



**THE UNIVERSITY  
OF QUEENSLAND**  
AUSTRALIA

This exam paper must not be removed from the venue

Venue \_\_\_\_\_

Seat Number \_\_\_\_\_

Student Number

--	--	--	--	--	--	--	--	--	--

Family Name \_\_\_\_\_

First Name \_\_\_\_\_

## School of Mathematics & Physics

### EXAMINATION

Semester Two Final Examinations, 2017

### STAT1201 Analysis of Scientific Data

*This paper is for St Lucia Campus, Gatton Campus (External) and Gatton Campus students.*

Examination Duration: 120 minutes

Reading Time: 10 minutes

#### Exam Conditions:

This is a Central Examination

This is a Closed Book Examination - specified materials permitted

During reading time - Write only on rough paper provided

This examination paper will be released to the Library

#### Materials Permitted In The Exam Venue:

**(No electronic aids are permitted e.g. laptops, phones)**

Calculators - Casio FX82 series or UQ approved (labelled)

An annotated copy of *A Portable Introduction to Data Analysis* (any edition) is also permitted.

#### Materials To Be Supplied To Students:

None

#### Instructions To Students:

There are **52** marks available on this exam from **4** questions, each worth **13** marks.

Write your answers in the spaces provided in pages 2–15 of this examination paper. Show your working and state conclusions where appropriate.

Pages 16–20 give formulas and statistical tables. Those pages will not be marked.

The textbook can have any amount of annotation on its pages. Loose sheets of paper or post-it notes are not permitted. Page tabs are allowed.

#### For Examiner Use Only

Question Mark

1	
2	
3	
4	

Total \_\_\_\_\_

## Question 1

A study aimed to determine if grapefruit juice has beneficial effects on the pharmacokinetics of oral digoxin, a drug often prescribed for heart ailments. Seven healthy non-smoking volunteers participated in the study. Subjects took digoxin with water for 2 weeks, no digoxin for 2 weeks, and then digoxin with grapefruit juice for 2 weeks. The peak plasma digoxin concentrations ( $C_{max}$ ; ng/mL) when subjects took digoxin under the two conditions are given in the following table:

Subject	1	2	3	4	5	6	7
Water	2.34	2.46	1.87	3.09	5.59	4.05	6.21
Grapefruit Juice	3.03	3.46	1.97	3.81	3.07	2.62	3.44
Decrease	-0.69	-1.00	-0.10	-0.72	2.52	1.43	2.77

While small, note that the sample size was chosen carefully by the authors. In their paper they state that “assuming an  $\alpha$  level of 0.05, a sample size of seven subjects has a power of 85% to detect a 25% change in digoxin  $C_{max}$ ”.

Lower values of  $C_{max}$  are better since they imply that digoxin is available in the body for longer. Is there any evidence that grapefruit juice increases the effectiveness of oral digoxin, by decreasing  $C_{max}$ ?

- (a) Identify one issue in the design of this experiment that undermines being able to use the data to answer this question. How could the design be improved? [2 marks]

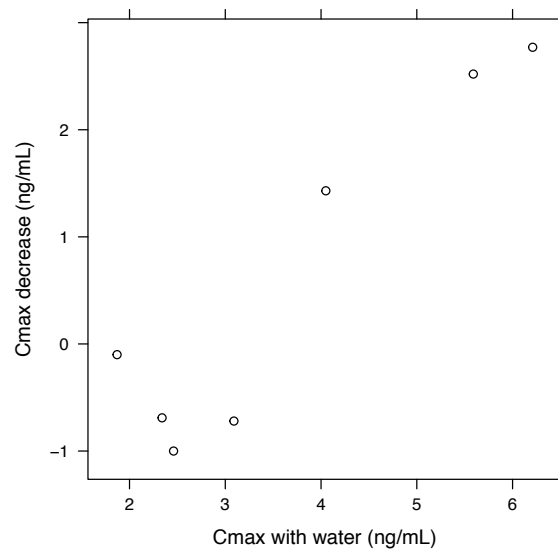
- (b) Suppose  $\mu$  is the true mean decrease in digoxin  $C_{max}$  with grapefruit compared to water. Define the null and alternative hypotheses for this study in symbols. [1 mark]

$H_0$ :

$H_1$ :

- (c) We have two sets of  $C_{\max}$  measurements, one for water and one for grapefruit. Briefly explain why we work with the differences rather than carrying out a two-sample  $t$ -test to compare the treatments. [1 mark]
- (d) The seven differences in  $C_{\max}$  have mean 0.601 ng/mL with standard deviation 1.609 ng/mL. Use these values to test the hypotheses in (b). What do you conclude? [2 marks]
- (e) How many of the seven subjects had a lower  $C_{\max}$  value with grapefruit juice? Use this to find the  $P$ -value for a sign test of whether grapefruit juice tends to lower  $C_{\max}$ . What do you conclude? [2 marks]
- (f) For a sign test from seven subjects, what is the minimum number of reductions in  $C_{\max}$  needed to give evidence of an effect at the 5% level? [1 mark]

- (g) The authors speculated that the decrease in  $C_{\max}$  with the grapefruit juice may be related to an individual's baseline response with water. They produced the following scatterplot of the relationship observed in the data:



They calculated the Pearson correlation coefficient for this relationship to be 0.9321. Does this give any evidence that there is an underlying association between the decrease in  $C_{\max}$  and the  $C_{\max}$  with water? [3 marks]

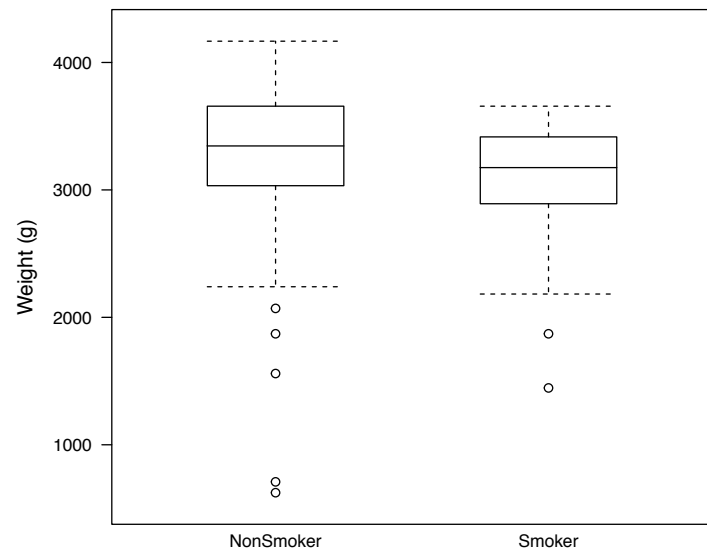
- (h) Briefly explain how your conclusion to (g) relates to your conclusion to (d). [1 mark]

## Question 2

In a study of factors thought to be associated with birth weight, a random sample of 100 births was selected from all the birth records at a hospital in 2001. Three variables were extracted from each record: the length of gestation (weeks), the smoking status of the mother (smoker/nonsmoker) and the birth weight (grams).

- (a) The mean length of gestation for the 100 births was 38.36 weeks with standard deviation 3.36 weeks. Construct a 95% confidence interval for the mean length of gestation for all births at this hospital in 2001. [2 marks]
- (b) There are claims that the mean length of gestation has been decreasing in recent decades from the traditionally held value of 40 weeks. Based on your interval in (a), or otherwise, do the data from this hospital support or contradict such claims? [1 mark]

The main question of interest is whether babies born to smoking mothers tend to have lower birth weights. The following figure shows a side-by-side comparison of the two distributions from this data set.



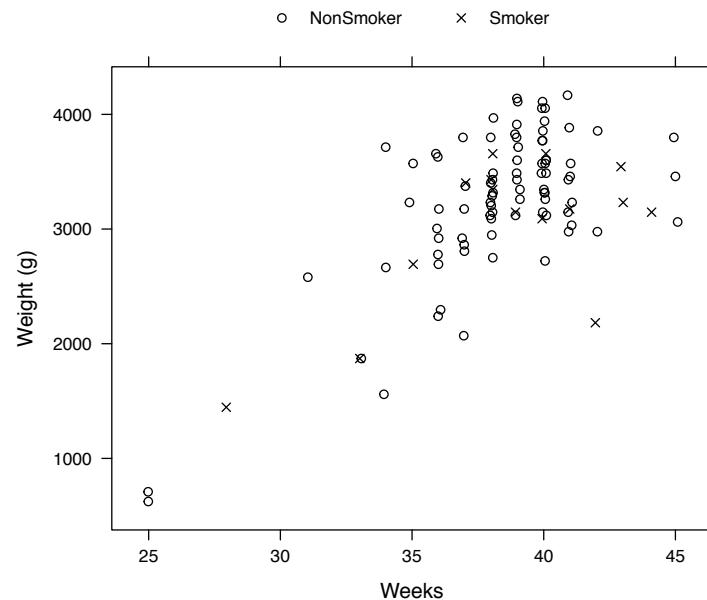
A two-sample t-test in R gave the following output:

```
Welch Two Sample t-test

data:  Weight by SmokingStatus
t = 1.3865, df = 18.42, p-value = 0.1817
alternative hypothesis: true difference in means is not equal to 0
95 percent confidence interval:
 -131.3122  646.6612
sample estimates:
mean in group NonSmoker    mean in group Smoker
      3258.941              3001.267
```

- (c) On average, how much lower are birth weights for smokers compared to non-smokers? [1 mark]
- (d) Summarise the conclusions from the t-test. [2 marks]

There is considerable variability in birth weights for both groups in the above figure. However, some of this variability may be explained by the different gestation lengths for each birth. The following scatterplot shows the relationship between birth weight and gestation length, with smokers and non-smokers shown by the plotting symbols:



The R output below gives the summary of the multiple regression model for birth weight based on both gestation length and smoking status:

```
lm(formula = Weight ~ Weeks + SmokingStatus, data = births)

Coefficients:
              Estimate Std. Error t value Pr(>|t|)
(Intercept)   -1724.42    558.84  -3.086  0.00265 **
Weeks           130.05     14.52   8.957 2.39e-14 ***
SmokingStatusSmoker -294.40    135.78  -2.168  0.03260 *
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 484.6 on 97 degrees of freedom
Multiple R-squared:  0.4636, Adjusted R-squared:  0.4525
F-statistic: 41.92 on 2 and 97 DF, p-value: 7.594e-14
```

- (e) Based on the model output, what is the estimated birth weight for a birth at 35 weeks gestation to a non-smoking mother? [1 mark]

- (f) Briefly interpret the value '130.05' in the output. [1 mark]
  
  
  
  
  
  
  
  
  
  
- (g) Why do the residuals have 97 degrees of freedom? [1 mark]
  
  
  
  
  
  
  
  
  
  
- (h) Based on the multiple regression model, is there any evidence of a difference in mean birth weight between smoking and non-smoking mothers? Justify your conclusion with reference to the R output above. [2 marks]
  
  
  
  
  
  
  
  
  
  
- (i) Briefly explain why the conclusion from the multiple regression model might be different to the conclusion from the two-sample t-test in (d). [2 marks]



### Question 3

Energy drinks have become widely popular among adolescents and are also consumed by athletes, particularly those who have just begun their sporting career. A recent paper presented a study on the consumption of energy drinks by teenagers engaged in sports, including quantity consumed and factors that might be associated with consumption.

A total of 707 students, selected randomly from sports classes at various schools, completed a questionnaire on energy drink consumption. The following table shows the cross-tabulation of regular energy drink consumption by gender:

Gender	Energy Drinks	
	Yes	No
Female	192	90
Male	296	129

- (a) Was this an observational or experimental study? Briefly justify your answer.  
[1 mark]
- (b) Overall, what proportion of the students consumed energy drinks? [1 mark]
- (c) What is the estimated difference in the proportions of females and males who consume energy drinks? [1 mark]

- (d) Assuming this sample is representative of all teenagers engaged in sports, give a 95% confidence interval for the true difference in the proportions of females and males who consume energy drinks. What does the interval say about the difference in energy drink consumption between genders? [3 marks]

Another factor recorded on the questionnaire was the frequency of practising sports. The following table gives the summary of results for the 681 students who indicated they practised at least once a week:

Practising Sports	Energy Drinks	
	Yes	No
Daily	328	146
2-3 times per week	28	13
Once per week	114	52

- (e) Based on this table, is there evidence of an association between energy drink consumption and frequency of practising sports? [5 marks]

The schools in the study were from different cities so the authors were also interested in whether there might be differences in energy drink consumption between the cities. They used R to conduct a chi-squared analysis, giving the following results:

Pearson's Chi-squared test

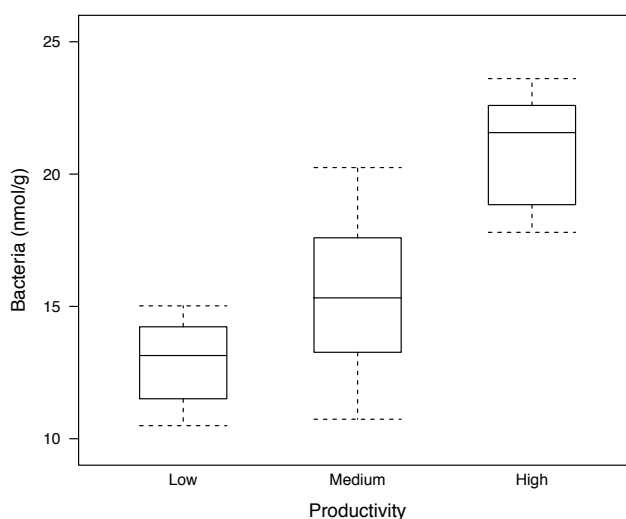
```
data: table(data$City, data$EnergyDrink)
X-squared = 9.8619, df = 1, p-value = 0.001687
```

- (f) How many cities were represented in the study? [1 mark]
- (g) Briefly summarise the conclusion from this chi-squared test. [1 mark]

## Question 4

Albic soil is a common low-yielding soil in eastern China and increasing its productivity is crucial to ongoing food security in the region. A study investigated factors that affect soil quality to gain a better understanding of those factors that are ultimately related to crop yields. They took soil samples from 36 sites and classified each site as either 'Low', 'Medium' or 'High' productivity on the basis of the mean annual crop yield over the previous five years.

One measurement made was the total concentration of bacterial communities in the soil (nmol/g). The figure below shows side-by-side box plots of the bacterial concentration for each productivity level:



The table below shows the observed sample means and standard deviations of the bacterial concentrations for each of the productivity levels.

Productivity	n	Mean	SD
Low	12	12.9	1.51
Medium	12	15.5	2.98
High	12	20.9	2.05

The researchers used one-way analysis of variance (ANOVA) to compare the mean bacterial concentrations between the productivity levels.

- (a) An important assumption for one-way ANOVA is that groups and observations are independent of each other. Give one example of how that assumption might be compromised in a study like this. [1 mark]

- (b) List two assumptions for one-way ANOVA that can be assessed using the above figure. Briefly comment on the validity of each assumption for these data. [2 marks]

- (c) The one-way ANOVA gave a total sum of squares of 569.6 and a residual sum of squares of 168.9. What are the units of these values? [1 mark]

- (d) Using the values from (c), complete the ANOVA table below. [3 marks]

Source	DF	SS	MS	F
Residuals				
Total				

- (e) Give bounds on the  $P$ -value for the  $F$ -test. What do you conclude? [1 mark]

- (f) What is the  $R^2$  value for this model of bacterial concentration? Briefly interpret the value. [2 marks]
- (g) Based on this model and the  $R^2$  value, would you recommend increasing the bacteria in the soil to improve productivity? Briefly explain why or why not. [1 mark]
- (h) In addition to measuring bacterial communities, the researchers also made 25 other measurements of physical, chemical, biochemical and biological properties of each soil sample. Briefly explain an issue in comparing all these measurements between the different productivity levels. What would you recommend to overcome this issue? [2 marks]

**END OF EXAMINATION**

## Formulas and Statistical Tables

### BASICS

$$\bar{x} = \frac{\sum x_j}{n} \quad s = \sqrt{\frac{\sum (x_j - \bar{x})^2}{n-1}}$$

### STANDARDISING

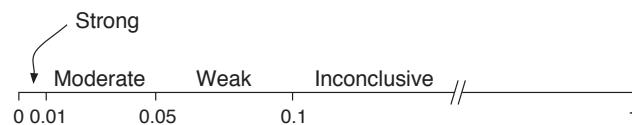
If  $X \sim \text{Normal}(\mu, \sigma)$  then  $Z = \frac{X - \mu}{\sigma} \sim \text{Normal}(0, 1)$

### BINOMIAL RANDOM VARIABLES

$$P(X = x) = \binom{n}{x} p^x (1-p)^{n-x} \quad (\text{but is usually available in tables})$$

$$E(X) = np \quad \text{sd}(X) = \sqrt{np(1-p)} \quad \hat{P} = \frac{X}{n} \quad E(\hat{P}) = p \quad \text{sd}(\hat{P}) = \sqrt{\frac{p(1-p)}{n}}$$

### P-VALUES AND ERRORS



	Decision	
	Retain	Reject
$H_0$ is true	Correct ( $1 - \alpha$ )	Type I Error ( $\alpha$ )
$H_0$ is false	Type II Error ( $\beta$ )	Correct ( $1 - \beta$ )

### TESTS AND CONFIDENCE INTERVALS BASED ON STANDARD ERRORS

$$t = \frac{\text{estimate} - \text{hypothesised}}{\text{se}(\text{estimate})} \quad \text{estimate} \pm t^* \text{se}(\text{estimate})$$

$$\text{se}(\bar{x}) = \frac{s}{\sqrt{n}} \quad \text{se}(\bar{x}_1 - \bar{x}_2) = \sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}} \quad \text{se}(r) = \sqrt{\frac{1-r^2}{n-2}}$$

$$\text{se}(\hat{p}) = \sqrt{\frac{\hat{p}(1-\hat{p})}{n}} \quad \text{se}(\hat{p}_1 - \hat{p}_2) = \sqrt{\frac{\hat{p}_1(1-\hat{p}_1)}{n_1} + \frac{\hat{p}_2(1-\hat{p}_2)}{n_2}}$$

Use  $t$  for means, correlation and regression. Use  $z$  for proportions.



## REGRESSION

$$y = b_0 + b_1x \quad y = b_0 + b_1x + b_2x_1 \quad x_1 = \begin{cases} 1, & \text{if Group B} \\ 0, & \text{if Group A} \end{cases}$$

## POOLED VARIANCE

$$s_p^2 = \frac{(n_1 - 1)s_1^2 + (n_2 - 1)s_2^2}{n_1 + n_2 - 2}$$

## ANOVA TABLES

$$\text{DFT} = n - 1 \quad \text{DFG} = k - 1$$

$$MS = \frac{SS}{DF} \quad R^2 = \frac{SSG}{SST} \quad s_p = \sqrt{MSR} \quad F = \frac{MSG}{MSR}$$

BONFERRONI CORRECTION FOR  $k$  COMPARISONS

$$\alpha = \frac{0.05}{k}$$

## ODDS AND ODDS RATIOS

$$\text{Odds} = \frac{p}{1-p} \quad \text{OR} = \frac{\text{Odds for group B}}{\text{Odds for group A}} \quad \text{se}(\ln(\text{OR})) = \sqrt{\frac{1}{a} + \frac{1}{b} + \frac{1}{c} + \frac{1}{d}}$$

## CHI-SQUARED TESTS

$$\text{expected} = \frac{(\text{row total}) \times (\text{column total})}{\text{overall total}} \quad \chi^2 = \sum \frac{(\text{observed} - \text{expected})^2}{\text{expected}}$$

$$\text{df} = (\# \text{ rows} - 1) \times (\# \text{ columns} - 1)$$

## SIGN TEST

$$\text{Count of positive values is } X \sim \text{Binomial}(n, 0.5)$$

## SIGNED-RANK TEST

$$S = \text{sum of ranks corresponding to positive differences}$$

$$E(S) = \frac{n(n+1)}{4} \quad \text{sd}(S) = \sqrt{\frac{n(n+1)(2n+1)}{24}}$$

## Binomial Distribution

This table gives  $P(X = x)$ , where  $X \sim \text{Binomial}(n, p)$ .

$n$	$x$	$p$							
		0.01	0.05	0.10	0.20	0.25	0.30	0.40	0.50
7	0	0.932	0.698	0.478	0.210	0.133	0.082	0.028	0.008
	1	0.066	0.257	0.372	0.367	0.311	0.247	0.131	0.055
	2	0.002	0.041	0.124	0.275	0.311	0.318	0.261	0.164
	3		0.004	0.023	0.115	0.173	0.227	0.290	0.273
	4			0.003	0.029	0.058	0.097	0.194	0.273
	5				0.004	0.012	0.025	0.077	0.164
	6					0.001	0.004	0.017	0.055
	7							0.002	0.008

Critical values of the  $\chi^2$  distribution

This table gives  $x^*$  such that  $P(X^2 \geq x^*) = p$ , where  $X^2 \sim \chi^2(\text{df})$ .

df	Probability $p$								
	0.975	0.95	0.25	0.10	0.05	0.025	0.01	0.005	0.001
1	0.001	0.004	1.323	2.706	3.841	5.024	6.635	7.879	10.83
2	0.051	0.103	2.773	4.605	5.991	7.378	9.210	10.60	13.82
3	0.216	0.352	4.108	6.251	7.815	9.348	11.34	12.84	16.27
4	0.484	0.711	5.385	7.779	9.488	11.14	13.28	14.86	18.47
5	0.831	1.145	6.626	9.236	11.07	12.83	15.09	16.75	20.52

Critical values of the  $F$  distribution

This table gives  $f^*$  such that  $P(F_{n,d} \geq f^*) = p$ .

$d$	$p$	$n$								
		1	2	3	4	5	6	7	8	9
33	0.100	2.86	2.47	2.26	2.12	2.03	1.96	1.91	1.86	1.83
	0.050	4.14	3.28	2.89	2.66	2.50	2.39	2.30	2.23	2.18
	0.010	7.47	5.31	4.44	3.95	3.63	3.41	3.24	3.11	3.00
	0.001	13.0	8.58	6.88	5.97	5.38	4.98	4.67	4.44	4.26
34	0.100	2.86	2.47	2.25	2.12	2.02	1.96	1.90	1.86	1.82
	0.050	4.13	3.28	2.88	2.65	2.49	2.38	2.29	2.23	2.17
	0.010	7.44	5.29	4.42	3.93	3.61	3.39	3.22	3.09	2.98
	0.001	13.0	8.52	6.83	5.92	5.34	4.93	4.63	4.40	4.22
35	0.100	2.85	2.46	2.25	2.11	2.02	1.95	1.90	1.85	1.82
	0.050	4.12	3.27	2.87	2.64	2.49	2.37	2.29	2.22	2.16
	0.010	7.42	5.27	4.40	3.91	3.59	3.37	3.20	3.07	2.96
	0.001	12.9	8.47	6.79	5.88	5.30	4.89	4.59	4.36	4.18

## Probabilities for the Standard Normal distribution

This table gives  $P(Z \geq z)$  for  $Z \sim \text{Normal}(0, 1)$ .

$z$	Second decimal place of $z$									
	0	1	2	3	4	5	6	7	8	9
0.0	0.500	0.496	0.492	0.488	0.484	0.480	0.476	0.472	0.468	0.464
0.1	0.460	0.456	0.452	0.448	0.444	0.440	0.436	0.433	0.429	0.425
0.2	0.421	0.417	0.413	0.409	0.405	0.401	0.397	0.394	0.390	0.386
0.3	0.382	0.378	0.374	0.371	0.367	0.363	0.359	0.356	0.352	0.348
0.4	0.345	0.341	0.337	0.334	0.330	0.326	0.323	0.319	0.316	0.312
0.5	0.309	0.305	0.302	0.298	0.295	0.291	0.288	0.284	0.281	0.278
0.6	0.274	0.271	0.268	0.264	0.261	0.258	0.255	0.251	0.248	0.245
0.7	0.242	0.239	0.236	0.233	0.230	0.227	0.224	0.221	0.218	0.215
0.8	0.212	0.209	0.206	0.203	0.200	0.198	0.195	0.192	0.189	0.187
0.9	0.184	0.181	0.179	0.176	0.174	0.171	0.169	0.166	0.164	0.161
1.0	0.159	0.156	0.154	0.152	0.149	0.147	0.145	0.142	0.140	0.138
1.1	0.136	0.133	0.131	0.129	0.127	0.125	0.123	0.121	0.119	0.117
1.2	0.115	0.113	0.111	0.109	0.107	0.106	0.104	0.102	0.100	0.099
1.3	0.097	0.095	0.093	0.092	0.090	0.089	0.087	0.085	0.084	0.082
1.4	0.081	0.079	0.078	0.076	0.075	0.074	0.072	0.071	0.069	0.068
1.5	0.067	0.066	0.064	0.063	0.062	0.061	0.059	0.058	0.057	0.056
1.6	0.055	0.054	0.053	0.052	0.051	0.049	0.048	0.047	0.046	0.046
1.7	0.045	0.044	0.043	0.042	0.041	0.040	0.039	0.038	0.038	0.037
1.8	0.036	0.035	0.034	0.034	0.033	0.032	0.031	0.031	0.030	0.029
1.9	0.029	0.028	0.027	0.027	0.026	0.026	0.025	0.024	0.024	0.023
2.0	0.023	0.022	0.022	0.021	0.021	0.020	0.020	0.019	0.019	0.018
2.1	0.018	0.017	0.017	0.017	0.016	0.016	0.015	0.015	0.015	0.014
2.2	0.014	0.014	0.013	0.013	0.013	0.012	0.012	0.012	0.011	0.011
2.3	0.011	0.010	0.010	0.010	0.010	0.009	0.009	0.009	0.009	0.008
2.4	0.008	0.008	0.008	0.008	0.007	0.007	0.007	0.007	0.007	0.006
2.5	0.006	0.006	0.006	0.006	0.006	0.005	0.005	0.005	0.005	0.005
2.6	0.005	0.005	0.004	0.004	0.004	0.004	0.004	0.004	0.004	0.004
2.7	0.003	0.003	0.003	0.003	0.003	0.003	0.003	0.003	0.003	0.003
2.8	0.003	0.002	0.002	0.002	0.002	0.002	0.002	0.002	0.002	0.002
2.9	0.002	0.002	0.002	0.002	0.002	0.002	0.002	0.001	0.001	0.001
3.0	0.001	0.001	0.001	0.001	0.001	0.001	0.001	0.001	0.001	0.001
3.1	0.001	0.001	0.001	0.001	0.001	0.001	0.001	0.001	0.001	0.001
3.2	0.001	0.001	0.001	0.001	0.001	0.001	0.001	0.001	0.001	0.001
3.3										

Critical values of Student's  $T$  distribution

This table gives  $t^*$  such that  $P(T \geq t^*) = p$ , where  $T \sim \text{Student}(\text{df})$ .

df	Probability $p$								
	0.25	0.10	0.05	0.025	0.01	0.005	0.001	0.0005	0.0001
1	1.000	3.078	6.314	12.71	31.82	63.66	318.3	636.6	3183
2	0.816	1.886	2.920	4.303	6.965	9.925	22.33	31.60	70.70
3	0.765	1.638	2.353	3.182	4.541	5.841	10.21	12.92	22.20
4	0.741	1.533	2.132	2.776	3.747	4.604	7.173	8.610	13.03
5	0.727	1.476	2.015	2.571	3.365	4.032	5.893	6.869	9.678
6	0.718	1.440	1.943	2.447	3.143	3.707	5.208	5.959	8.025
7	0.711	1.415	1.895	2.365	2.998	3.499	4.785	5.408	7.063
8	0.706	1.397	1.860	2.306	2.896	3.355	4.501	5.041	6.442
9	0.703	1.383	1.833	2.262	2.821	3.250	4.297	4.781	6.010
10	0.700	1.372	1.812	2.228	2.764	3.169	4.144	4.587	5.694
11	0.697	1.363	1.796	2.201	2.718	3.106	4.025	4.437	5.453
12	0.695	1.356	1.782	2.179	2.681	3.055	3.930	4.318	5.263
13	0.694	1.350	1.771	2.160	2.650	3.012	3.852	4.221	5.111
14	0.692	1.345	1.761	2.145	2.624	2.977	3.787	4.140	4.985
15	0.691	1.341	1.753	2.131	2.602	2.947	3.733	4.073	4.880
16	0.690	1.337	1.746	2.120	2.583	2.921	3.686	4.015	4.791
17	0.689	1.333	1.740	2.110	2.567	2.898	3.646	3.965	4.714
18	0.688	1.330	1.734	2.101	2.552	2.878	3.610	3.922	4.648
18.42	0.688	1.329	1.732	2.097	2.547	2.871	3.597	3.905	4.623
19	0.688	1.328	1.729	2.093	2.539	2.861	3.579	3.883	4.590
20	0.687	1.325	1.725	2.086	2.528	2.845	3.552	3.850	4.539
30	0.683	1.310	1.697	2.042	2.457	2.750	3.385	3.646	4.234
40	0.681	1.303	1.684	2.021	2.423	2.704	3.307	3.551	4.094
50	0.679	1.299	1.676	2.009	2.403	2.678	3.261	3.496	4.014
60	0.679	1.296	1.671	2.000	2.390	2.660	3.232	3.460	3.962
$\infty$	0.674	1.282	1.645	1.960	2.326	2.576	3.090	3.291	3.719