

Name:

USC ID:

Notes:

- Write your name and ID number in the solution you submit.
- No books, cell phones or other notes are permitted. Only one letter size cheat sheet (back and front), a tablet for writing down your solutions, and a calculator are allowed.
- Problems are not sorted in terms of difficulty. Please avoid guess work and long and irrelevant answers.
- Show all your work and your final answer. Simplify your answer as much as you can.
- Open your exam only when you are instructed to do so.
- The exam has 5 questions, 11 pages, and 13 points extra credit.

Problem	Score	Earned
1	22	
2	25	
3	22	
4	22	
5	22	
Total	113	

1. As director of the local tourist board, you are interested in determining the factors that influence the hotel occupancy rate in your city each month. Hotel occupancy can be measured as the percentage of available hotel rooms that are occupied by paying customers. You develop the following model: $Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \beta_3 X_3 + \beta_4 X_4 + \epsilon$, where Y is the hotel occupancy rate, X_1 is the total number of passengers arriving at the airport, X_2 is a price index of local hotel room rates, X_3 is the consumer confidence index, and X_4 is a dummy variable that shows being in the months of June, July, and August. You look at data from the past 35 months and obtain the following results: $y = 67.1 + 0.02x_1 - 0.055x_2 + 0.08x_3 + 12.3x_4$. Also, assume that the standard errors of $\hat{\beta}_i$'s are calculated as $SE(\hat{\beta}_0) = 58.3, SE(\hat{\beta}_1) = 0.008, SE(\hat{\beta}_2) = 0.01, SE(\hat{\beta}_3) = 0.06, SE(\hat{\beta}_4) = 4.7$. Moreover, the regression sum of squares is $RegSS = 1,169.45$, and the residual (error) sum of squares is $RSS = 576$.
 - (a) Interpret the estimated regression coefficient $\hat{\beta}_2$.
 - (b) Interpret the estimated regression coefficient $\hat{\beta}_4$.
 - (c) Test the hypothesis $H_0 : \beta_1 = 0$ at $\alpha_1 = 0.05$ and $\alpha_2 = 0.01$ and interpret your results. Based on your results, explain why level of significance has to be determined *before* testing a hypothesis.
 - (d) Test the hypothesis $H_0 : \beta_1 = \beta_2 = \beta_3 = \beta_4 = 0$ at $\alpha = 0.05$. Interpret your results.
 - (e) Build a 98% confidence interval for β_3 . What happens to the confidence interval if the level of confidence is increased? Is it desirable to have a very high level of confidence?

2. Choose either T (True) or F (False) (no need to explain why):

- (a) When the assumption of conditional independence of features holds, the Naïve Bayes' classifier provides the best accuracy among all possible classifiers. T F
- (b) The F1 score is not an appropriate measure for evaluating binary classifiers when data are not imbalanced. T F
- (c) Leave-One-Out Cross Validation has less bias in estimating the error of a classifier for a large data set than 5 fold cross validation. T F
- (d) When classifying imbalanced data into two classes, we can decrease the threshold on class conditional probability $\Pr(Y = k|X_1 = x_1, \dots, X_p = x_p]$ to increase the true positive rate at the expense of increasing the false negative rate. T F
- (e) Logistic regression assumes that the conditional odds of the outcome Y given the features, $\mathbb{O}[Y = k|X_1 = x_1, \dots, X_p = x_p]$, is a logistic function of the features. T F

3. Assume that we have a binary classification problem with only one feature in which the conditional distribution of the feature in class $k = 1$ is a normal with mean $\mu_1 = 1$ and standard deviation $\sigma_1 = 1$ and the conditional distribution of the feature in class $k = 2$ is normal with mean $\mu_2 = 2$ and $\sigma_2 = 1$. Determine the values of x that are classified in each class by the Bayes' optimal classifier. Assume the classes are balanced, i.e. $\pi_1 = \pi_2 = 0.5$.

Important note: you must *derive* the decision rule from scratch, i.e. you must write down the posterior probabilities.

4. Consider the logistic regression method for binary classification ($Y = 0$ or $Y = 1$) with two features $\mathbf{X} = (X_1, X_2)$, formulated by

$$P(Y = 1|\mathbf{X}) = \frac{e^{\beta_0 + \beta_1 X_1 + \beta_2 X_2}}{1 + e^{\beta_0 + \beta_1 X_1 + \beta_2 X_2}}$$

Assume that using a dataset of 200 observations, we obtained the following estimates:

	Coefficient	Standard Error
β_0	1	0.2
β_1	2	0.1
β_2	1	s

- (a) Find all values of s that makes the coefficient β_2 statistically insignificant. You can consider the significance level to be $\alpha = 0.05$.
- (b) Determine the equation for the decision boundary for this classifier.
- (c) In what class will $\mathbf{X} = (-1, 1)$ will be classified?

5. In a weird simulated world, we have three types of creatures. Each creature has between 1 to 100 legs, 1 to 100 teeth, and 1 to 100 noses. The fraction of creatures type-1, type-2, and type-3 are respectively $l/(l+t+n)$, $t/(l+t+n)$, and $n/(l+t+n)$, where l, t, n are respectively the number of legs, teeth, and noses of a creature.
- (a) If a creature has 10 teeth, 25 legs, and 30 noses, what is your best guess about the type of the creature?
 - (b) What type of supervised learning problem are you solving in this question? Explain.

Scratch paper

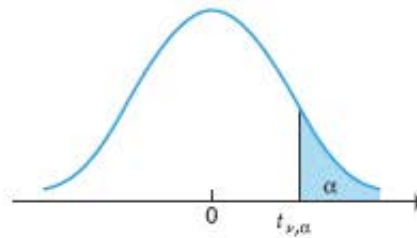
Name:

USC ID:

Scratch paper

Name:

USC ID:

Upper Critical Values of Student's t Distribution with ν Degrees of Freedom

For selected probabilities, α , the table shows the values $t_{\nu, \alpha}$ such that $P(t_{\nu} > t_{\nu, \alpha}) = \alpha$, where t_{ν} is a Student's t random variable with ν degrees of freedom. For example, the probability is .10 that a Student's t random variable with 10 degrees of freedom exceeds 1.372.

PROBABILITY OF EXCEEDING THE CRITICAL VALUE						
ν	0.10	0.05	0.025	0.01	0.005	0.001
1	3.078	6.314	12.706	31.821	63.657	318.313
2	1.886	2.920	4.303	6.965	9.925	22.327
3	1.638	2.353	3.182	4.541	5.841	10.215
4	1.533	2.132	2.776	3.747	4.604	7.173
5	1.476	2.015	2.571	3.365	4.032	5.893
6	1.440	1.943	2.447	3.143	3.707	5.208
7	1.415	1.895	2.365	2.998	3.499	4.782
8	1.397	1.860	2.306	2.896	3.355	4.499
9	1.383	1.833	2.262	2.821	3.250	4.296
10	1.372	1.812	2.228	2.764	3.169	4.143
11	1.363	1.796	2.201	2.718	3.106	4.024
12	1.356	1.782	2.179	2.681	3.055	3.929
13	1.350	1.771	2.160	2.650	3.012	3.852
14	1.345	1.761	2.145	2.624	2.977	3.787
15	1.341	1.753	2.131	2.602	2.947	3.733
16	1.337	1.746	2.120	2.583	2.921	3.686
17	1.333	1.740	2.110	2.567	2.898	3.646
18	1.330	1.734	2.101	2.552	2.878	3.610
19	1.328	1.729	2.093	2.539	2.861	3.579
20	1.325	1.725	2.086	2.528	2.845	3.552
21	1.323	1.721	2.080	2.518	2.831	3.527
22	1.321	1.717	2.074	2.508	2.819	3.505
23	1.319	1.714	2.069	2.500	2.807	3.485
24	1.318	1.711	2.064	2.492	2.797	3.467
25	1.316	1.708	2.060	2.485	2.787	3.450
26	1.315	1.706	2.056	2.479	2.779	3.435
27	1.314	1.703	2.052	2.473	2.771	3.421
28	1.313	1.701	2.048	2.467	2.763	3.408
29	1.311	1.699	2.045	2.462	2.756	3.396
30	1.310	1.697	2.042	2.457	2.750	3.385
40	1.303	1.684	2.021	2.423	2.704	3.307
60	1.296	1.671	2.000	2.390	2.660	3.232
100	1.290	1.660	1.984	2.364	2.626	3.174
∞	1.282	1.645	1.960	2.326	2.576	3.090
ν	0.10	0.05	0.025	0.01	0.005	0.001

F - Distribution ($\alpha = 0.05$ in the Right Tail)

df ₂ \ df ₁		Numerator Degrees of Freedom								
		1	2	3	4	5	6	7	8	9
1	161.45	199.50	215.71	224.58	230.16	233.99	236.77	238.88	240.54	
2	18.513	19.000	19.164	19.247	19.296	19.330	19.353	19.371	19.385	
3	10.128	9.5521	9.2766	9.1172	9.0135	8.9406	8.8867	8.8452	8.8123	
4	7.7086	9.9443	6.5914	6.3882	6.2561	6.1631	6.0942	6.0410	6.9988	
5	6.6079	5.7861	5.4095	5.1922	5.0503	4.9503	4.8759	4.8183	4.7725	
6	5.9874	5.1433	4.7571	4.5337	4.3874	4.2839	4.2067	4.1468	4.0990	
7	5.5914	4.7374	4.3468	4.1203	3.9715	3.8660	3.7870	3.7257	3.6767	
8	5.3177	4.4590	4.0662	3.8379	3.6875	3.5806	3.5005	3.4381	3.3881	
9	5.1174	4.2565	3.8625	3.6331	3.4817	3.3738	3.2927	3.2296	3.1789	
10	4.9646	4.1028	3.7083	3.4780	3.3258	3.2172	3.1355	3.0717	3.0204	
11	4.8443	3.9823	3.5874	3.3567	3.2039	3.0946	3.0123	2.9480	2.8962	
12	4.7472	3.8853	3.4903	3.2592	3.1059	2.9961	2.9134	2.8486	2.7964	
13	4.6672	3.8056	3.4105	3.1791	3.0254	2.9153	2.8321	2.7669	2.7144	
14	4.6001	3.7389	3.3439	3.1122	2.9582	2.8477	2.7642	2.6987	2.6458	
15	4.5431	3.6823	3.2874	3.0556	2.9013	2.7905	2.7066	2.6408	2.5876	
16	4.4940	3.6337	3.2389	3.0069	2.8524	2.7413	2.6572	2.5911	2.5377	
17	4.4513	3.5915	3.1968	2.9647	2.8100	2.6987	2.6143	2.5480	2.4943	
18	4.4139	3.5546	3.1599	2.9277	2.7729	2.6613	2.5767	2.5102	2.4563	
19	4.3807	3.5219	3.1274	2.8951	2.7401	2.6283	2.5435	2.4768	2.4227	
20	4.3512	3.4928	3.0984	2.8661	2.7109	2.5990	2.5140	2.4471	2.3928	
21	4.3248	3.4668	3.0725	2.8401	2.6848	2.5727	2.4876	2.4205	2.3660	
22	4.3009	3.4434	3.0491	2.8167	2.6613	2.5491	2.4638	2.3965	2.3419	
23	4.2793	3.4221	3.0280	2.7955	2.6400	2.5277	2.4422	2.3748	2.3201	
24	4.2597	3.4028	3.0088	2.7763	2.6207	2.5082	2.4226	2.3551	2.3002	
25	4.2417	3.3852	2.9912	2.7587	2.6030	2.4904	2.4047	2.3371	2.2821	
26	4.2252	3.3690	2.9752	2.7426	2.5868	2.4741	2.3883	2.3205	2.2655	
27	4.2100	3.3541	2.9604	2.7278	2.5719	2.4591	2.3732	2.3053	2.2501	
28	4.1960	3.3404	2.9467	2.7141	2.5581	2.4453	2.3593	2.2913	2.2360	
29	4.1830	3.3277	2.9340	2.7014	2.5454	2.4324	2.3463	2.2783	2.2229	
30	4.1709	3.3158	2.9223	2.6896	2.5336	2.4205	2.3343	2.2662	2.2107	
40	4.0847	3.2317	2.8387	2.6060	2.4495	2.3359	2.2490	2.1802	2.1240	
60	4.0012	3.1504	2.7581	2.5252	2.3683	2.2541	2.1665	2.0970	2.0401	
120	3.9201	3.0718	2.6802	2.4472	2.2899	2.1750	2.0868	2.0164	1.9588	
∞	3.8415	2.9957	2.6049	2.3719	2.2141	2.0986	2.0096	1.9384	1.8799	

Denominator Degrees of Freedom

Cumulative Distribution Function, $F(z)$, of the Standard Normal Distribution Table

z	0	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
z	0	0.01	0.02	0.03	0.04	0.05	0.06	0.07	0.08	0.09

Cumulative Distribution Function, $F(z)$, of the Standard Normal Distribution Table