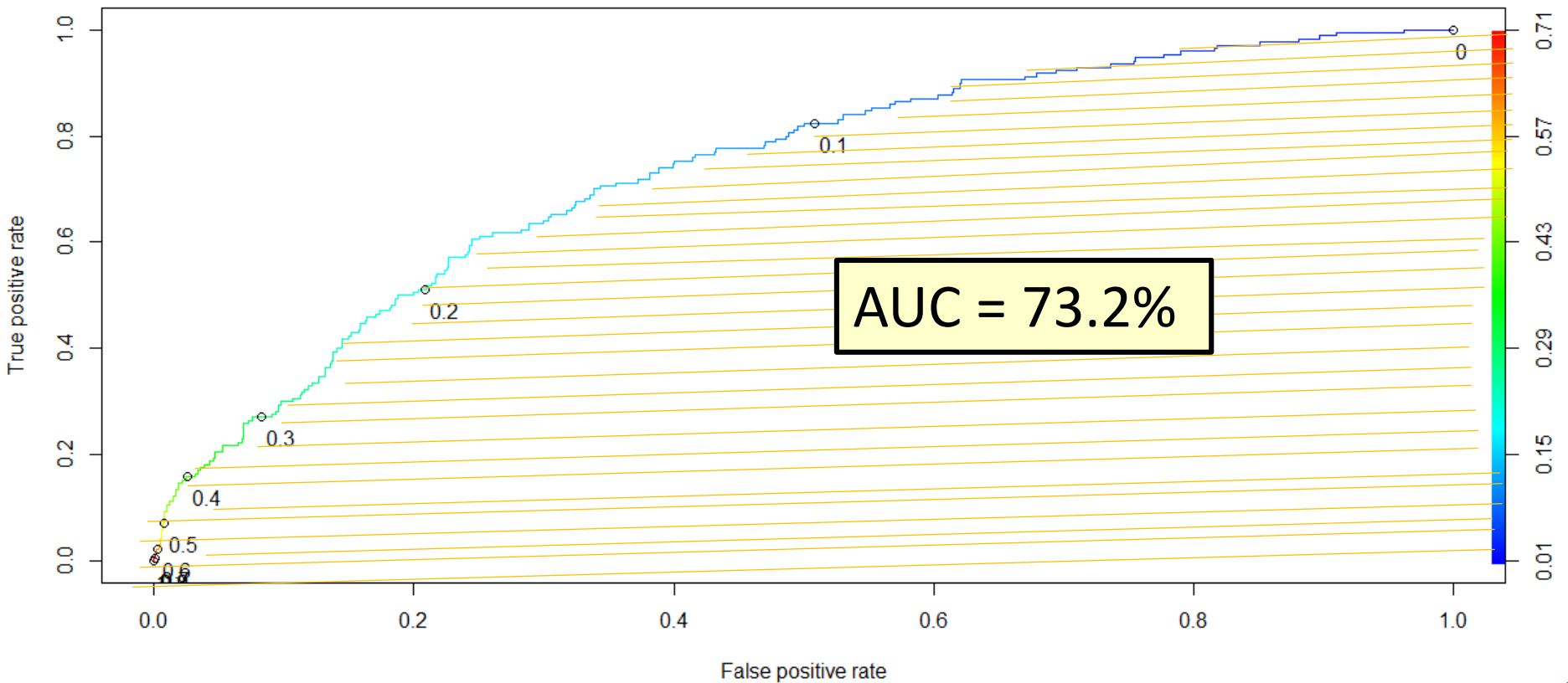
- 0.90 - 1.0 = Excellent
- 0.80 - 0.90 = Good
- 0.70 - 0.80 = Fair
- 0.60 - 0.70 = Poor
- 0.50 - 0.60 = Fail

• < 0.50 – You are better off doing a coin toss than working hard to build a model 😊



ROC Curves and AUC

- The model does a fair job of discrimination between high risk and low risk people.
- Useful for comparing different models.



Gains and Lift Charts

- In some business problems, it is not good enough to just classify. For example, in direct mail or phone marketing campaigns, where it costs money to send a mail to each prospect, it is better to be able to rank the prospective buyers by their probability to buy. That way, you can order them and start calling or mailing them in their decreasing order of propensity to buy.
- **Lift** is a measure of the effectiveness of a predictive model calculated as the ratio between the results obtained with and without the predictive model (random selection).



Gains and Lifts Charts

- A Lift Chart describes how well a model ranks samples in a particular class.
- The greater the area between the lift curve and the baseline (random selection), the better the model.



Gains and Lift Charts

- A company sends mail catalogs to prospective buyers. It costs the company \$1 to print and mail one catalog.
- From past data, they know the response rate is 5%, i.e., if 100,000 prospective customers are contacted, 5000 buy.
- This means that if there is no model and the company randomly contacts the prospects, they will have the following result.

No. of customers contacted	No. of responses
1000	500
2000	1000
3000	1500
.	.
.	.
.	.
10000	5000



Gains and Lift Charts

With a predictive model, where the model assigns a probability to each customer, the customers are ordered and divided into deciles (or any other quantiles). They are then called in decreasing order of probability to buy.

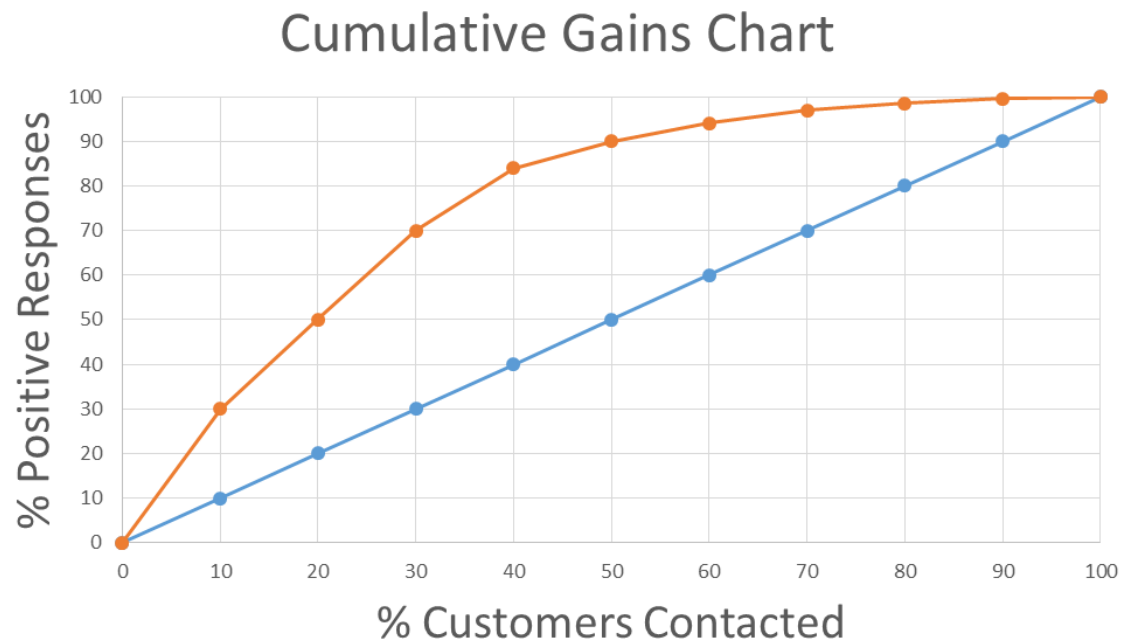
Cost	Decile contacted	Cumulative responses
1000	10	1500
2000	9	2500
3000	8	3500
4000	7	4200
5000	6	4500
6000	5	4700
7000	4	4850
8000	3	4925
9000	2	4975
10000	1	5000



Gains and Lift Charts

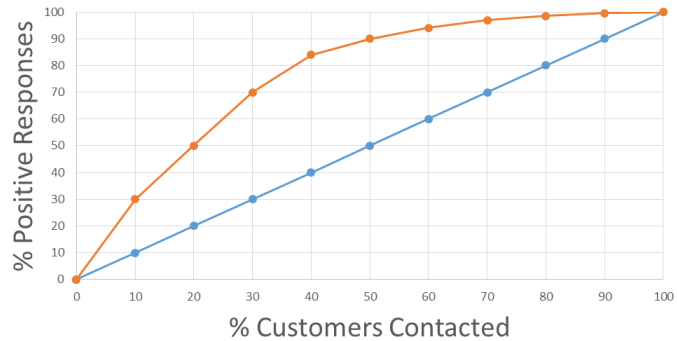
Cost	Decile contacted	Cumulative responses
1000	10	1500
2000	9	2500
3000	8	3500
4000	7	4200
5000	6	4500
6000	5	4700
7000	4	4850
8000	3	4925
9000	2	4975
10000	1	5000

% Called	Called at Random	Called According to Model Score
0	0	0
10	10	30
20	20	50
30	30	70
40	40	84
50	50	90
60	60	94
70	70	97
80	80	98.5
90	90	99.5
100	100	100



Gains and Lift Charts

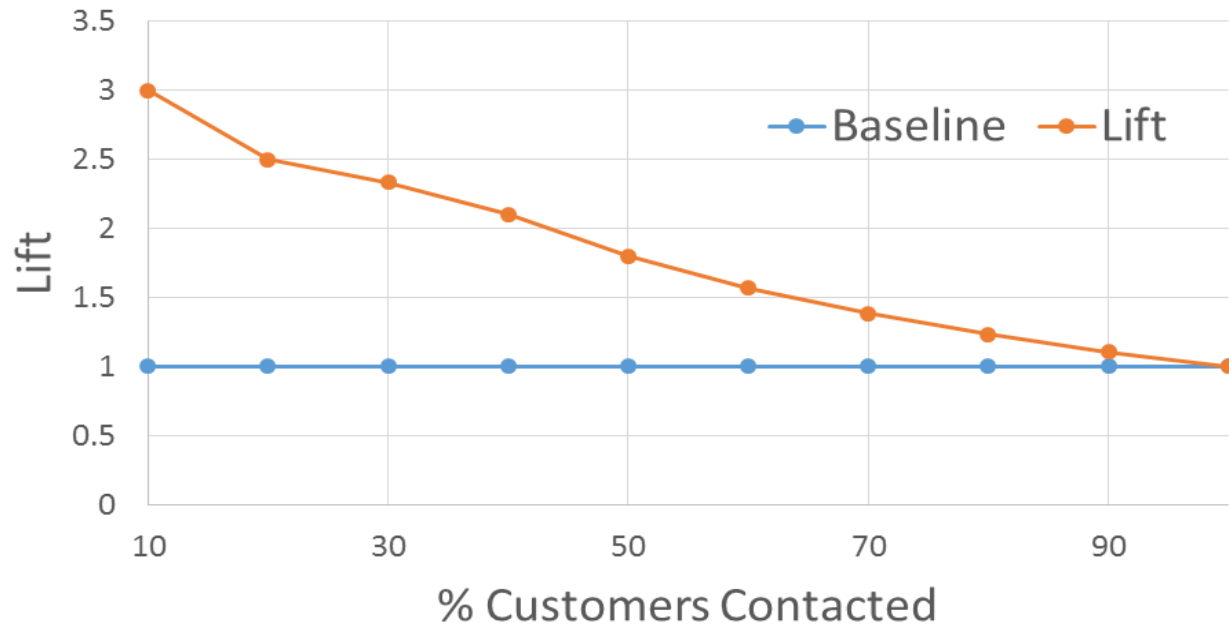
Cumulative Gains Chart



Max lift of 3 at the top decile.

Model advantage diminishes as more customers are contacted, especially in lower deciles.

Lift Chart



Useful to compare different models.



Evaluating Model Accuracy and Bias-Variance Tradeoff



The Ultimate Test of Model Accuracy

- Holdout set: Split data into train, validation and test sets (in 70:20:10 or 60:20:20, etc. ratios), and **ensure model performance is similar.**
 - Training Set: For fitting a model
 - Validation Set: For selecting a model based on estimated prediction errors
 - Test Set: For assessing selected model's performance on "new" data
- k-fold cross-validation: Same as holdout but useful when the data size is small.



Appropriate Error Measures for Evaluating Model Accuracy

- Use accurate measures of prediction error, experiment with different models and use the model with minimum error.
- Some measures for comparing models within the same technique (e.g., Linear Regression):
 - R^2
 - AIC



Appropriate Error Measures for Evaluating Model Accuracy

Some measures for comparing models across techniques:

- MAE (Mean Absolute Error): Mean of the absolute value of the difference between the predicted and actual values.
- MAPE (Mean Absolute Percentage Error): Same as above but converted into percentages to allow for comparison across different scales (e.g., comparing accuracies of forecasts on BSE vs NSE).
- RMSE (Root Mean Square Error): Accounts for infrequent large errors, whose impact may be understated by the mean-based error measures.



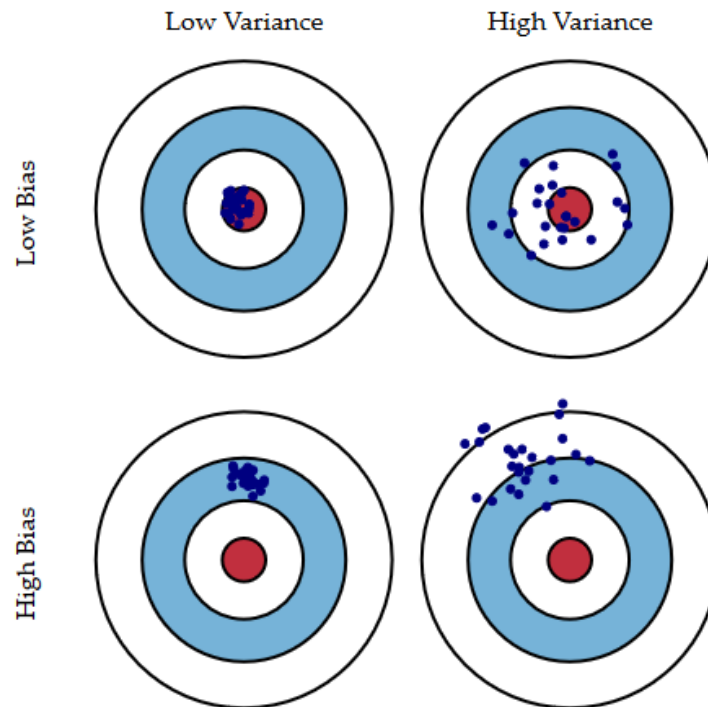
Bias-Variance Tradeoff

- Total error is composed of Bias, Variance and a Random irreducible error. Bias and Variance can be managed.
- If the model performance on training and testing data sets is inconsistent, it indicates a problem either with Bias or Variance.

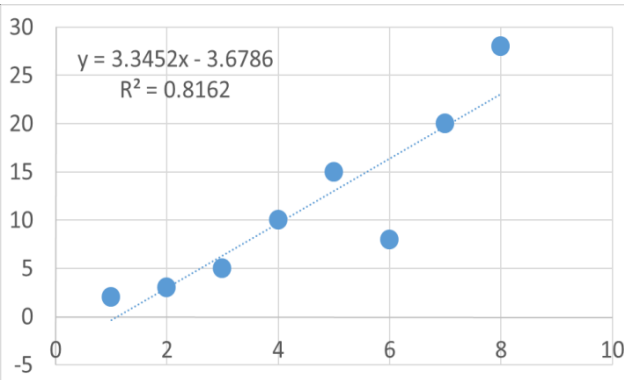


Bias – Variance Tradeoff

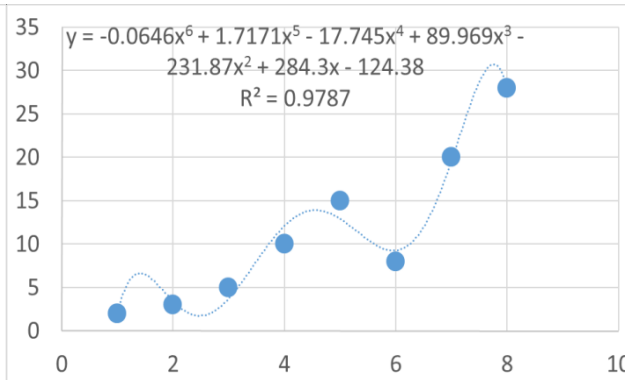
- Bulls-eye is a model that correctly predicts the real values.
- Each hit is a model based on chance variability in training datasets.



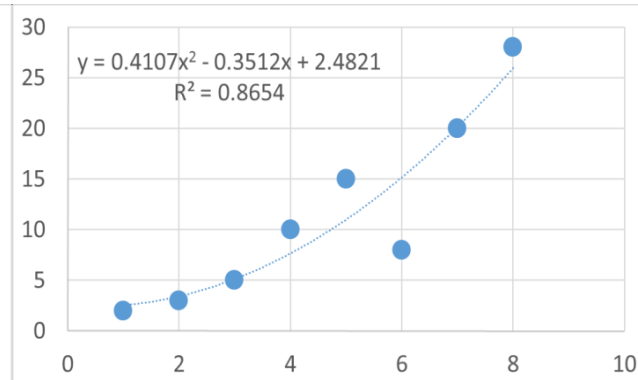
Bias-Variance Tradeoff and Under fitting vs Over fitting Excel



Too Simple a Model
Underfit



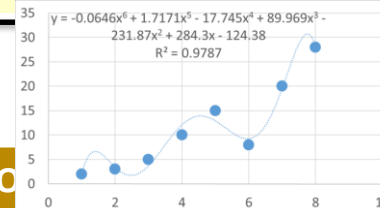
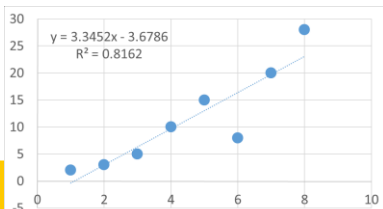
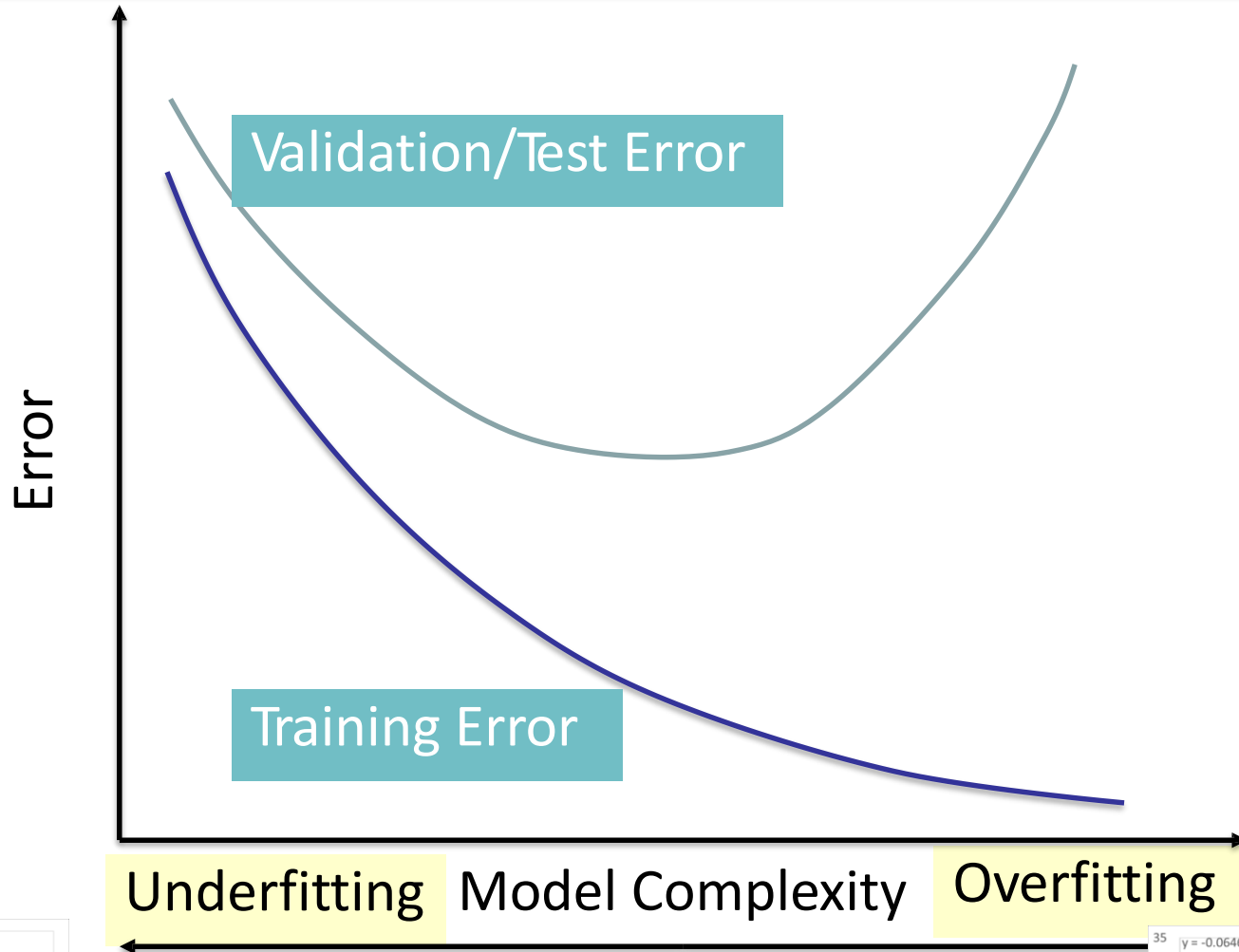
Too Complex a Model
Overfit



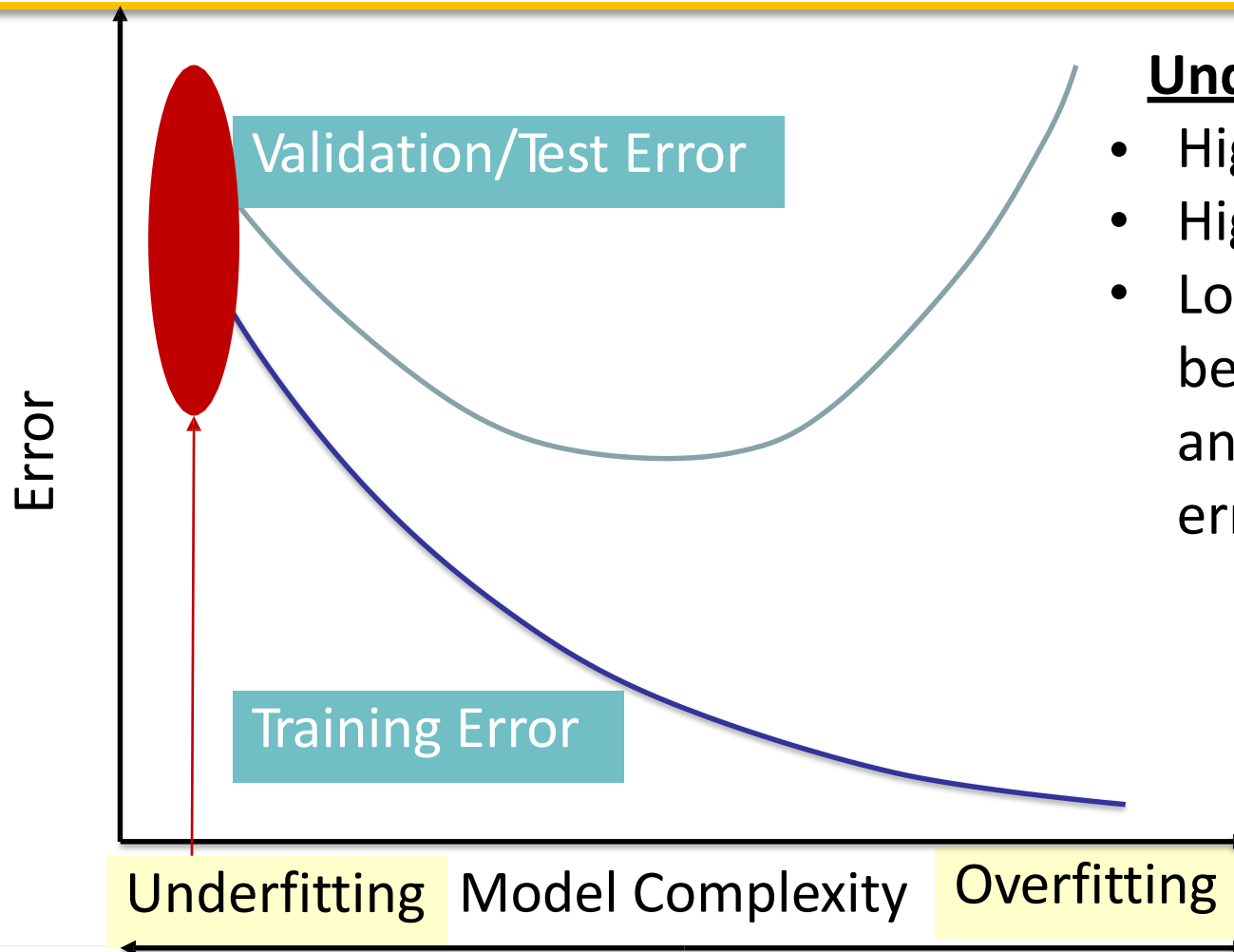
Right Model
Reasonable fit



Bias-Variance Tradeoff and Under fitting vs Over fitting

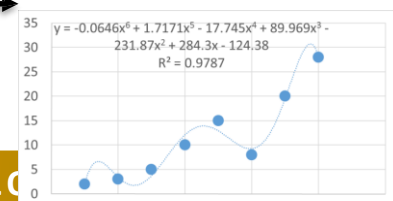
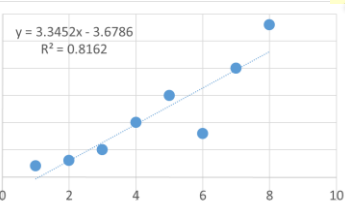


Diagnosing Bias and Variance

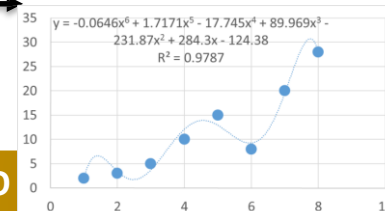
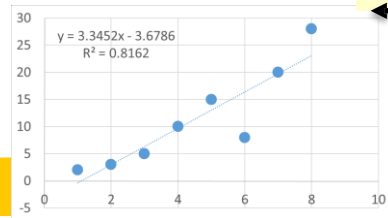
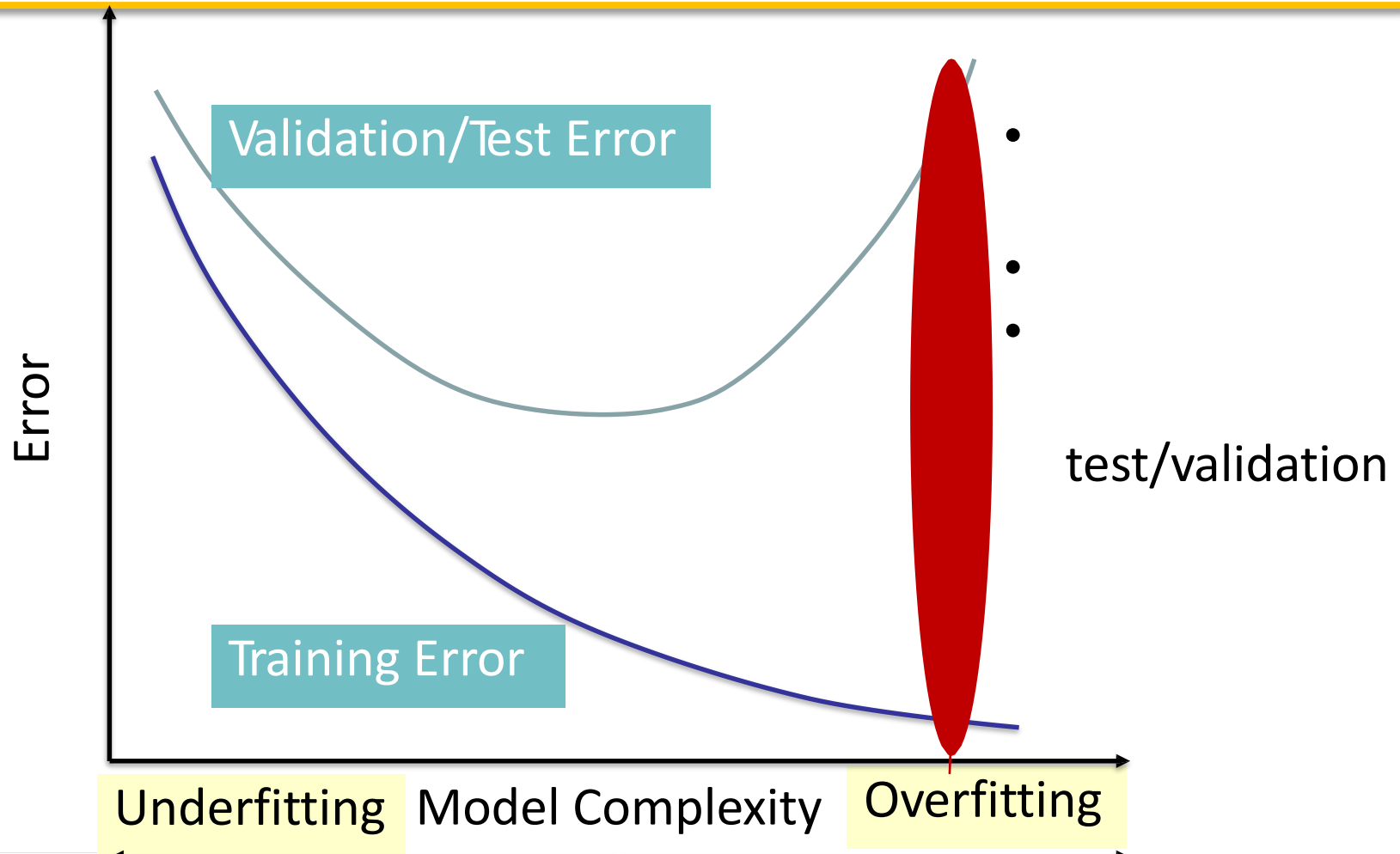


Underfitting (Bias)

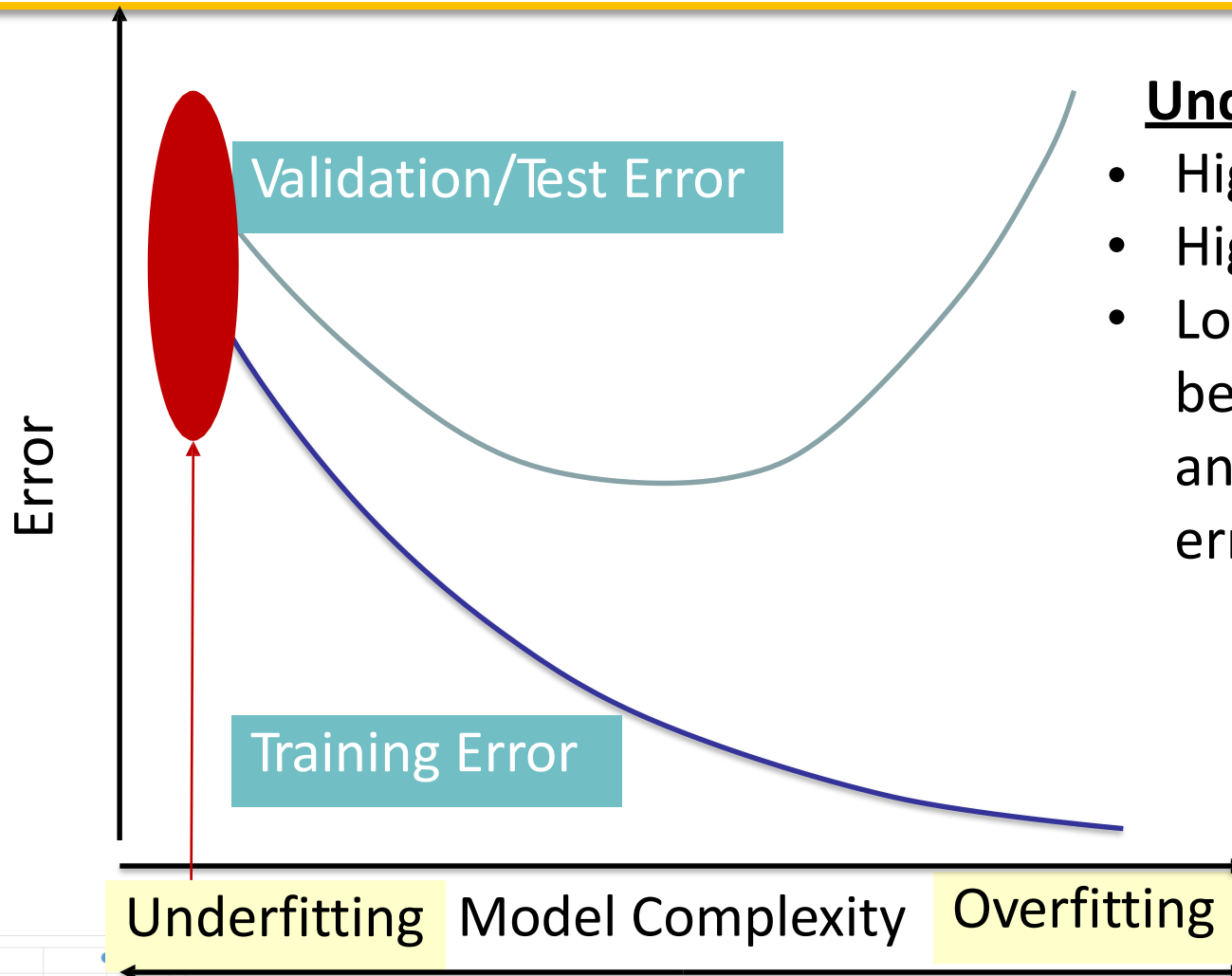
- High Bias problem
- High training error
- Low difference between training and test/validation errors



Diagnosing Bias and Variance

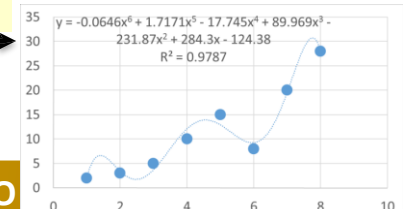
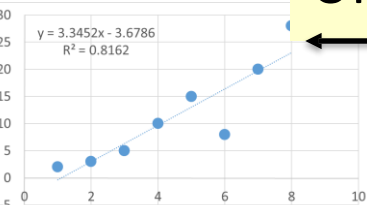


Diagnosing Bias and Variance



Underfitting (Bias)

- High Bias problem
- High training error
- Low difference between training and test/validation errors



Bias-Variance Tradeoff

Ways of detecting and minimizing Bias and Variance

Outliers and Influential Observations can cause statistical bias. Can be identified using various methods like Box plots, points outside ± 2 or ± 3 standard deviations/errors, residual plots, etc.

Bias cannot be corrected by increasing training sample size.

Variance or standard error can be minimized by increasing training sample size.

Bagging (bootstrap aggregating) techniques (*taught later in the program*) can be used to minimize errors.



Reference

