

STAT 200 Final Exam Study Guide Questions for Chapters 42-56

About 80% of the Final covers Chapters 42 - 56 and is both Multiple Choice and calculated answers (similar to Exam 1 and 2).

Formulas to Know:

Confidence Intervals for Transformed Variables (asymmetrical CI's)

3 forms of logistic regression model: ln(odds), odds, probability

Odds, and OR

Z and Chi square tests

Rank sums and U for Wilcoxon Mann Whitney, Z test

Rank sums for Kruskal Wallis, Chi square test

Spearman r, Z test

Only 3 Formulas that will be given to you:

$$SE_{R_A} = SE_{R_B} = SE_U = \sqrt{\frac{n_A n_B (N+1)}{12}}$$

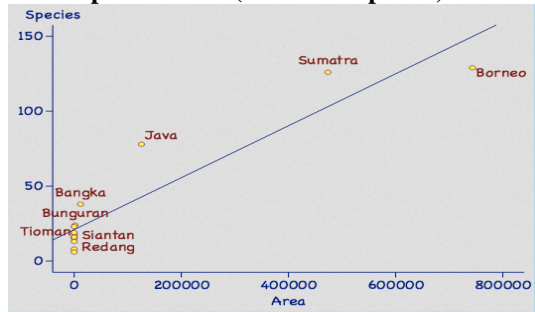
$$H = \frac{12}{N(N+1)} \sum_{i=1}^g \frac{(\text{obs}R_i - \text{exp}R_i)^2}{n_i}$$

$$SE_{r_s} = \frac{1}{\sqrt{n-1}}$$

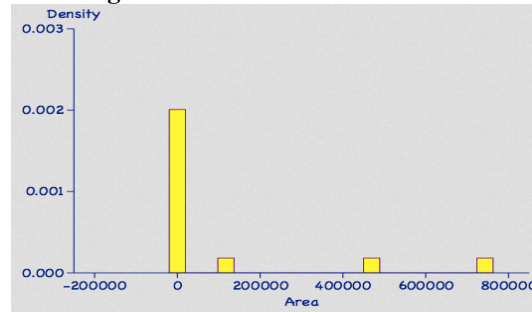
STAT 200 Final Exam Study Guide Questions

Question 1 pertains to the **Area** (in km²) and the **number of mammal species** for 13 islands in Southeast Asia. How does the size of the island predict the number of species on the island?

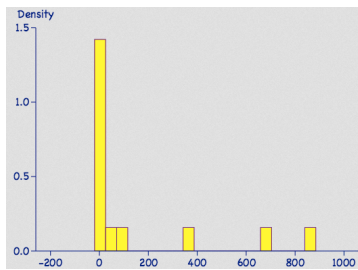
Scatter plot of Area (in km² vs Species)



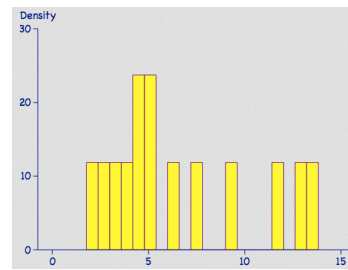
Histogram of Area



- a) Notice how most of the islands are all squished together in the corner. Also look how skewed the Area histogram is. I want to transform the X variable (Area) to make the histogram more normal. Which transformations should I try? Circle ALL that might work. i) X^2 ii) X^3 iii) e^X iv) \sqrt{X} v) $\ln(X)$
- b) You tried one of the transformations and it was a step in the right direction but it didn't go far enough. You tried another and it worked much better well. Below each histogram circle the transformation it depicts.

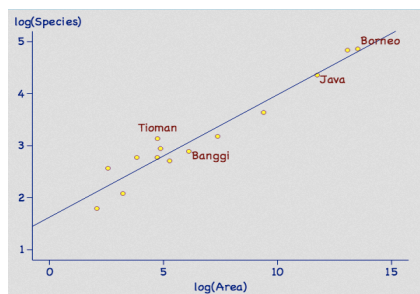


- i) X^2 ii) X^3 iii) e^X iv) \sqrt{X} v) $\ln(X)$



- i) X^2 ii) X^3 iii) e^X iv) \sqrt{X} v) $\ln(X)$

c) Below is the scatter plot of $\ln(\text{Species})$ vs $\ln(\text{Area})$ where Species= the number of mammal species on each island and Area= area of each island in km² The regression equation is: **Predicted $\ln(\text{Species}) = 1.6 + 0.23 \ln(\text{Area})$** $SD_{\text{errors}} = 0.2$



i) Banggai has an area= 450 km². Use the regression equation to predict the **$\ln(\text{Species})$** and **Species** number for Banggai

- a) $\ln(\text{Species}) =$ _____ b) Number of species= _____

c) 95% Confidence Interval for part(b) above= (_____, _____)
(Use $Z=2$ for 95% CI)

ii) Another island has a 95% confidence interval = **(11.23, 25)** for the predicted number of species. What is the predicted number of species? _____

iii) Change the regression equation **$\ln(\text{Species}) = 1.6 + 0.23 \ln(\text{Area})$** to an equation in terms of **species** and **Area**, not **$\ln(\text{Area})$** .

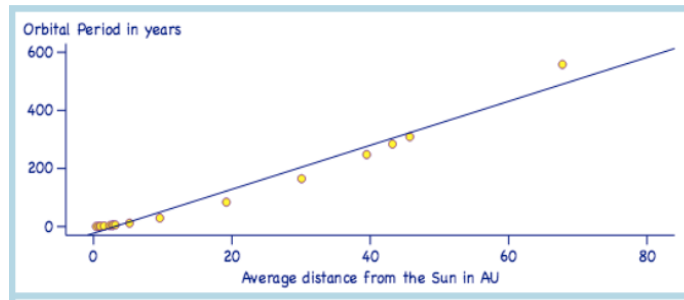
Species = _____

iv) One island has twice the area of another island. The regression estimate for the number of species on the smaller island is 9. What is the regression estimate for the number of species on the larger island? _____

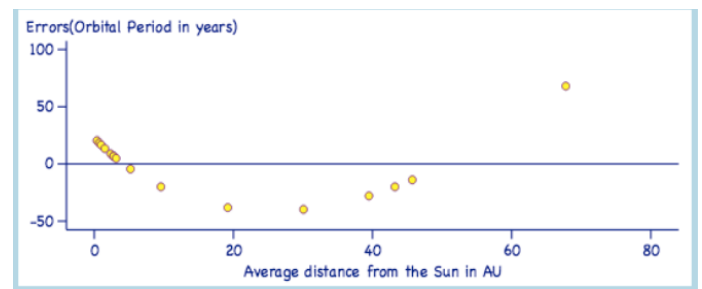
STAT 200 Final Exam Study Guide Questions

Question 2 The scatter plot below shows the average distance from the Sun in AU (astronomical units) on the X axis and the Orbital period in years (length of time to orbit sun) on the Y axis of 16 solar systems objects. (Imagine these 16 objects were randomly chosen from a large collection of objects orbiting the sun.)

Scatter Plot



Residual Plot



Here's the regression equation: **Predicted Orbital Period = -23.12 + 7.57(Distance from Sun) $r = 0.9864$ and $SD_{\text{errors}} = 26.04$**

- a) Why do the 16 points closely follow a line in the scatter plot but follow a curve in the residual plot?
- Residual plots always transform linear plots into curves that either point up or down depending on the whether the correlation is positive or negative.
 - It's because the correlation is so high, the higher the correlation the stronger the curvature.
 - It's because the scale on the Y axis for the residual plot has been changed, making it easier to see the curvature.
- b) Is it appropriate to use the regression equation above to describe the relation between distance from the sun and orbital period for all the objects ?
- Yes, because the scatter plot follows a line very closely.
 - No, because the residual plot shows a clear pattern violating the assumptions needed to use a linear model.
 - Yes, because the 16 objects were randomly selected so there is no need to check whether assumptions were violated.

Question 3

For each of the following is it appropriate to use logistic regression? **Circle Yes or No.**

- Predicting income based on years of college. YES NO
- Predicting $\ln(\text{income})$ based on years of college YES NO
- Predicting graduating college based on family income. YES NO
- Predicting getting a scholarship based on gender and ethnicity. YES NO
- Predicting favorite color based on gender YES NO

Question 4 Circle True or False for each statement below.

- The logistic regression model only handles X values that can be coded as 1's and 0's. i)True ii)False
- Transforming non-linear scatter plots into linear ones by converting Y to $\ln(Y)$ is called logistic regression. i)True ii)False
- The assumptions needed to make inferences for linear and logistic regression are the same i)True ii)False

Question 5

How are the parameters chosen in logistic regression and linear regression?

Fill in the first blank below with "logistic" or "linear" and the second blank with "minimize" or "maximize".

- In _____ regression, the parameters are chosen to _____ the sum of the squared errors
- In _____ regression, the parameters are chosen to _____ the likelihood of getting our sample data.

STAT 200 Final Exam Study Guide Questions

Question 6 Are F and t tests ever appropriate to test significance in Logistic regression models?

Choose one:

- a) Yes, when the sample size is small the F and t tests give more accurate results.
- b) No, because F and t tests can never be done on variables that have undergone log transformations.
- c) No, because F and t tests are never done when we are predicting counts (when Y is binary), since the SD can be estimated directly from the count.

Question 7 Part I On our survey, 178 students anonymously answered these 2 questions:

“Would you volunteer to be randomly assigned to either the online or in person section?”(No = 0, Yes =1)

“Which section are you in?” (L1=0, online=1)

To predict the probability of volunteering from section, we fit a logistic regression model. Here’s the ln(odds) form of the

$$\text{regression equation: } \ln\left(\frac{\hat{p}}{1-\hat{p}}\right) = -0.5261 + -0.7267(\text{Section})$$

- a) Are online students more or less likely to volunteer? **Choose one:** i) More ii) Less iii) Same iv) Not enough info
- b) What is the probability that an L1 student would volunteer? p=_____
- c) What is the probability that an online student would volunteer? p=_____
- d) The Odds Ratio = _____.
- e) If we switched the coding for section to online = 0 and L1 =1 what would change? **Choose one:**
i) Odds ii) Probabilities iii) Odds Ratio iv) All v) None
- f) Look at the table showing the 178 responses to the 2 questions.

Use the table to compute the odds for an L1 and online student volunteering. **Please leave your answers in fraction form.**

	No	Yes	Total
L1	44	26	70
Online	84	24	108
Total	128	50	178

i) Odds for L1 = _____

ii) Odds for Online = _____

iii) Should you get the same OR as in (d) above? (Assuming you compute the ratio of Online odds to L1 odds.)

- a) Yes, within rounding error
- b) No

STAT 200 Final Exam Study Guide Questions

Question 7 Part II A third question on the same survey was: “How many people have you been in a serious relationship with?” Adding relationships to the model gives us: $\ln\left(\frac{\hat{p}}{1-\hat{p}}\right) = -1.33 + -1.03(\text{Section}) + 0.64(\text{Relationships})$

- a) The χ^2 test for the overall regression effect: H_0 : All β 's = 0 yielded a χ^2 stat = 26.
How many **degrees of freedom?** = _____
- b) The p value < 0.1%. This means that the probability that ... **Choose only one:**
 i) the null is true < 0.1% ii) the null is false > 99.9% iii) we'd get a χ^2 stat ≥ 26 if the null was true < 0.1%
- c) The relationship slope has a SE = 0.14. To test $H_0: \beta_{\text{relationship}} = 0$ against $H_A: \beta_{\text{relationship}} \neq 0$ compute the Z stat.

$$Z = \underline{\hspace{2cm}}$$

- d) Since p _____5%, a 95% Confidence interval for the Relationship slope _____include _____.
Fill in the first blank with > or < , the second with “does” or “does not”, and the third blank with a number.
- e) The OR for Relationship = _____ and the OR for Section = _____
- f) Comparing two people in the same section, the person with 2 more relationships has _____ times the odds of volunteering. **Fill in the blank with a number.**
- g) Comparing an L1 student with 4 relationships to an online student with 2 relationships, the L1 student has _____ times the odds of volunteering. **Fill in the blank with a number.**
- h) What's the probability that an L1 student with 10 relationships will volunteer? _____.
- i) Would the $\ln(\text{odds})$ equation for Part II change if we reversed the coding for Section so that L1=1 and online=0 and kept everything else the same? If so, write the new equation in the blank provided.
- a) No, it would not change. b) Yes, it would change to $\ln\left(\frac{\hat{p}}{1-\hat{p}}\right) = \underline{\hspace{4cm}}$

STAT 200 Final Exam Study Guide Questions

Question 8 A predictor of whether esophageal cancer has not metastasized to the lymph nodes is the diameter of the tumor. Below is the log odds regression equation predicting the probability of no metastasis from the diameter of the tumor (measured in cm) from a hypothetical study of 200 patients.

$$\ln(p/(1-p)) = 2 - 0.5 (\text{Diameter})$$

- a) Use the equation to estimate the **odds** and **probability** of no metastasis for a tumor of diameter = 8 cm. *Show work.*

i) Odds= _____

ii) Probability = _____

- b) How do the estimated *odds* of no metastasis change if the tumor increases in diameter by 1 cm ?

i) odds are multiplied by 0.61 ii) the odds decrease by 0.5 iii) not enough info

- c) How does the estimated *probability* of no metastasis change if the tumor increases in diameter by 1 cm?

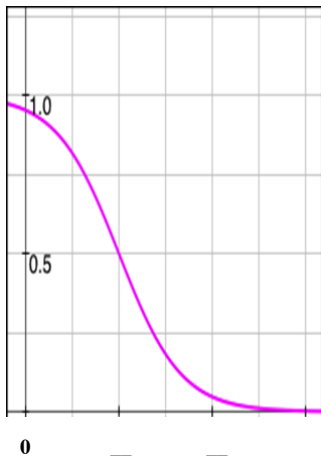
i) the probability is multiplied by 0.61 ii) the probability decreases by 0.5 iii) not enough info

- d) How big a tumor would give a 50% probability of metastasis? _____

- e) How big a tumor would give a 40% probability of no metastasis? _____

- f) Below is a graph of the probability form of the model.

Write its equation: $p =$ _____ and fill in the 2 blanks on the X-axis with the correct diameter values (in cm).



Fill in the 2 blanks above with the correct numbers.

STAT 200 Final Exam Study Guide Questions

Question 9 pertains to the Wilcoxon Mann Whitney test

A randomized double-blind test was done to test the effectiveness of a drug to cure warts. The subjects were 8 people with lots of warts. 4 subjects took the drug and 4 took the placebo. The number of warts that disappeared for each of the 8 subjects is recorded below.

Drug Group: 0, 10, 11, 40 Placebo group: 5, 6, 8, 9

Part 1

Fill out the chart below. Show work for how you got the observed rank sum for each group.

No partial credit since you should know what the totals should be and you can check your work.

	Observed Rank Sum	Expected Rank Sum	Observed - Expected
Drug Group			
Placebo Group			
Total should be....			

Question 9 Part II

The sample sizes in Part I are too small to use the Normal Approximation but let's just assume for the purpose of this exam that you can use the Normal Approximation anyway.

H_0 : The drug works no better than the placebo in the population

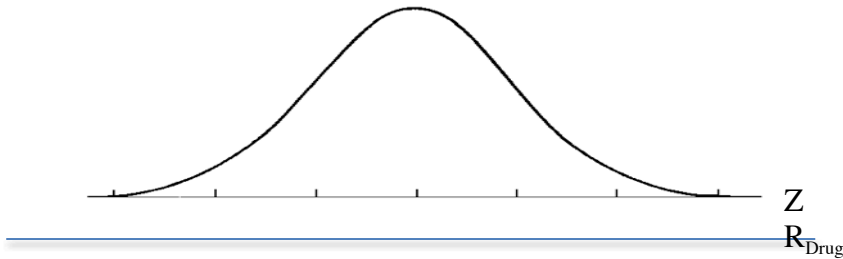
H_A : The drug *does work better* than the placebo in the population for some segments of the population.

a) Compute the Z stat for the drug group.

$$\text{Use } SE_R = \sqrt{\frac{n_1 n_2 (N + 1)}{12}}$$

Z=

b) Label the Observed and Expected Value for both the Z and R_{Drug} axes below. Calculate the p-value and shade the appropriate area. (R_{Drug} is the Rank Sum for the Drug group.)



c) What do you conclude? (Remember, we're assuming the sample size was large enough so the normal approximation is valid).

- i) Reject the null, we're sure the drug works.
- ii) Reject the null, we have strong evidence the drug works.
- iii) Cannot reject the null, it's plausible the drug works no better than a placebo.
- iv) There's over a 95% chance the drug didn't work.

STAT 200 Final Exam Study Guide Questions

Question 9 cont.

Drug Group: 0, 10, 11, 40 Placebo group: 5, 6, 8, 9

d) What's the U statistic for the Drug Group? For the Placebo group?

$U_{\text{drug}} =$

$U_{\text{placebo}} =$

e) The sum of the 2 group U statistics must = _____ for any 2 groups with 4 members each.
(Check that your $U_{\text{drug}} + U_{\text{placebo}}$ is correct.)

f) Would you get the same Z stat and p-value using U_{drug} as you did using R_{drug} in part (a)?

i) Yes, exactly the same.

ii) Exactly the same values but the Z-scores would be opposite signs.

iii) No, the p-value would be smaller using U.

iv) No, the p-value would be larger using U.

Question 10 pertains to the Kruskal Wallis test (6 pts)

There are 3 forms of this Final. Suppose at the grading meeting I randomly select 9 Finals and grade them with these results:

Form A: 80, 81, 82

Form B: 83, 84, 85

Form C: 86, 87, 89

Null Hypothesis: No difference in difficulty of the exams in the population. We just happen to observe differences in our sample due to chance variation.

Alternative Hypothesis: At least one of the exams is of different difficulty in the population.

a) The Rank Sum for Form A= _____, Form B=_____ and Form C=_____

b) The total Rank Sum for any set of 9 numbers is always= _____. (give a number.)

a) The H-stat = 7.2 Would any other arrangement of 9 numbers into 3 groups of 3 yield a higher H-stat?

i) No

ii) Yes

iii) Not enough info

b) For large enough samples we can best approximate the distribution of the H stat with

i) Z stat

ii) t stat

iii) Chi-square stat

iv) the F stat

STAT 200 Final Exam Study Guide Questions

Question 11

a) If we decide to do a non-parametric test and use the Spearman correlation coefficient to test the null hypothesis that the population correlation is 0 then the appropriate test-statistic for small samples (<7) is ...

- i) a t-statistic
- ii) Spearman correlation tables that calculate the exact probability distribution
- iii) 2 sample t-statistic
- iv) an F-test
- v) a Chi Square test

b) For large enough samples the appropriate test statistic is

- i) Z-test
- ii) t-test
- iii) either
- iv) F-test
- v) none of the above

Question 12

Look at the 3 data sets below:

Data Set 1: (1,2), (2, 4) , (3, 6) , (4,8)

Data Set 2: (-1, 5) (-2, 4) (-3, 3)

Data Set 3: (1,1) (8,9) (103,10)

For which data set(s) is $r \neq r_s$?

STAT 200 Final Exam Study Guide Questions

Critical Values for F distribution at $p = 5\%$ and $p = 1\%$

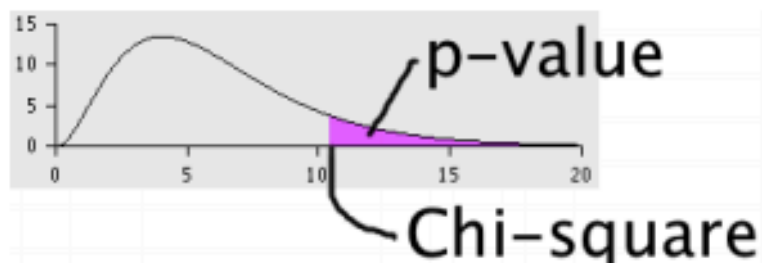
F Distribution critical values for $P=0.05$

F Denominator														
	Numerator DF													
DF	1	2	3	4	5	7	10	15	20	30	60	120	500	1000
1	161.45	199.50	215.71	224.58	230.16	236.77	241.88	245.95	248.01	250.10	252.20	253.25	254.06	254.19
2	18.513	19.000	19.164	19.247	19.296	19.353	19.396	19.429	19.446	19.462	19.479	19.487	19.494	19.495
3	10.128	9.5522	9.2766	9.1172	9.0135	8.8867	8.7855	8.7028	8.6602	8.6165	8.5720	8.5493	8.5320	8.5292
4	7.7086	6.9443	6.5915	6.3882	6.2560	6.0942	5.9644	5.8579	5.8026	5.7458	5.6877	5.6580	5.6352	5.6317
5	6.6078	5.7862	5.4095	5.1922	5.0504	4.8759	4.7351	4.6187	4.5582	4.4958	4.4314	4.3985	4.3731	4.3691
7	5.5914	4.7375	4.3469	4.1202	3.9715	3.7871	3.6366	3.5108	3.4445	3.3758	3.3043	3.2675	3.2388	3.2344
10	4.9645	4.1028	3.7082	3.4780	3.3259	3.1354	2.9782	2.8450	2.7741	2.6996	2.6210	2.5801	2.5482	2.5430
15	4.5431	3.6823	3.2874	3.0556	2.9013	2.7066	2.5437	2.4035	2.3275	2.2467	2.1601	2.1141	2.0776	2.0718
20	4.3512	3.4928	3.0983	2.8660	2.7109	2.5140	2.3479	2.2032	2.1241	2.0391	1.9463	1.8962	1.8563	1.8498

F Distribution critical values for $P=0.01$

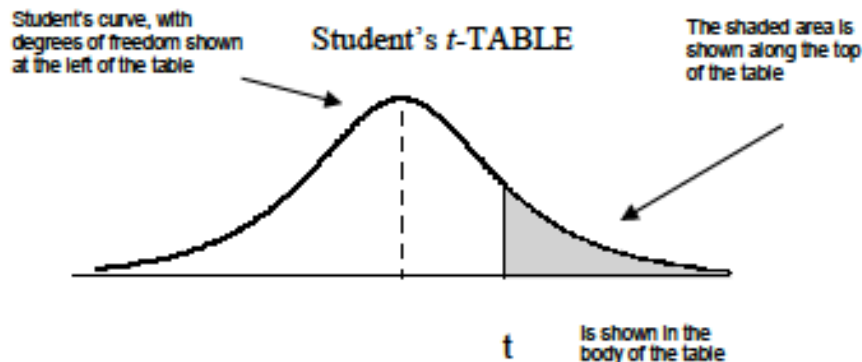
F Denominator														
	Numerator DF													
DF	1	2	3	4	5	7	10	15	20	30	60	120	500	1000
1	4052.2	4999.5	5403.4	5624.6	5763.6	5928.4	6055.8	6157.3	6208.7	6260.6	6313.0	6339.4	6359.5	6362.7
2	98.503	99.000	99.166	99.249	99.299	99.356	99.399	99.433	99.449	99.466	99.482	99.491	99.497	99.498
3	34.116	30.817	29.457	28.710	28.237	27.672	27.229	26.872	26.690	26.504	26.316	26.221	26.148	26.137
4	21.198	18.000	16.694	15.977	15.522	14.976	14.546	14.198	14.020	13.838	13.652	13.558	13.486	13.474
5	16.258	13.274	12.060	11.392	10.967	10.455	10.051	9.7222	9.5526	9.3793	9.2020	9.1118	9.0424	9.0314
7	12.246	9.5467	8.4513	7.8466	7.4605	6.9929	6.6201	6.3143	6.1554	5.9920	5.8236	5.7373	5.6707	5.6601
10	10.044	7.5594	6.5523	5.9944	5.6363	5.2001	4.8492	4.5582	4.4055	4.2469	4.0818	3.9964	3.9303	3.9195
15	8.6831	6.3588	5.4169	4.8932	4.5557	4.1416	3.8049	3.5223	3.3719	3.2141	3.0471	2.9594	2.8906	2.8796
20	8.0960	5.8489	4.9382	4.4306	4.1027	3.6987	3.3682	3.0880	2.9377	2.7785	2.6078	2.5167	2.4446	2.4330
30	7.5624	5.3903	4.5098	4.0179	3.6990	3.3046	2.9791	2.7002	2.5486	2.3859	2.2078	2.1108	2.0321	2.0192
60	7.0771	4.9774	4.1259	3.6491	3.3388	2.9530	2.6318	2.3522	2.1978	2.0284	1.8362	1.7264	1.6328	1.6169
120	6.8509	4.7865	3.9490	3.4795	3.1736	2.7918	2.4720	2.1914	2.0345	1.8600	1.6557	1.5330	1.4215	1.4015
500	6.6858	4.6479	3.8210	3.3569	3.0539	2.6751	2.3564	2.0746	1.9152	1.7353	1.5175	1.3774	1.2317	1.2007
1000	6.6603	4.6264	3.8012	3.3379	3.0356	2.6571	2.3387	2.0564	1.8967	1.7158	1.4953	1.3513	1.1947	1.1586

Chi-square table



Degrees of freedom ↓	30%	10%	5%	1%	0.1%	← p-value
1	1.07	2.71	3.84	6.63	10.83	
2	2.41	4.61	5.99	9.21	13.82	
3	3.66	6.25	7.81	11.34	16.27	
4	4.88	7.78	9.49	13.28	18.47	
5	6.06	9.24	11.07	15.09	20.52	
6	7.23	10.64	12.59	16.81	22.46	
7	8.38	12.02	14.07	18.48	24.32	
8	9.52	13.36	15.51	20.09	26.12	
9	10.66	14.68	16.92	21.67	27.88	
10	11.78	15.99	18.31	23.21	29.59	
11	12.90	17.28	19.68	24.72	31.26	
12	14.01	18.55	21.03	26.22	32.91	← Chi-square
13	15.12	19.81	22.36	27.69	34.53	
14	16.22	21.06	23.68	29.14	36.12	
15	17.32	22.31	25.00	30.58	37.70	
16	18.42	23.54	26.30	32.00	39.25	
17	19.51	24.77	27.59	33.41	40.79	
18	20.60	25.99	28.87	34.81	42.31	
19	21.69	27.20	30.14	36.19	43.82	
20	22.77	28.41	31.41	37.57	45.31	
21	23.86	29.62	32.67	38.93	46.80	
22	24.94	30.81	33.92	40.29	48.27	
23	26.02	32.01	35.17	41.64	49.73	
24	27.10	33.20	36.42	42.98	51.18	

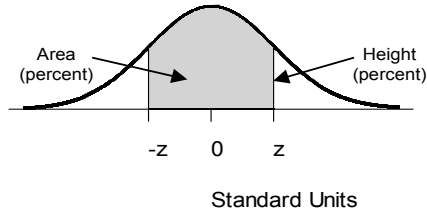
STAT 200 Final Exam Study Guide Questions



<i>Degrees of freedom</i>	25%	10%	5%	2.5%	1%	0.5%
1	1.00	3.08	6.31	12.71	31.82	63.66
2	0.82	1.89	2.92	4.30	6.96	9.92
3	0.76	1.64	2.35	3.18	4.54	5.84
4	0.74	1.53	2.13	2.78	3.75	4.60
5	0.73	1.48	2.02	2.57	3.36	4.03
6	0.72	1.44	1.94	2.45	3.14	3.71
7	0.71	1.41	1.89	2.36	3.00	3.50
8	0.71	1.40	1.86	2.31	2.90	3.36
9	0.70	1.38	1.83	2.26	2.82	3.25
10	0.70	1.37	1.81	2.23	2.76	3.17
11	0.70	1.36	1.80	2.20	2.72	3.11
12	0.70	1.36	1.78	2.18	2.68	3.05
13	0.69	1.35	1.77	2.16	2.65	3.01
14	0.69	1.35	1.76	2.14	2.62	2.98
15	0.69	1.34	1.75	2.13	2.60	2.95
16	0.69	1.34	1.75	2.12	2.58	2.92
17	0.69	1.33	1.74	2.11	2.57	2.90
18	0.69	1.33	1.73	2.10	2.55	2.88
19	0.69	1.33	1.73	2.09	2.54	2.86
20	0.69	1.33	1.72	2.09	2.53	2.85
21	0.69	1.32	1.72	2.08	2.52	2.83
22	0.69	1.32	1.72	2.07	2.51	2.82
23	0.69	1.32	1.71	2.07	2.50	2.81
24	0.68	1.32	1.71	2.06	2.49	2.80
25	0.68	1.32	1.71	2.06	2.49	2.79

STAT 200 Final Exam Study Guide Questions

STANDARD NORMAL TABLE



<i>z</i>	<i>Area</i>		<i>z</i>	<i>Area</i>		<i>z</i>	<i>Area</i>
0.00	0.00		1.50	86.64		3.00	99.730
0.05	3.99		1.55	87.89		3.05	99.771
0.10	7.97		1.60	89.04		3.10	99.806
0.15	11.92		1.65	90.11		3.15	99.837
0.20	15.85		1.70	91.09		3.20	99.863
0.25	19.74		1.75	91.99		3.25	99.885
0.30	23.58		1.80	92.81		3.30	99.903
0.35	27.37		1.85	93.57		3.35	99.919
0.40	31.08		1.90	94.26		3.40	99.933
0.45	34.73		1.95	94.88		3.45	99.944
0.50	38.29		2.00	95.45		3.50	99.953
0.55	41.77		2.05	95.96		3.55	99.961
0.60	45.15		2.10	96.43		3.60	99.968
0.65	48.43		2.15	96.84		3.65	99.974
0.70	51.61		2.20	97.22		3.70	99.978
0.75	54.67		2.25	97.56		3.75	99.982
0.80	57.63		2.30	97.86		3.80	99.986
0.85	60.47		2.35	98.12		3.85	99.988
0.90	63.19		2.40	98.36		3.90	99.990
0.95	65.79		2.45	98.57		3.95	99.992
1.00	68.27		2.50	98.76		4.00	99.9937
1.05	70.63		2.55	98.92		4.05	99.9949
1.10	72.87		2.60	99.07		4.10	99.9959
1.15	74.99		2.65	99.20		4.15	99.9967
1.20	76.99		2.70	99.31		4.20	99.9973
1.25	78.87		2.75	99.40		4.25	99.9979
1.30	80.64		2.80	99.49		4.30	99.9983
1.35	82.30		2.85	99.56		4.35	99.9986
1.40	83.85		2.90	99.63		4.40	99.9989
1.45	85.29		2.95	99.68		4.45	99.9991