

What we discussed last week?

□ Exploratory data analysis

- What is it?
- Comparing two/multiple categorical variables?
- Exploring numeric variables?
- Deal with correlated variables?

□ Regression analysis

Statistical Analysis

Data Mining tasks in KDD

□ Chapter 1 explained Description, along with other Data Mining tasks – in Description:

- Analyst try to find patterns and trends
- Look for possible explanations for these
- Recommend possible policy changes

□ Description can be accomplished by

- Exploratory Data Analysis – Chapter 3
- Descriptive statistics – Chapter 4 & 5
 - Sample proportion or Regression equation

□ Data mining method can be applied to more than one task

- Example: Decision Trees used for Classification, Estimation and Prediction

Task	We Learn about This Task in
Description	Chapter 3: Exploratory Data Analysis Chapter 4: Univariate Statistical Analysis Chapter 5: Multivariate Statistical Analysis
Estimation	Chapter 4: Univariate Statistical Analysis Chapter 5: Multivariate Statistical Analysis
Prediction	Chapter 4: Univariate Statistical Analysis Chapter 5: Multivariate Statistical Analysis
Classification	Chapter 7: k-Nearest Neighbor Algorithm Chapter 8: Decision Trees Chapter 9: Neural Networks
Clustering	Chapter 10: Hierarchical and k-Means Clustering Chapter 11: Kohonen Networks
Association	Chapter 12: Association Rules

□ Table 4.1 presents a general outline on tasks and methods used to accomplish them

Univariate Analysis

Statistical Inference

▮ A widespread tool for performing **estimation** and **prediction** is *statistical inference*.

▮ Population

▮ A **population** is a collection of all elements of interest for a particular study

▮ A **parameter** is a characteristic of a population, such as the mean number of customer service calls of all cell phone customers

▮ Example: Cell phone company wants actionable results for all their present and future customers (population), not only the 3333 customers for which they gathered the data (sample)

▮ Sample

▮ A **sample** is a representative subset of the population

▮ If the sample characteristics deviate systematically from the population characteristics, *statistical inference should not be applied*

▮ A **statistic** is a characteristic of a sample, such as the mean number of customer service calls of the 3333 customers in the sample (1.563).

▮ Statistical Inference

▮ Methods for estimating and testing hypotheses about **population** characteristics based on information contained in a sample

Statistical Inference (*cont'd*)

- ▮ The values of population parameters are *unknown* for most problems

	Sample Statistics	...Estimates...	Population Parameters
Mean			
Standard deviation			
Proportion			

▮ Example: Churn Data Customer Service Calls

- ▮ Sample $n = 3333$
- ▮ Population = unknown
- ▮ Sample mean = 1.56
- ▮ Sample standard deviation $s = 1.315$

Statistical Inference (*cont'd*)

□ Example

- The proportion of those in the population who “churn” is unknown
- The proportion of those in the population who churn is estimated from the proportion of those in the sample who churn
- The sample proportion $p = 0.145$ is used to estimate the population proportion who churn
- In these examples, our estimation is only valid when the sample is truly representative of the entire population

How Confident Are We in Our Estimates?

- Is the population number of customer service calls the same as the sample mean =1.563? **Probably not.**
 - Population contains more information than the sample
 - Our point estimates will nearly always “miss” the target parameter by a certain amount
- Sampling Error
 - Distance between the observed value of the point estimate and the unknown value of its target parameter
 - Defined as: $|statistic - parameter|$
 - e.g., the sampling error for the mean is: $|\bar{X} - \mu|$ (always positive)
 - Value of the sampling error is usually unknown
 - For continuous variables, the probability that the observed value of a point estimate exactly equals its target parameter is precisely zero
- Point estimates have no measure of confidence in their accuracy because we don't know the population.

How Confident Are We in Our Estimates? (*cont'd*)

- ▮ Point estimation can be compared to throwing darts with **infinitesimally small tips** to a **vanishingly small bull's-eye**
 - ▮ Worse, the bull's-eye is hidden, and the thrower will never know for sure how close the darts are to the target
- ▮ As an alternative the dart thrower could throw a wide beer mug to the target
 - ▮ There does indeed exist a positive probability that some portion of the mug has hit the hidden bull's-eye
 - ▮ While not sure, we can have a **certain degree of confidence** that the target has been hit
 - ▮ **Very roughly, the beer mug represents our next estimation method, *confidence intervals***

Confidence Interval Estimation

- Consists of an interval of numbers produced by a point estimate
- Includes a confidence level specifying the probability the interval contains the population parameter
- Confidence intervals have the general form:

Point Estimate +/- Margin of Error

- The margin of error is a measure of precision, and smaller values indicate greater precision

Confidence Interval Estimation (*cont'd*)

□ Example

□ The t -interval for the population mean is:

$$\bar{x} \pm t_{\alpha/2} (s / \sqrt{n})$$

Where, the Point Estimate $= \bar{x}$, and Margin of Error $= t_{\alpha/2} (s / \sqrt{n})$
(s / \sqrt{n})

is the standard error of sample mean

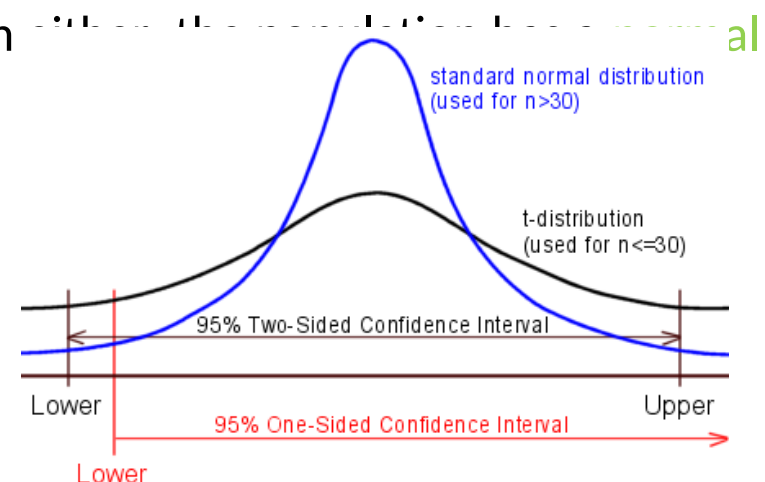
□ The t -interval for the mean is used when **normal distribution**; or the sample size n is large

□ For normal distribution ($n > 30$)

□ = 1.645 for 90% confidence

□ = 1.960 for 95% confidence

□ = 2.576 for 99% confidence

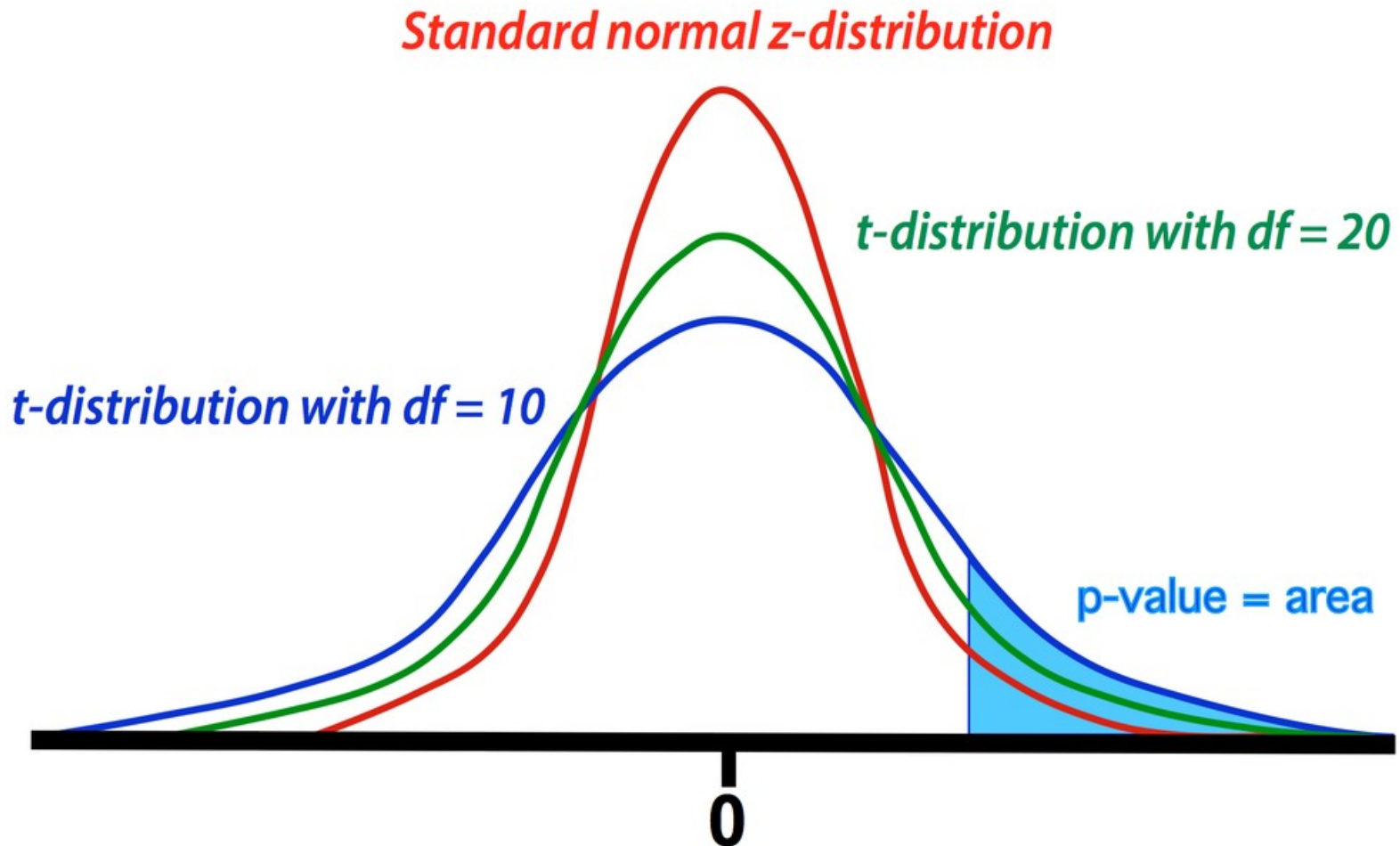


t-Distribution Table

***t* Table**

cum. prob	<i>t</i> _{.50}	<i>t</i> _{.75}	<i>t</i> _{.80}	<i>t</i> _{.85}	<i>t</i> _{.90}	<i>t</i> _{.95}	<i>t</i> _{.975}	<i>t</i> _{.99}	<i>t</i> _{.995}	<i>t</i> _{.999}	<i>t</i> _{.9995}
one-tail	0.50	0.25	0.20	0.15	0.10	0.05	0.025	0.01	0.005	0.001	0.0005
two-tails	1.00	0.50	0.40	0.30	0.20	0.10	0.05	0.02	0.01	0.002	0.001
df											
1	0.000	1.000	1.376	1.963	3.078	6.314	12.71	31.82	63.66	318.31	636.62
2	0.000	0.816	1.061	1.386	1.886	2.920	4.303	6.965	9.925	22.327	31.599
3	0.000	0.765	0.978	1.250	1.638	2.353	3.182	4.541	5.841	10.215	12.924
4	0.000	0.741	0.941	1.190	1.533	2.132	2.776	3.747	4.604	7.173	8.610
5	0.000	0.727	0.920	1.156	1.476	2.015	2.571	3.365	4.032	5.893	6.869
6	0.000	0.718	0.906	1.134	1.440	1.943	2.447	3.143	3.707	5.208	5.959
7	0.000	0.711	0.896	1.119	1.415	1.895	2.365	2.998	3.499	4.785	5.408
8	0.000	0.706	0.889	1.108	1.397	1.860	2.306	2.896	3.355	4.501	5.041
9	0.000	0.703	0.883	1.100	1.383	1.833	2.262	2.821	3.250	4.297	4.781
10	0.000	0.700	0.879	1.093	1.372	1.812	2.228	2.764	3.169	4.144	4.587
11	0.000	0.697	0.876	1.088	1.363	1.796	2.201	2.718	3.106	4.025	4.437
12	0.000	0.695	0.873	1.083	1.356	1.782	2.179	2.681	3.055	3.930	4.318
13	0.000	0.694	0.870	1.079	1.350	1.771	2.160	2.650	3.012	3.852	4.221
14	0.000	0.692	0.868	1.076	1.345	1.761	2.145	2.624	2.977	3.787	4.140
15	0.000	0.691	0.866	1.074	1.341	1.753	2.131	2.602	2.947	3.733	4.073
16	0.000	0.690	0.865	1.071	1.337	1.746	2.120	2.583	2.921	3.686	4.015
17	0.000	0.689	0.863	1.069	1.333	1.740	2.110	2.567	2.898	3.646	3.965
18	0.000	0.688	0.862	1.067	1.330	1.734	2.101	2.552	2.878	3.610	3.922
19	0.000	0.688	0.861	1.066	1.328	1.729	2.093	2.539	2.861	3.579	3.883
20	0.000	0.687	0.860	1.064	1.325	1.725	2.086	2.528	2.845	3.552	3.850
21	0.000	0.686	0.859	1.063	1.323	1.721	2.080	2.518	2.831	3.527	3.819
22	0.000	0.686	0.858	1.061	1.321	1.717	2.074	2.508	2.819	3.505	3.792
23	0.000	0.685	0.858	1.060	1.319	1.714	2.069	2.500	2.807	3.485	3.768
24	0.000	0.685	0.857	1.059	1.318	1.711	2.064	2.492	2.797	3.467	3.745
25	0.000	0.684	0.856	1.058	1.316	1.708	2.060	2.485	2.787	3.450	3.725
26	0.000	0.684	0.856	1.058	1.315	1.706	2.056	2.479	2.779	3.435	3.707
27	0.000	0.684	0.855	1.057	1.314	1.703	2.052	2.473	2.771	3.421	3.690
28	0.000	0.683	0.855	1.056	1.313	1.701	2.048	2.467	2.763	3.408	3.674
29	0.000	0.683	0.854	1.055	1.311	1.699	2.045	2.462	2.756	3.396	3.659
30	0.000	0.683	0.854	1.055	1.310	1.697	2.042	2.457	2.750	3.385	3.646
40	0.000	0.681	0.851	1.050	1.303	1.684	2.021	2.423	2.704	3.307	3.551
60	0.000	0.679	0.848	1.045	1.296	1.671	2.000	2.390	2.660	3.232	3.460
80	0.000	0.678	0.846	1.043	1.292	1.664	1.990	2.374	2.639	3.195	3.416
100	0.000	0.677	0.845	1.042	1.290	1.660	1.984	2.364	2.626	3.174	3.390
1000	0.000	0.675	0.842	1.037	1.282	1.646	1.962	2.330	2.581	3.098	3.300
∞	0.000	0.674	0.842	1.036	1.282	1.645	1.960	2.326	2.576	3.090	3.291
	0%	50%	60%	70%	80%	90%	95%	98%	99%	99.8%	99.9%
	Confidence Level										

t-distribution vs z-distribution



Confidence Interval Estimation (*cont'd*)

□ Example

- Find the 95% confidence t -interval for the mean number of Customer Service Calls
- Recall that the sample mean = 1.563,
- and sample standard deviation = 1.315

$$\begin{array}{l} n = 3333 \\ \bar{x} = 1.563 \end{array} \qquad \begin{array}{l} \bar{x} \pm t_{\alpha/2}(s/\sqrt{n}) \\ 1.563 \pm 1.96(1.315/\sqrt{3333}) \\ 1.563 \pm 0.045 \\ (1.518, 1.608) \end{array}$$

- We are ~95% confident the population mean number of Customer Service Calls is between 1.518 and 1.608. Our margin of error is **0.045**

Confidence Interval Estimation (*cont'd*)

- Data miners are often called upon to perform ***subgroup analyses***
 - This is, to estimate the behavior of specific subsets of customers instead of the entire customer base
 - Example: With the 95% confidence, estimate the mean number of customer service calls for customers who have both the International Plan and the VoiceMail Plan and who have more than 220 day minutes

Count = 28	$\bar{x} \pm t_{\alpha/2}(s/\sqrt{n})$
= 1.607	$1.607 \pm 2.048(1.892/\sqrt{28})$
= 1.892	1.607 ± 0.732
)

- We are 95% confident the population mean number of Customer Service Calls is between 0.875 and 2.339. Our margin of error is 0.732
- Here, our estimate for this specific subset of customers is **less precise** than our estimate obtained for the entire customer base in the previous Example, **why?**

Confidence Interval Estimation (*cont'd*)

□ Conclusion

- Confidence interval estimation can be applied to **any** desired target parameter.
- Most widely used in statistical practice to estimate the population mean, population standard deviation, and population proportion of successes, etc.
- It is important for data miners to know the confidence interval of their data if it is a sample of the population.

How to reduce the margin of error

- The margin of error E for a 95% confidence interval for the population mean can be interpreted as:

“We can estimate to within E units with 95% confidence.”

- The smaller the margin of error, the more precise our estimation
- *How can we **reduce** our margin of error?*

- Margin of error E contains three quantities

$$\bar{x} \pm t_{\alpha/2} (s / \sqrt{n})$$

- , which depends on the confidence level
- , which is a characteristic of the data, and may not be changed
- , the sample size

- Thus, we can decrease the margin of error (E) by

- Decreasing the confidence level, - NOT RECOMMENDED
- Increasing sample size, - RECOMMENDED

- **Increasing the sample size** is the **only** way to decrease the margin of error while maintaining a constant level of confidence

Confidence Interval Estimation of the **Proportion**

- An estimate of the *population proportion* π of all company's customers who churn is:
- But we have no measure of confidence on the accuracy of this estimate
- A *confidence interval for the population proportion* π is
- Where the sampling proportion p is the point estimate of π and represents the margin of error.
- represents the confidence level (1.645 for 90%, 1.96 for 95% and 2.576 for 99%)
- The Z-interval for π may be used whenever both and .

Confidence Interval Estimation of the Proportion (*cont'd*)

- Example: Calculate 95% confidence interval for the proportion π of churners among the entire population of the company's customers
- We are 95% confident that this interval captures the population proportion π . The confidence interval for π takes the form
- The margin of error can be interpreted as follows:
 - “We can estimate π to within E with 95% of confidence”
- In this case, we estimate proportion of churners to within 0.012 (or 1.2%) with 95% confidence.
- Again, given a confidence level, margin of error can be reduced by **using a larger sample**

Hypothesis Testing for the Mean

- Hypothesis testing *evaluates* claims about a population parameter using evidence from the sample
- Two competing statements (hypotheses) are crafted about the parameter value:
 - Null Hypothesis (H_0) – represents assumed value
 - Alternative Hypothesis (H_a) – represents alternative claim about the value
- There are two possible conclusions:
 - Reject H_0
 - Do not reject H_0

Hypothesis Testing for the Mean (*cont'd*)

- Example: A criminal trial is a form of hypothesis test, with following hypotheses:
 - H_0 : Defendant is innocent
 - H_a : Defendant is guilty
 - Correct: Reject H_0 when H_0 is false – Jury convicts a guilty person
 - Correct: Do not reject H_0 when H_0 is true – Jury acquits an innocent person
 - **Type I error**: Reject H_0 when H_0 is true – Jury convicts innocent person
 - **Type II error**: Do not reject H_0 when H_0 is false – Jury acquits a guilty person

		Reality	
		H_0 true: Defendant did not commit crime	H_0 false: Defendant did commit crime
Jury's decision	Reject H_0 : Find defendant guilty	Type I error	Correct decision
	Do not reject H_0 : Find defendant not guilty	Correct decision	Type II error

Hypothesis Testing for the Mean (*cont'd*)

- ▢ Probability of Type I error is denoted α

- ▢ Probability of Type II error is denoted β
 - ▢ For constant sample size, decrease in α is associated with increase in β and vice versa

- ▢ Common to restrict the hypotheses to these three forms
 - ▢ Left-tailed test;
 - ▢ Right-tailed test;
 - ▢ Two-tailed test;where μ_0 represents a hypothesized value of μ

Hypothesis Testing for the Mean (*cont'd*)

- When sample size is large or population is normally distributed, the test statistic

follows a t distribution, with degrees of freedom.

- The value of t is interpreted as the number of standard errors above/below the hypothesized mean μ_0 that the sample mean resides, where the standard error equals
- When t is extreme, this indicates **conflict** between null hypothesis (with μ_0) and the observed data
 - The null hypothesis H_0 is rejected
- How extreme is extreme? This is measured using the p -value

Hypothesis Testing for the Mean (*cont'd*)

- ▮ The *p-value* is the **probability** of observing a sample statistic as extreme as the statistic actually observed
- ▮ *p-value* is a probability, and its value falls between 0 and 1
- ▮ The names of the forms of the hypothesis test indicate in which tail or tails of the t distribution the p-value will be found
- ▮ We will reject if the p-value is small
- ▮ Researchers set the **level of significance** at some small value (such as 0.05); *p*-value is small if it is less than :

Reject if the p-value is

Form of Hypothesis Test	p-Value
Left-tailed test.	
Right-tailed test.	
Two-tailed test.	If then p-value = . If then p-value = .

Hypothesis Testing for the Mean (*cont'd*)

- Example: Test whether mean number of Customer Service Calls for customers with International Plan and VoiceMail Plan and who have more than 220 day minutes differs from 2.4, with a level of significance $\alpha=0.05$
- Null hypothesis is rejected if $p\text{-value} < 0.05$
- We also know that:
- Then
- Since , we have
 $p\text{-}$
- Since $p\text{-}$, **we reject H_0**
- At level of significance, there is evidence that the population mean number of Customer Service Calls of such customers differs from 2.4

Assessing the strength of evidence against the null hypothesis

- For the example with , and
 - We would not provide a conclusion
 - We would state that there is *solid evidence against the null hypothesis*

p-Value	Strength of evidence against
	Extremely strong evidence
	Very strong evidence
	Solid evidence
	Mild evidence
	Slight evidence
	No evidence

Using confidence intervals to perform hypothesis tests

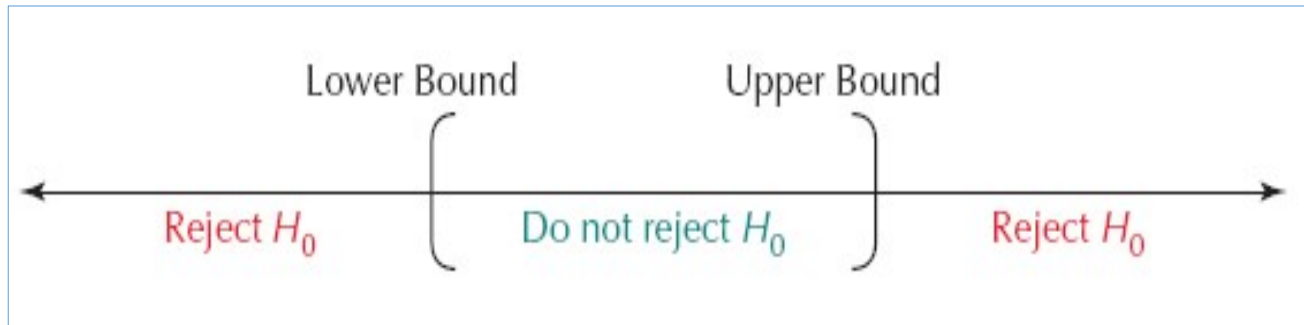
- One confidence interval is worth 1000 hypothesis tests
- The t confidence interval and the t hypothesis test are both based on the same distribution with the same assumptions

A confidence interval for μ is equivalent to a two-tailed hypothesis test for μ , with level of significance α .

Confidence level	Level of significance
90%	0.10
95%	0.05
99%	0.01

Using confidence intervals to perform hypothesis tests

- **Outside of the Confidence Interval** – then the two tailed hypothesis test will reject H_0 for that value of μ_0
- **Inside of the Confidence Interval** – then the two tailed hypothesis test will not reject H_0 for that value of μ_0



Using confidence intervals to perform hypothesis tests (*cont'd*)

- Example: The 95% confidence interval for the population mean of customer service calls was:

(lower bound, upper bound) = (0.875, 2.339)

- Perform two-tailed confidence interval test for the following values of μ_0

A. $\mu_0=0.5$ B. $\mu_0=1.0$ C. $\mu_0=2.4$

- The solution requires the following hypothesis tests

A.

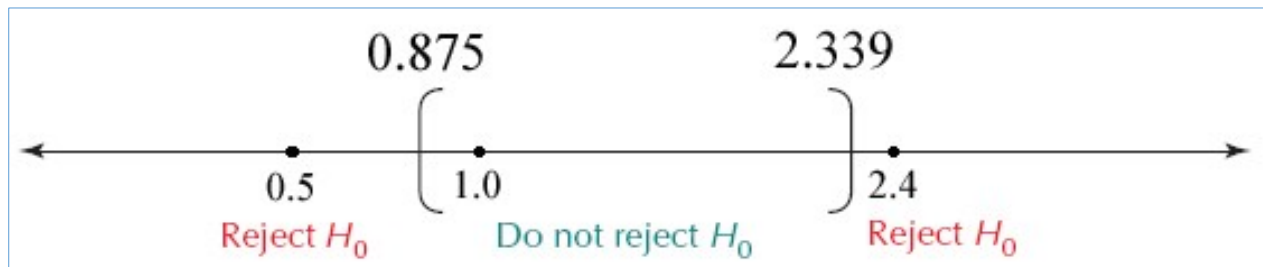
B.

C.

- Place the proposed values of μ_0 in the number line to determine their position with respect to the the 95% confidence interval of this variable (see Figure 4.4 in next slide)

Using confidence intervals to perform hypothesis tests (*cont'd*)

- Placement confirms that only $\mu_0=1.0$ is within the 95% confidence interval for the two-tailed hypothesis test with level of significance $\alpha=0.05$



	Hypotheses with	Position in relation to 95% confidence interval	Conclusion
0.5		Outside	Reject
1.0		Inside	Do not reject
2.4		Outside	Reject

Hypothesis testing for the proportion

□ Hypothesis test also applies to the population proportion π as:

where π_0 is the hypothesized value of π and p is the sample proportion

Hypotheses with	p-Value
Left-tailed test.	
Right-tailed test.	
Two-tailed test.	If then p-value = . If then p-value = .

Table of Standard Normal Probabilities for Negative Z-scores

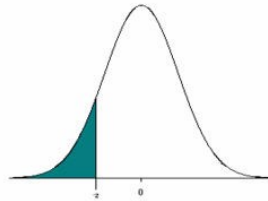
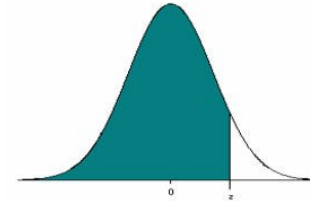


Table of Standard Normal Probabilities for Positive Z-scores



z	0.00	0.01	0.02	0.03	0.04	0.05	0.06	0.07	0.08	0.09
-3.4	0.0003	0.0003	0.0003	0.0003	0.0003	0.0003	0.0003	0.0003	0.0003	0.0002
-3.3	0.0005	0.0005	0.0005	0.0004	0.0004	0.0004	0.0004	0.0004	0.0004	0.0003
-3.2	0.0007	0.0007	0.0006	0.0006	0.0006	0.0006	0.0006	0.0005	0.0005	0.0005
-3.1	0.0010	0.0009	0.0009	0.0009	0.0008	0.0008	0.0008	0.0008	0.0007	0.0007
-3.0	0.0013	0.0013	0.0013	0.0012	0.0012	0.0011	0.0011	0.0011	0.0010	0.0010
-2.9	0.0019	0.0018	0.0018	0.0017	0.0016	0.0016	0.0015	0.0015	0.0014	0.0014
-2.8	0.0026	0.0025	0.0024	0.0023	0.0023	0.0022	0.0021	0.0021	0.0020	0.0019
-2.7	0.0035	0.0034	0.0033	0.0032	0.0031	0.0030	0.0029	0.0028	0.0027	0.0026
-2.6	0.0047	0.0045	0.0044	0.0043	0.0041	0.0040	0.0039	0.0038	0.0037	0.0036
-2.5	0.0062	0.0060	0.0059	0.0057	0.0055	0.0054	0.0052	0.0051	0.0049	0.0048
-2.4	0.0082	0.0080	0.0078	0.0075	0.0073	0.0071	0.0069	0.0068	0.0066	0.0064
-2.3	0.0107	0.0104	0.0102	0.0099	0.0096	0.0094	0.0091	0.0089	0.0087	0.0084
-2.2	0.0139	0.0136	0.0132	0.0129	0.0125	0.0122	0.0119	0.0116	0.0113	0.0110
-2.1	0.0179	0.0174	0.0170	0.0166	0.0162	0.0158	0.0154	0.0150	0.0146	0.0143
-2.0	0.0228	0.0222	0.0217	0.0212	0.0207	0.0202	0.0197	0.0192	0.0188	0.0183
-1.9	0.0287	0.0281	0.0274	0.0268	0.0262	0.0256	0.0250	0.0244	0.0239	0.0233
-1.8	0.0359	0.0351	0.0344	0.0336	0.0329	0.0322	0.0314	0.0307	0.0301	0.0294
-1.7	0.0446	0.0436	0.0427	0.0418	0.0409	0.0401	0.0392	0.0384	0.0375	0.0367
-1.6	0.0548	0.0537	0.0526	0.0516	0.0505	0.0495	0.0485	0.0475	0.0465	0.0455
-1.5	0.0668	0.0655	0.0643	0.0630	0.0618	0.0606	0.0594	0.0582	0.0571	0.0559
-1.4	0.0808	0.0793	0.0778	0.0764	0.0749	0.0735	0.0721	0.0708	0.0694	0.0681
-1.3	0.0968	0.0951	0.0934	0.0918	0.0901	0.0885	0.0869	0.0853	0.0838	0.0823
-1.2	0.1151	0.1131	0.1112	0.1093	0.1075	0.1056	0.1038	0.1020	0.1003	0.0985
-1.1	0.1357	0.1335	0.1314	0.1292	0.1271	0.1251	0.1230	0.1210	0.1190	0.1170
-1.0	0.1587	0.1562	0.1539	0.1515	0.1492	0.1469	0.1446	0.1423	0.1401	0.1379
-0.9	0.1841	0.1814	0.1788	0.1762	0.1736	0.1711	0.1685	0.1660	0.1635	0.1611
-0.8	0.2119	0.2090	0.2061	0.2033	0.2005	0.1977	0.1949	0.1922	0.1894	0.1867
-0.7	0.2420	0.2389	0.2358	0.2327	0.2296	0.2266	0.2236	0.2206	0.2177	0.2148
-0.6	0.2743	0.2709	0.2676	0.2643	0.2611	0.2578	0.2546	0.2514	0.2483	0.2451
-0.5	0.3085	0.3050	0.3015	0.2981	0.2946	0.2912	0.2877	0.2843	0.2810	0.2776
-0.4	0.3446	0.3409	0.3372	0.3336	0.3300	0.3264	0.3228	0.3192	0.3156	0.3121
-0.3	0.3821	0.3783	0.3745	0.3707	0.3669	0.3632	0.3594	0.3557	0.3520	0.3483
-0.2	0.4207	0.4168	0.4129	0.4090	0.4052	0.4013	0.3974	0.3936	0.3897	0.3859
-0.1	0.4602	0.4562	0.4522	0.4483	0.4443	0.4404	0.4364	0.4325	0.4286	0.4247
-0.0	0.5000	0.4960	0.4920	0.4880	0.4840	0.4801	0.4761	0.4721	0.4681	0.4641

z	0.00	0.01	0.02	0.03	0.04	0.05	0.06	0.07	0.08	0.09
0.0	0.5000	0.5040	0.5080	0.5120	0.5160	0.5199	0.5239	0.5279	0.5319	0.5359
0.1	0.5398	0.5438	0.5478	0.5517	0.5557	0.5596	0.5636	0.5675	0.5714	0.5753
0.2	0.5793	0.5832	0.5871	0.5910	0.5948	0.5987	0.6026	0.6064	0.6103	0.6141
0.3	0.6179	0.6217	0.6255	0.6293	0.6331	0.6368	0.6406	0.6443	0.6480	0.6517
0.4	0.6554	0.6591	0.6628	0.6664	0.6700	0.6736	0.6772	0.6808	0.6844	0.6879
0.5	0.6915	0.6950	0.6985	0.7019	0.7054	0.7088	0.7123	0.7157	0.7190	0.7224
0.6	0.7257	0.7291	0.7324	0.7357	0.7389	0.7422	0.7454	0.7486	0.7517	0.7549
0.7	0.7580	0.7611	0.7642	0.7673	0.7704	0.7734	0.7764	0.7794	0.7823	0.7852
0.8	0.7881	0.7910	0.7939	0.7967	0.7995	0.8023	0.8051	0.8078	0.8106	0.8133
0.9	0.8159	0.8186	0.8212	0.8238	0.8264	0.8289	0.8315	0.8340	0.8365	0.8389
1.0	0.8413	0.8438	0.8461	0.8485	0.8508	0.8531	0.8554	0.8577	0.8599	0.8621
1.1	0.8643	0.8665	0.8686	0.8708	0.8729	0.8749	0.8770	0.8790	0.8810	0.8830
1.2	0.8849	0.8869	0.8888	0.8907	0.8925	0.8944	0.8962	0.8980	0.8997	0.9015
1.3	0.9032	0.9049	0.9066	0.9082	0.9099	0.9115	0.9131	0.9147	0.9162	0.9177
1.4	0.9192	0.9207	0.9222	0.9236	0.9251	0.9265	0.9279	0.9292	0.9306	0.9319
1.5	0.9332	0.9345	0.9357	0.9370	0.9382	0.9394	0.9406	0.9418	0.9429	0.9441
1.6	0.9452	0.9463	0.9474	0.9484	0.9495	0.9505	0.9515	0.9525	0.9535	0.9545
1.7	0.9554	0.9564	0.9573	0.9582	0.9591	0.9599	0.9608	0.9616	0.9625	0.9633
1.8	0.9641	0.9649	0.9656	0.9664	0.9671	0.9678	0.9686	0.9693	0.9699	0.9706
1.9	0.9713	0.9719	0.9726	0.9732	0.9738	0.9744	0.9750	0.9756	0.9761	0.9767
2.0	0.9772	0.9778	0.9783	0.9788	0.9793	0.9798	0.9803	0.9808	0.9812	0.9817
2.1	0.9821	0.9826	0.9830	0.9834	0.9838	0.9842	0.9846	0.9850	0.9854	0.9857
2.2	0.9861	0.9864	0.9868	0.9871	0.9875	0.9878	0.9881	0.9884	0.9887	0.9890
2.3	0.9893	0.9896	0.9898	0.9901	0.9904	0.9906	0.9909	0.9911	0.9913	0.9916
2.4	0.9918	0.9920	0.9922	0.9925	0.9927	0.9929	0.9931	0.9932	0.9934	0.9936
2.5	0.9938	0.9940	0.9941	0.9943	0.9945	0.9946	0.9948	0.9949	0.9951	0.9952
2.6	0.9953	0.9955	0.9956	0.9957	0.9959	0.9960	0.9961	0.9962	0.9963	0.9964
2.7	0.9965	0.9966	0.9967	0.9968	0.9969	0.9970	0.9971	0.9972	0.9973	0.9974
2.8	0.9974	0.9975	0.9976	0.9977	0.9977	0.9978	0.9979	0.9979	0.9980	0.9981
2.9	0.9981	0.9982	0.9982	0.9983	0.9984	0.9984	0.9985	0.9985	0.9986	0.9986
3.0	0.9987	0.9987	0.9987	0.9988	0.9988	0.9989	0.9989	0.9989	0.9990	0.9990
3.1	0.9990	0.9991	0.9991	0.9991	0.9992	0.9992	0.9992	0.9992	0.9993	0.9993
3.2	0.9993	0.9993	0.9994	0.9994	0.9994	0.9994	0.9994	0.9995	0.9995	0.9995
3.3	0.9995	0.9995	0.9995	0.9996	0.9996	0.9996	0.9996	0.9996	0.9996	0.9997
3.4	0.9997	0.9997	0.9997	0.9997	0.9997	0.9997	0.9997	0.9997	0.9997	0.9998

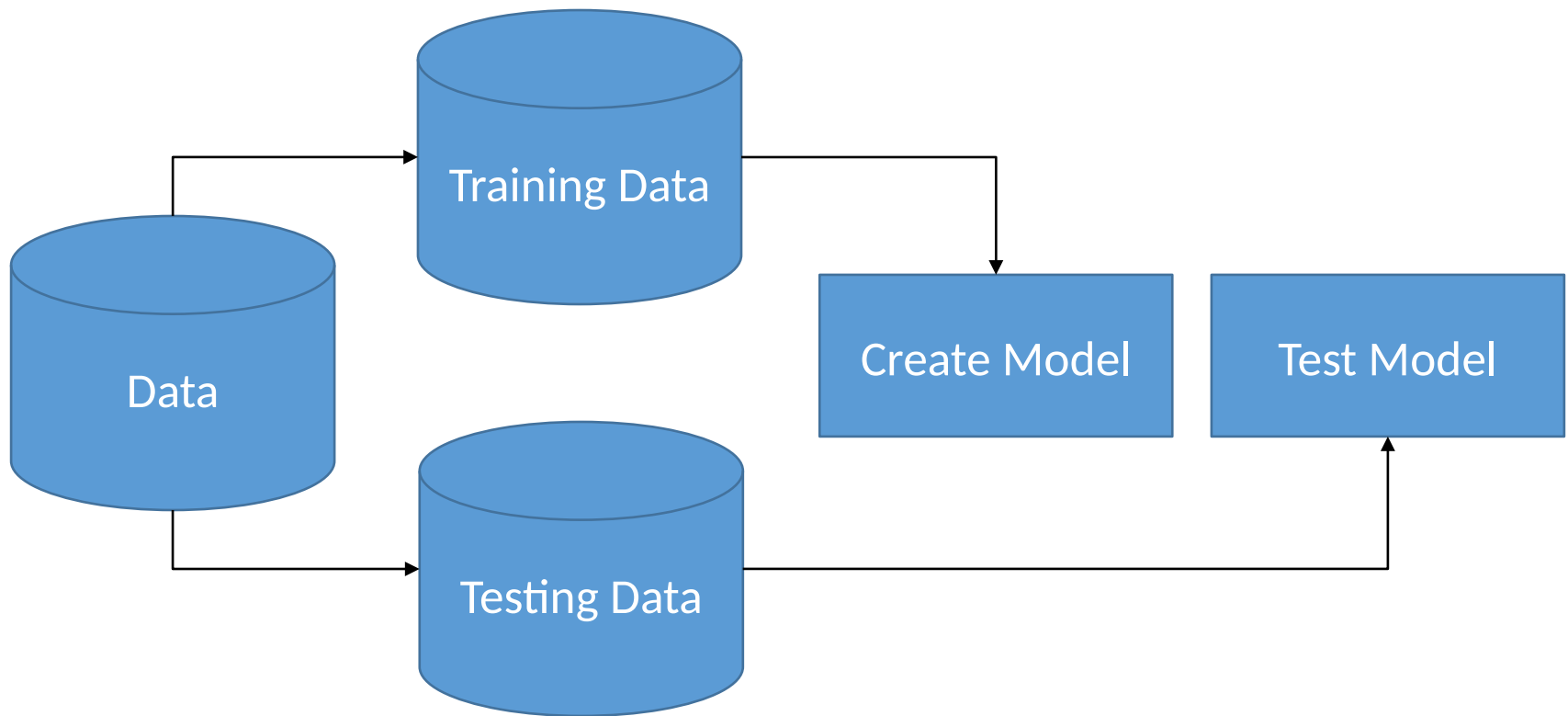
**Note that the probabilities given in this table represent the area to the LEFT of the z-score.
The area to the RIGHT of a z-score = 1 – the area to the LEFT of the z-score**

Hypothesis testing for the proportion

- Example: As 483 of 3333 customer churned, the estimate of the population proportion π of all the customers that churned is:
- Testing at $\alpha=0.10$ whether π differs from 0.15, the hypotheses are:
- The test statistics is:
- As , the p -=
- Since the p -value is not less than $\alpha=0.10$ we would not reject H_0

Multivariate Analysis

Supervised Learning



Multivariate Analysis

- Important when splitting dataset into training and test data sets
 - Bivariate hypothesis tests shown here can be used to determine whether significant differences exist between the means of various variables in the **training and test data sets**
 - If such differences exists, training set is **not representative** of the test set
 - For a continuous variable, use the **two-sample t test** for the difference in means.
 - For a flag variable, use the **two-sample Z test** for the difference in proportions.
 - For a multinomial variable, use the **test for the homogeneity** of proportions.
- Even if the data set has more than two variables, spot-checking of a few randomly chosen variables is usually sufficient.

Two-sample t Test for difference in means

□ The test statistic for the difference in population mean is:

which follows an approximate t distribution with degrees of freedom the smaller of $n_1 - 1$ and $n_2 - 1$ when both populations are normally distributed or both samples are large

□ Example: We divided the churn dataset into a training and a test data set

□ [Assess the validity](#) of the partition by testing whether the population mean number of customer service calls differs between the two data sets

Data Set	Sample Mean	Sample Standard Deviation	Sample Size
Training Set	1.5714	1.3126	2529
Test Set	1.5361	1.3251	804

Two-sample t Test for difference in means (*cont'd*)

- ▮ Need to perform hypothesis test to make sure the validity
- ▮ Hypothesis is:
- ▮ The test statistic is:
- ▮ The two-tailed *p-value* for is:
- ▮ *p-value* is large
- ▮ There is **no evidence** that mean number of customer service calls differs between test and training data sets
- ▮ For this variable, the partition seems **valid**

Two-sample Z Test for difference in proportions

- Not all variables are numeric/continuous
- For a **flag** variable (like 1/0) we need the two-sample Z test for the difference in proportions

where n_i and p_i represents the number of and proportion of records with value 1 (for example) for sample i , respectively

Two-sample Z Test for difference in proportions (*cont'd*)

- Example: The training partition in the previous example resulted in $x_1=707$ of $n_1=2529$ customers belonging to Voice Mail Plan, while the Test set has $x_2=215$ of $n_2=804$
- ,
-
-
- The hypotheses are:
- The test statistic is:
- The *p-value* is:
- There is **no evidence** that the proportion of Voice Mail Plan differs between the training and test data sets.
- For this variable, the partition is **valid**

Test for the homogeneity of proportions

- ▮ *Multinomial data* is an extension of binomial data to $k > 2$ categories
 - ▮ Example multinomial variable: *marital status* can be *married*, *single*, *other*
 - ▮ Training set of 1000 people and test set of 250 people
- ▮ Test for the homogeneity of proportions
 - ▮ To determine whether significant differences exist between multimodal proportions
- ▮ Hypotheses are:

Data Set	Married	Single	Other	Total
Training set	410	340	250	1000
Test set	95	85	70	250
Total	505	425	320	1250

Test for the homogeneity of proportions (*cont'd*)

- Compare **observed frequencies** against **expected frequencies** if H_0 were true
- Example:
 - A. Find overall proportion of married people in whole dataset (training+test sets):
 - B. Multiply this overall proportion by the number of people in training set, 1000, yields the expected proportion of married people in the training set to be:
- Step A above uses the overall proportion because H_0 states that both partitions are equal

Test for the homogeneity of proportions (*cont'd*)

□ Generalizing, the expected frequencies are calculated as follows:

□ *Observed frequencies (O)* and *expected frequencies (E)* are compared using test statistics from the (chi-square distribution:

Data Set	Married	Single	Other	Total
Training set	404	340	256	1000
Test set	101	85	64	250
Total	505	425	320	1250

Test for the homogeneity of proportions (*cont'd*)

- Large differences between observed and expected frequencies, and large value for χ^2 , leads to small p-value, and rejection of null hypothesis
- The *p-value* is the **area** to the right of χ^2 under the χ^2 curve with degrees of freedom equal to $(\text{number of rows} - 1)(\text{number of columns} - 1) = (1)(2) = 2$:
- As p-value is large, there is no evidence that the frequencies significantly differ between the training and the test data sets
- The partition is **valid**

Chi-square test for goodness of fit of multinomial data

- ▮ Assume that marital status of the population is married=40%, single=35%, other=25%
- ▮ Determine whether the sample is representative of the population
- ▮ Use χ^2 (chi-square) goodness of fit test. Hypotheses are:
- ▮ Sample size $n=100$ yields the following observed frequencies
- ▮ Need to compare observed frequencies against the expected frequencies assuming that H_0 is true

Chi-square test for goodness of fit of multinomial data (*cont'd*)

- Comparing the frequencies using the test statistic

Cell	Observed Frequency	Expected Frequency	
Married	36	40	
Single	35	35	
Other	29	25	
			=1.04

Chi-square test for goodness of fit of multinomial data (*cont'd*)

- The p-value is the area to the right of χ^2 under the χ^2 curve with $k-1$ degrees of freedom, where k -number of categories (here $k=3$)
- There is no evidence that the observed frequencies represent proportions that differ from those in the null hypothesis
- Our sample is representative of the population

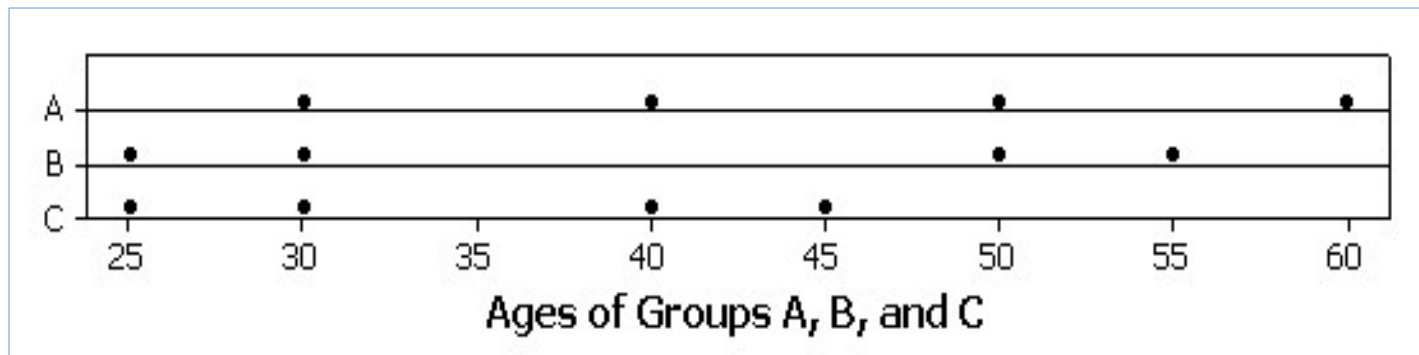
Analysis of Variance

- ▮ One-way analysis of **variance** (ANOVA)
 - ▮ Used as an extension of the situation for the two-sample t test
 - ▮ For example: Check whether mean value of a continuous variable is the same across three partitions of a data set
 - ▮ For the samples for Age variable from Groups A, B and C, we set the following hypotheses
-
- ▮ Sample mean ages are , , and

Group A	Group B	Group C
30	25	25
40	30	30
50	50	40
60	55	45

Analysis of Variance (*cont'd*)

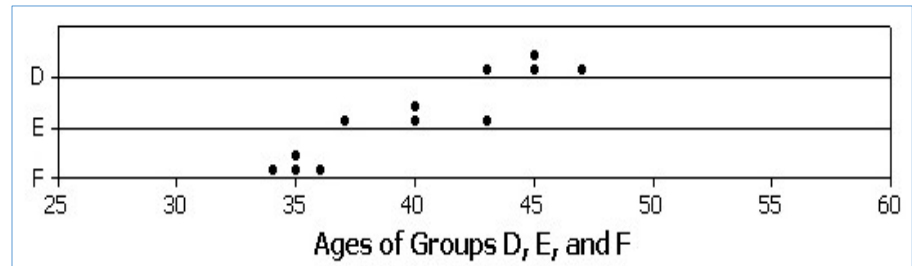
- Dot plot of the data shows considerable amount of overlap among the groups
- Dot plot offers little evidence to reject the null hypotheses that the population means are all equal



Analysis of Variance (*cont'd*)

- Consider now the groups D, E and F
- Sample means are again , , and
- But the dot plot shows little overlap among data sets
 - There is good evidence to reject the null hypotheses

Group D	Group E	Group F
43	37	34
45	40	35
45	40	35
47	43	36



Analysis of Variance (*cont'd*)

- When spread within each sample is large, the difference in sample means seems small
- When the spread within each sample is small, the difference in sample means seems large
- Analysis of variance compares:
 1. Between-sample variability (measured by the variability of the sample means)
 2. Within-sample variability (measured, for example, by the sample standard deviation)
- When #1 is much larger than #2, this represents evidence that the population means are not equal.
- The analysis depends on measuring variability, hence the name *analysis of variance*

Analysis of Variance (*cont'd*)

- Let $\bar{\bar{y}}$ represent the overall sample mean
- Measure the **between-sample variability** by finding the variance of the k sample means, weighted by sample size
 - Expressed as the *mean square treatment* (MSTR):
- Measure the **within-sample variability** by finding the weighted mean of the sample variances
 - Expressed as the *mean square error* (MSE):
- We compare these two quantities by taking their ratio:

which follows an F distribution, with degrees of freedom $k-1$ and $n-k$

Analysis of Variance (*cont'd*)

- The numerator of MSTR is the *sum of squares treatment*, SSTR
- The numerator of MSE is the *sum of squares error*, SSE.
- The total sum of squares (SST) is the sum of SSTR and SSE
- ANOVA table is a convenient way of displaying this information
- ❖ F_{data} will be large when the between-sample variability is much greater than the within-sample variability – this calls for **rejection** of null hypothesis
- The p-value is ; reject the null hypothesis when the p-value is small

Source of Variation	Sum of Squares	Degrees of Freedom	Mean Square	F
Treatment	SSTR			
Error	SSE			
Total	SST			

Analysis of Variance (*cont'd*)

- Example: ANOVA results for Groups A, B, and C/Groups D, E, and F
- The p-value=0.548 indicates that there is **no evidence** against the null hypothesis that all population means are equal
- The p-value=0.000 indicates that there is **strong evidence** that not all population means ages are equal

Group A	Group B	Group C
30	25	25
40	30	30
50	50	40
60	55	45

Source	DF	SS	MS	F	P
Factor	2	200	100	0.64	0.548
Error	9	1400	156		
Total	11	1600			

Group D	Group E	Group F
43	37	34
45	40	35
45	40	35
47	43	36

Source	DF	SS	MS	F	P
Factor	2	200.00	100.00	32.14	0.000
Error	9	28.00	3.11		
Total	11	228.00			

Regression Analysis

□ Cereals Dataset

Cereals dataset variables	
Cereal manufacturer	Grams of sugar
Type (hot or cold)	Milligrams of potassium
Calories per serving	Percentage of recommended daily allowance of vitamins (0% 25%, or 100%)
Grams of protein	Weight of one serving
Grams of fat	Number of cups per serving
Milligrams of sodium	Shelf location (1 = bottom, 2 = middle, 3 = top)
Grams of fiber	Nutritional rating, calculated by Consumer Reports
Grams of carbohydrates	Grams of carbohydrates in a cereal

Regression Analysis (*cont'd*)

□ Caveat: This dataset contains missing data:

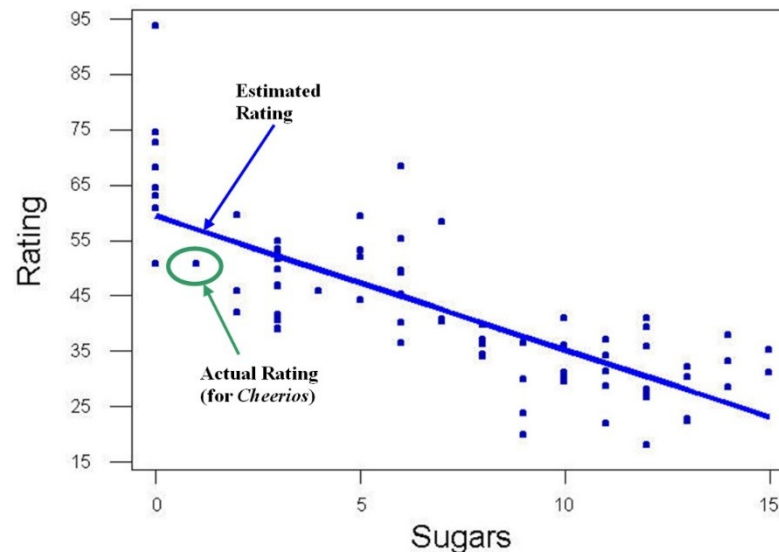
- Potassium content of Almond Delight
- Potassium content of Cream of Wheat
- Carbohydrates
- Sugars content of Quaker Oatmeal

□ Will not be able to use sugar content in Quaker Oatmeal to build model

Cereal Name	Manuf.	Sugars	Calories	Protein	Fat	Sodium	Rating
100% Bran	N	6	70	4	1	130	68.4030
100% Natural Bran	Q	8	120	3	5	15	33.9837
All-Bran	K	5	70	4	1	260	59.4255
All-Bran Extra Fiber	K	0	50	4	0	140	93.7049
Almond Delight	R	8	110	2	2	200	34.3848
Apple Cinnamon Cheerios	G	10	110	2	2	180	29.5095

Regression Analysis (*cont'd*)

- Scatter plot of the nutritional rating vs sugar content, along with the least-squares regression line
- The regression equation is , where:
 - is the estimated value of the response variable
 - is the y-intercept of the regression line
 - is the slope of the regression line
 - and , together, are called the regression coefficients



Regression Analysis (cont'd)

- The top regression equation $\text{Rating} = 59.9 - 2.46 \text{ Sugars}$ shows rounded coefficients. More digits in the “Coef” column
 - $b_0 = 59.852$
 - $b_1 = -2.4614$
- Thus,
- Interpreted as “*The estimated cereal rating equals 59.853 minus 2.4614 times the sugar content in grams*”
- The regression line and equations are linear approximations of the relationship between x (**predictor**) and y (**response**) variables – sugar and nutritional rating in this case
- Regression equation can be used to make estimates and predictions

The regression equation is
 $\text{Rating} = 59.9 - 2.46 \text{ Sugars}$

76 cases used, 1 cases contain missing values

Predictor	Coef	SE Coef	T	P
Constant	59.853	1.998	29.96	0.000
Sugars	-2.4614	0.2417	-10.18	0.000

S = 9.16616 R-Sq = 58.4% R-Sq(adj) = 57.8%

Analysis of Variance

Source	DF	SS	MS	F	P
Regression	1	8711.9	8711.9	103.69	0.000
Residual Error	74	6217.4	84.0		
Total	75	14929.3			

Unusual Observations

Obs	Sugars	Rating	Fit	SE Fit	Residual	St Resid
1	6.0	68.40	45.08	1.08	23.32	2.56R
4	0.0	93.70	59.85	2.00	33.85	3.78R

R denotes an observation with a large standardized residual.

Predicted Values for New Observations

New Obs	Fit	SE Fit	95% CI	95% PI
1	57.39	1.80	(53.81, 60.97)	(38.78, 76.00)

Values of Predictors for New Observations

New Obs	Sugars
1	1.00

Regression Analysis (*cont'd*)

- Example: Estimate the nutritional rating for new cereal that contains $x = 1$ gram of sugar
 - Using the regression equation:
 - Estimated value lies directly on the regression line, at $(x=1,)$
 - For any value of x (sugar content) the estimated value for y (nutritional rating) lies precisely in the regression line
- Cheerios cereal also has sugar=1 gram, but rating is 50.765, not the estimated 57.3916. Point $(x=1,)$
- Prediction using regression line **deviated** by $50.765 - 57.3916 = -6.6266$ rating points
 - This vertical distance in general $()$ is known as **prediction error**, **estimation error**, or **residual**
- Two unusual observations: cereal 1 (100% Bran) and cereal 4 (All-Bran with Extra Fiber), which have large positive residuals, indicating nutrition rating unexpectedly high, given their sugar level

Regression Analysis (*cont'd*)

- Regression objective is to **minimize the overall size of the prediction errors**
 - Least-squares regression works by choosing the regression that minimizes the sum of squared errors (SSE) – this is the most common method
- y-intercept b_0 is the location on the y-axis where the regression line intercepts the y-axis; the value when the predictor is zero
 - For the cereal data set, b_0 represents the nutritional content of a cereal with no sugar
- The slope of the regression line indicates the change in y per unit increase in x
 - $b_1 = -2.4614$ means, “For each increase of 1 gram in sugar content, the estimated nutritional rating decreases by 2.4614 rating points”
 - For example: For cereal A with 5 more grams of sugar than cereal B would have estimated nutritional rating of $5(2.4614) = 12.307$ rating points lower than cereal B.

Hypothesis Testing in Regression

- Objective: to use the known value of the slope b_1 to perform inference for the unknown value of the slope β_1 of the population regression equation
 - Just like when in univariate analysis we used sample mean to infer the unknown value of the population mean μ
- The population regression equation represents the relationship between, say, nutritional ratings and sugar content for the entire population of cereals – not just the ones in the sample

where ε represents a random variable for modeling the errors

- Notice that when $\beta_1 = 0$, the population equations becomes and there is **no relationship** between the predictor x and the response y
- For any other value of β_1 there is a linear relationship between x and y

Hypothesis Testing in Regression (*cont'd*)

- We wish to test for the existence of a linear relationship between x and y
- Thus, perform the following hypothesis test

No relationship between x and y
Linear relationship between x and y

- The test statistic for this hypothesis test is:

where represents the standard error for

- Large values of indicated too much variability in the slope of b_1 – this makes precise inference difficult
- Per the formula above, large values of reduce the size of the t -statistic

Measuring the Quality of a Regression Model

- If we cannot reject the null hypothesis that , the regression is not useful
- If we find that , then there are two statistics to measure the quality of the regression
 - Standard error of the estimate: s
 - R-squared: r^2

Measuring the Quality of a Regression Model

- Standard error of the estimate: s
 - Should not be confused with s , the sample standard deviation for univariate statistics, or the standard error of the slope coefficient
 - Defined as:

where SSE is the sum of squared errors
 - Its value indicates the size of the “typical” prediction error

Measuring the Quality of a Regression Model

□ R-squared statistic: r^2

- Measures how closely the linear regression fits the data
- Values near 100% indicate a more perfect fit
- Defined as:

SST: the variability in the y-variable

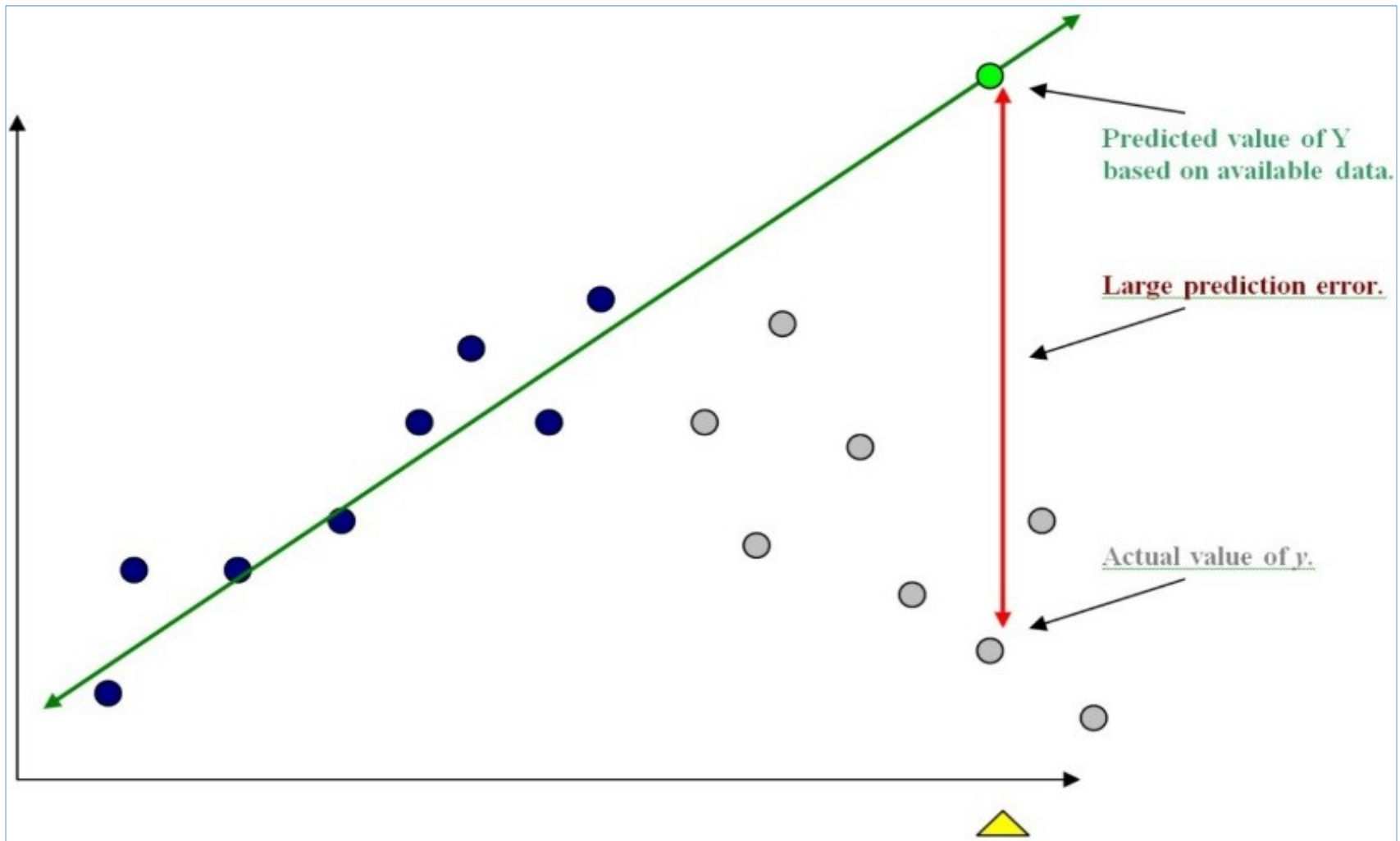
SSR: the improvement in the estimation as compared to just using

- Interpreted as the ratio of total variability in y that is accounted for by the linear relationship between x and y

Dangers of extrapolation

- Using the estimated regression on a new cereal with 30 grams of sugar per serving yields:
- The minimum rating in the dataset is 18, but this new cereal has a negative rating?
- This is an example of the dangers of extrapolation
- Analyst should restrict estimates and predictions to the values within the **range of the values of x** in dataset
 - Example: The range of sugar content in the dataset is from 0 to 15
 - Predictions of nutritional rating for cereals with 0-15 grams of sugar be appropriate
 - Prediction outside of this range would be dangerous, since we don't know the nature of the relationship outside of this range
- Extrapolation should be avoided, and end user should be informed that no x-data is available to support a prediction
- Relationship outside of this range may no longer be linear

Dangers of extrapolation (*cont'd*)



Confidence intervals for the mean value of y given x

- Point estimates suffer from the lack of a probability statement associated with their accuracy
- Use confidence intervals for the mean of y for a given x :

x_p = a particular x for which the prediction is being made

= the point estimate of y for a particular x

= a multiplier associated with the sample size and confidence level

s = the standard error of the estimate

SSE = the sum of squared residuals

Predictions for a randomly chosen value of y given x

- It is 'easier' to predict a mean value than a randomly chosen value of a variable
 - Example: predicting the school class exam score average vs predicting a specific student exam score
- Data miners are **more interested in predicting an individual value**, rather than mean of all the values, given x
- **Prediction intervals** are used to estimate the value of a randomly chosen value of y, given x; as follows:
 - This formula is the same as the formula for the confidence interval of the mean value of y, given x, **except for the presence of the "1+" inside the square root**
 - The "1+" ensures a wider prediction interval than the analogous confidence interval

Multiple Regression

- Some data sets include hundreds of variables, which may have a linear relationship with the target variable
- Multiple regression modeling provides a method for describing those relationships
- Example: For the *Cereal* dataset, add *sodium* (in addition to existing variable *sugar*) for predicting the *rating*, and observe whether the quality of the model improved or not
- The equation for multiple regression with two predictors is:

The regression equation is

$$\text{Rating} = 69.2 - 2.39 \text{ Sugars} - 0.0606 \text{ Sodium}$$

76 cases used, 1 cases contain missing values

Predictor	Coef	SE Coef	T	P
Constant	69.180	2.373	29.15	0.000
Sugars	-2.3944	0.2041	-11.73	0.000
Sodium	-0.06057	0.01086	-5.58	0.000

S = 7.72769 R-Sq = 70.8% R-Sq(adj) = 70.0%

Analysis of Variance

Source	DF	SS	MS	F	P
Regression	2	10569.9	5285.0	88.50	0.000
Residual Error	73	4359.4	59.7		
Total	75	14929.3			

Source	DF	Seq SS
Sugars	1	8711.9
Sodium	1	1858.0

Unusual Observations

Obs	Sugars	Rating	Fit	SE Fit	Residual	St Resid
1	6.0	68.403	46.940	0.970	21.463	2.80R
2	8.0	33.984	49.116	1.845	-15.133	-2.02R
3	5.0	59.426	41.461	1.465	17.964	2.37R
4	0.0	93.705	60.701	1.691	33.004	4.38R

R denotes an observation with a large standardized residual.

Predicted Values for New Observations

New Obs	Fit	SE Fit	95% CI	95% PI
1	49.222	2.107	(45.023, 53.421)	(33.258, 65.185)

Values of Predictors for New Observations

New Obs	Sugars	Sodium
1	1.00	290

Multiple Regression (*cont'd*)

□ Multiple regression equation is:

where $b_2 = -0.06057$ is interpreted so that for each additional milligram of sodium, the estimated decrease in nutritional rating is 0.06057, ***when sugars is held constant***

□ The point estimate of rating for Cheerios, with 1 gram of sugar and 290 mg of sodium is:

□ The prediction error for Cheerios is the difference between actual rating y and the predicted rating : $(y-) = 50.765 - 49.22 = 1.545$

□ This prediction error is smaller than the obtained using sugars only:

□ This, because our model uses double the data (two predictors instead of one)

□ The standard error of the estimate has been reduced from $s \approx 9.2$ to $s \approx 7.7$

□ The value for r^2 has increased from 58.4% to 70.8%

□ 70.8% of our variability in nutritional rating is explained by our regression model

Verifying model assumptions

- Before implementing a model, the prerequisite model assumptions must be verified
- Making predictions using a model where assumptions are violated may lead to erroneous results
- The assumptions are: *Linearity, Independence, Normality and Constant Variance*
- Assumptions might be checked with:
 - Normality plot of residuals (Figure 5.9, upper left)
 - Checks whether there are systematic deviations from linearity
 - In there are deviations, the values do not follow a normal distribution
 - Plot of standardized residuals against the fitted (predicted) values
 - If obvious curvature exists in the scatter plot, the linearity assumption is violated
 - If the vertical spread is systematically non-uniform, the constant variance assumption is violated

Verifying model assumptions (*cont'd*)

□ Important

- We must be careful **not to find patterns where no such pattern exists!**
- Departure from linearity must be systematic and significant
- The huge data sets used in data mining do not always follow perfect normality
 - so that analyst should give the data the benefit of the doubt unless good evidence exists for the contrary

□ Outliers have an **outsized influence** in Least squares regression

- Should identify **extreme** outliers, and *if necessary and appropriate*, omit them
- <http://r-statistics.co/Linear-Regression.html>
- <https://jakevdp.github.io/PythonDataScienceHandbook/05.06-linear-regression.html>
- https://scikit-learn.org/stable/modules/generated/sklearn.linear_model.LinearRegression.html