

统计学大作业——利用 SPSS 探索学生入学申请指标的分布规律以及对研究生申请成功概率的影响（以 UCLA 为例）

陈冠华 2017202095

一、研究背景与意义：

每年平均有 40 万不同教育背景的中国学生选择留学深造，即使是在疫情肆虐，大国关系微妙的 2020 年，留学依旧还是许多国内学生的出路选项之一。在研究生领域，每年出国的学生数量达到了约 23 万，其中不乏名校毕业的应届生。对于所有准备出国的学生而言，“刷指标”是必修功课，因此如何分配有限的精力，着重提升哪些指标便是一个开放的，依赖智慧的难题。这份报告从数据的角度洞察各种申请指标的分布规律以及对申请成功概率的影响，能使学生对国外大学的偏好有更深入的认识，进而拥有对竞争者的认知优势，在准备阶段更高效地完成申请工作。

二、数据来源与介绍

数据集的样本量为 500，数据来源是 UCLA 的研究生信息数据库，下载于数据竞赛网站 www.kaggle.com。对数据集的介绍见表一

表一

变量名称	中文释义	数据解释	数据单位	数据类型
GRE Scores (out of 340)	GRE 分数	GRE 分数	分	Int
TOEFL Scores (out of 120)	TOEFL 分数	TOEFL 分数	分	Int
University Rating (out of 5)	学校档次	学校档次（共分为 5 档）	无	Int

Statement of Purpose	个人申请分数	个人申请获得的分数	分	Float
Letter of Recommendation Strength (out of 5)	推荐信分数	推荐信获得的分数	分	Float
Undergraduate GPA (out of 10)	毕业成绩	毕业成绩	无	Float
Research Experience (either 0 or 1)	研究经历	是否拥有科研经历	无	Int
Chance of Admit (ranging from 0 to 1)	被接受的概率	最终申请成功的概率	无	Float

三、数据分析

1.正态性检验， 样本比例检验

我们采用 K-S 检验方法完成对连续数据的正态性检验，用卡方分布，二项分布检验类别数据中的各个类别满足等概率发生的假设。假设检验的汇总结果如下图所示。

假设检验汇总				
	原假设	测试	Sig.	决策者
1	由 Research = 1.00 和 0.00 定义类别以 0.5 和 0.5 的概率发生。	单样本 Binomial 检验	.008	拒绝原假设。
2	UniversityRating 的类别以相同概率发生。	单样本卡方检验	.000	拒绝原假设。
3	TOEFLScore 的分布为正态分布，平均值为 107.19，标准差为 6.08。	单样本 Kolmogorov-Smirnov 检验	.023	拒绝原假设。
4	StatementofPurpose 的分布为正态分布，平均值为 3.37，标准差为 0.99。	单样本 Kolmogorov-Smirnov 检验	.000	拒绝原假设。
5	LetterofRecommendationStrength 的分布为正态分布，平均值为 3.48，标准差为 0.93。	单样本 Kolmogorov-Smirnov 检验	.000	拒绝原假设。
6	CGPA 的分布为正态分布，平均值为 8.58，标准差为 0.60。	单样本 Kolmogorov-Smirnov 检验	.238	保留原假设。
7	ChanceofAdmit 的分布为正态分布，平均值为 0.72，标准差为 0.14。	单样本 Kolmogorov-Smirnov 检验	.254	保留原假设。
8	GREScore 的分布为正态分布，平均值为 316.47，标准差为 11.30。	单样本 Kolmogorov-Smirnov 检验	.168	保留原假设。

显示渐进显著性。显著性水平是 .05。

显然有三个连续指标不符合正态分布，而两个类别数据均不满足类发生概率相等的假设。由于 UCLA 在美国大学中排名较高，因此我们可以合理推断在托福成绩，个人陈述与推荐信得分三个指标上分布呈左偏分布的，而考虑学校档次以及是否拥有科研经历时，我们更倾向于那些高档次学校以及拥有科研经历的学生占比会更高一些。

2.参数估计

A . 均值

考虑到一部分特征不服从正态分布，因此只对那些通过了正态性验证的连续特征进行均值的参数估计，同时统一取置信区间为 95%。均值的参数估计对学生的启示是在 GRE 以及在校成绩上至少要达到均值下限的水准才能在该指标上对至少 50%的竞争者形成优势，在该指标还未达到均值下限的学生应当把精力着重放在相应部分上。

描述

			统计量	标准误
GREScore	均值		316.47	.505
	均值的 95% 置信区间	下限	315.48	
		上限	317.46	
	5% 修整均值		316.43	
	中值		317.00	
	方差		127.580	
	标准差		11.295	
	极小值		290	
	极大值		340	
	范围		50	
	四分位距		17	
	偏度		-.040	.109
	峰度		-.711	.218
CGPA	均值		8.5764	.02705
	均值的 95% 置信区间	下限	8.5233	
		上限	8.6296	
	5% 修整均值		8.5779	
	中值		8.5600	
	方差		.366	
	标准差		.60481	
	极小值		6.80	
	极大值		9.92	
	范围		3.12	
	四分位距		.92	
	偏度		-.027	.109
	峰度		-.561	.218
ChanceofAdmit	均值		.7217	.00631
	均值的 95% 置信区间	下限	.7093	
		上限	.7341	
	5% 修整均值		.7257	
	中值		.7200	
	方差		.020	
	标准差		.14114	
	极小值		.34	
	极大值		.97	
	范围		.63	
	四分位距		.19	
	偏度		-.290	.109
	峰度		-.455	.218

接着我们考虑双总体均值之差的参数估计。我们按照托福分数划分，以 105 分为界限，计算出 105 分及以下学生的申请成功概率以及 105 以上学生申请成

功概率之间均值差的置信区间。

独立样本检验									
		方差方程的 Levene 检验		均值方程的 t 检验					
		F	Sig.	t	df	Sig.(双侧)	均值差值	标准误差值	差分的 95% 置信区间
申请成功概率	假设方差相等	.441	.507	21.505	711	.000	.17896501	.00832211	.16262616 .19530386
	假设方差不相等			21.292	592.348	.000	.17896501	.00840520	.16245739 .19547263

可以看出假设两者方差相等，sig 值并不显著，因此不拒绝该假设。同时在两种假设下，托福成绩在 105 分以上的同学申请成功的概率比 105 分以下的同学申请成功的概率平均要高出 16%-19%的概率。对英语的应用能力向来是申请院校着重的考察对象，因此理论上更高的托福分数自然对应着更高的申请成功概率。该分析结果给申请者的托福成绩目标提供了指导，对于那些分数尚未达到 105 分的同学，可以通过提升托福分数达到提升申请成功概率的目的。

使用同样的方法，对推荐信得分超过 3.5 的学生以及得分小于等于 3.5 分的学生之间的申请成功概率均值之差进行参数估计，得到以下结果。

独立样本检验									
		方差方程的 Levene 检验		均值方程的 t 检验					
		F	Sig.	t	df	Sig.(双侧)	均值差值	标准误差值	差分的 95% 置信区间
申请成功概率	假设方差相等	1.836	.176	-11.821	359	.000	-.142393887	.012045693	-.166082904 -.118704871
	假设方差不相等			-11.958	358.612	.000	-.142393887	.011907456	-.165811103 -.118976671

推荐信往往能突出学生特点，给面试官留下更好的印象，同时业界与学界权威人士的推荐信也是对学生能力的背书。从参数估计的结果中可以看出推荐信得分高于 3.5 的学生申请成功概率大约要比小于等于 3.5 的学生平均高出 11%-16%。因此在推荐信上还有提升得分空间的同学可以考虑通过润色，请求老师或者老板的帮助来增加得分，进而提升申请成功概率。

B . 比例

对拥有研究经历的学生占总体样本的比例进行参数估计，得到以下结果。设置置信区间为 95%。

re / 样本总量 的比率统计量

均值	均值的 95% 置信区间		价格相关微分	离散系数	方差系数
	下限	上限			中值居中
.560	.516	.604	1.000	.440	66.4%

通过假设比率的正态分布构建置信区间。

可以看到有科研经历的学生比例估计区间大约在51%到60.4%之间，说明大多数的申请者均拥有至少一次研究经历。

再对学校档次进行参数估计，分别求出五个档次的学生分别占总体的比例区间（置信度取为95%）。

学校档次为1 / 样本总量 的比率统计量

均值	均值的 95% 置信区间		价格相关微分	离散系数	方差系数
	下限	上限			中值居中
.068	.046	.090	1.000	.	.

通过假设比率的正态分布构建置信区间。

学校档次为2 / 样本总量 的比率统计量

均值	均值的 95% 置信区间		价格相关微分	离散系数	方差系数
	下限	上限			中值居中
.252	.214	.290	1.000	.	.

通过假设比率的正态分布构建置信区间。

学校档次为3 / 样本总量 的比率统计量

均值	均值的 95% 置信区间		价格相关微分	离散系数	方差系数
	下限	上限			中值居中
.324	.283	.365	1.000	.	.

通过假设比率的正态分布构建置信区间。

学校档次为4 / 样本总量 的比率统计量

均值	均值的 95% 置信区间		价格相关微分	离散系数	方差系数
	下限	上限			中值居中
.210	.174	.246	1.000	.	.

通过假设比率的正态分布构建置信区间。

学校档次为5 / 样本总量 的比率统计量

均值	均值的 95% 置信区间		价格相关微分	离散系数	方差系数
	下限	上限			中值居中
.146	.115	.177	1.000	.	.

通过假设比率的正态分布构建置信区间。

可以看出当中三个档次的学生占比是较高的。

3.相关性分析

采用皮尔逊相关系数计算各个指标和申请成功概率之间的相关性与显著程度。

相关性

		GREScore	ChanceofAdmit
GREScore	Pearson 相关性	1	.810**
	显著性（双侧）		.000
	N	500	500
ChanceofAdmit	Pearson 相关性	.810**	1
	显著性（双侧）	.000	
	N	500	500

** 在 .01 水平（双侧）上显著相关。

相关性

		ChanceofAdmit	TOEFLScore
ChanceofAdmit	Pearson 相关性	1	.792**
	显著性（双侧）		.000
	N	500	500
TOEFLScore	Pearson 相关性	.792**	1
	显著性（双侧）	.000	
	N	500	500

** 在 .01 水平（双侧）上显著相关。

相关性

		ChanceofAdmit	UniversityRating
ChanceofAdmit	Pearson 相关性	1	.690**
	显著性（双侧）		.000
	N	500	500
UniversityRating	Pearson 相关性	.690**	1
	显著性（双侧）	.000	
	N	500	500

** 在 .01 水平（双侧）上显著相关。

相关性

		ChanceofAdmit	StatementofPurpose
ChanceofAdmit	Pearson 相关性	1	.684**
	显著性 (双侧)		.000
	N	500	500
StatementofPurpose	Pearson 相关性	.684**	1
	显著性 (双侧)	.000	
	N	500	500

** . 在 .01 水平 (双侧) 上显著相关。

相关性

		ChanceofAdmit	LetterofRecommendationStrength
ChanceofAdmit	Pearson 相关性	1	.645**
	显著性 (双侧)		.000
	N	500	500
LetterofRecommendationStrength	Pearson 相关性	.645**	1
	显著性 (双侧)	.000	
	N	500	500

** . 在 .01 水平 (双侧) 上显著相关。

相关性

		ChanceofAdmit	CGPA
ChanceofAdmit	Pearson 相关性	1	.882**
	显著性 (双侧)		.000
	N	500	500
CGPA	Pearson 相关性	.882**	1
	显著性 (双侧)	.000	
	N	500	500

** . 在 .01 水平 (双侧) 上显著相关。

相关性

		ChanceofAdmit	Research
ChanceofAdmit	Pearson 相关性	1	.546**
	显著性 (双侧)		.000
	N	500	500
Research	Pearson 相关性	.546**	1
	显著性 (双侧)	.000	
	N	500	500

** . 在 .01 水平 (双侧) 上显著相关。

所有的指标均和申请成功概率有着显著的线性关系，皮尔逊相关系数从高到低的指标依次是在校成绩，GRE 成绩，托福成绩，学校档次，个人陈述，推荐信以及科研经历。可见在校成绩，英语成绩对申请是否成功的影响至关重要，事实上这个排序结果也符合目前大多数市场上留学辅导机构的观点。皮尔逊相关系数反映了变量之间的线性关系，其研究结果为不同的指标派出了先后顺序，为准备出国但精力有限的中国学生提供了高效率准备材料，提升背景的理论依据。如果一个学生想要尽可能快速地准备出国的话，他首先应该保证自己的在校成绩水平以及英语标准化考试的成绩。

4.方差分析

选取相关性分析中与申请成功概率皮尔逊系数最高的在校成绩进行方差分析。将成绩分为小于 8.0，8.0 到 8.5，8.5 到 9.0，9.0 到 9.5 以及 9.5 到 10.0 五个类别，并重编码为 1，2，3，4，5。由于数据量较大因此假设五个总体均满足正态分布，而显然他们是彼此独立的。

下表显示五个总体拒绝了方差齐性的原假设，但上网查阅资料后发现 SPSS 的方差分析基于最小二乘法的框架，对于方差齐性的条件并不敏感，因此虽然数据在理论上不符合方差分析的条件，但是 SPSS 的分析结果依旧是可信的。

方差齐性检验

经过排名后的申请成功率

Levene 统计量	df1	df2	显著性
10.376	4	495	.000

单因素方差分析

经过排名后的申请成功率

	平方和	df	均方	F	显著性
组间	7.190	4	1.797	323.427	.000
组内	2.751	495	.006		
总数	9.940	499			

这里拒绝原假设，说明不同水平的在校成绩对申请成功率的影响是显著的。

进一步两两分析，发现每个水平对申请成功概率的影响都是显著的，从均值差可以看出，在校成绩档次越高，申请成功的概率越大。

多重比较

因变量: 经过排名后的申请成功率

		均值差 (I-J)	标准误	显著性	95% 置信区间	
(I) 在校成绩档次	(J) 在校成绩档次				下限	上限
Dunnett T3	1	-.12026224 [*]	.01194821	.000	-.1541634	-.0863611
		-.20025194 [*]	.01182991	.000	-.2338344	-.1666694
		-.32132883 [*]	.01156572	.000	-.3541963	-.2884613
		-.40674242 [*]	.01043821	.000	-.4365903	-.3768946
	2	.12026224 [*]	.01194821	.000	.0863611	.1541634
		-.07998970 [*]	.00935591	.000	-.1063884	-.0535910
		-.20106659 [*]	.00901955	.000	-.2265304	-.1756027
		-.28648019 [*]	.00751947	.000	-.3077895	-.2651709
	3	.20025194 [*]	.01182991	.000	.1666694	.2338344
		.07998970 [*]	.00935591	.000	.0535910	.1063884
		-.12107689 [*]	.00886223	.000	-.1461089	-.0960449
		-.20649049 [*]	.00733003	.000	-.2272848	-.1856962
	4	.32132883 [*]	.01156572	.000	.2884613	.3541963
		.20106659 [*]	.00901955	.000	.1756027	.2265304
		.12107689 [*]	.00886223	.000	.0960449	.1461089
		-.08541360 [*]	.00689555	.000	-.1050082	-.0658189
	5	.40674242 [*]	.01043821	.000	.3768946	.4365903
		.28648019 [*]	.00751947	.000	.2651709	.3077895
		.20649049 [*]	.00733003	.000	.1856962	.2272848
		.08541360 [*]	.00689555	.000	.0658189	.1050082
Games-Howell	1	-.12026224 [*]	.01194821	.000	-.1532341	-.0872904
		-.20025194 [*]	.01182991	.000	-.2329116	-.1675922
		-.32132883 [*]	.01156572	.000	-.3532886	-.2893690
		-.40674242 [*]	.01043821	.000	-.4357419	-.3777429
	2	.12026224 [*]	.01194821	.000	.0872904	.1532341
		-.07998970 [*]	.00935591	.000	-.1056839	-.0542955
		-.20106659 [*]	.00901955	.000	-.2258490	-.1762842
		-.28648019 [*]	.00751947	.000	-.3072087	-.2657516
	3	.20025194 [*]	.01182991	.000	.1675922	.2329116
		.07998970 [*]	.00935591	.000	.0542955	.1056839
		-.12107689 [*]	.00886223	.000	-.1454374	-.0967163
		-.20649049 [*]	.00733003	.000	-.2267152	-.1862657
	4	.32132883 [*]	.01156572	.000	.2893690	.3532886
		.20106659 [*]	.00901955	.000	.1762842	.2258490
		.12107689 [*]	.00886223	.000	.0967163	.1454374
		-.08541360 [*]	.00689555	.000	-.1044672	-.0663599
	5	.40674242 [*]	.01043821	.000	.3777429	.4357419
		.28648019 [*]	.00751947	.000	.2657516	.3072087
		.20649049 [*]	.00733003	.000	.1862657	.2267152
		.08541360 [*]	.00689555	.000	.0663599	.1044672

*. 均值差的显著性水平为 0.05。

方差分析是对相关性分析的补充，在验证了在校成绩对申请成功率有着较大影响的基础上，通过方差分析证实了每提升在校成绩的一个档次，便会有统计学意义上显著的申请成功率的提升。因此准备申请的学生应该将绝大部分精力放在学校的表现上，花更多的时间取得更好的成绩是对申请最有帮助的事情。

5.多元线性回归

选取所有的连续性数据作为自变量，申请成功概率作为因变量进行多元线性回归分析。假设各个连续型变量服从正态分布与方差齐性的特点，随即进行多重共线性的检验：

系数 ^a										
模型	非标准化系数			标准系数		相关性			共线性统计量	
	B	标准 误差	试用版	t	Sig.	零阶	偏	部分	容差	VIF
1 (常量)	-1.044	.042		-24.689	.000					
CGPA	.206	.005	.882	41.855	.000	.882	.882	.882	1.000	1.000
2 (常量)	-1.635	.091		-17.901	.000					
CGPA	.156	.008	.670	18.826	.000	.882	.645	.378	.318	3.145
GREScore	.003	.000	.257	7.206	.000	.810	.308	.145	.318	3.145
3 (常量)	-1.532	.091		-16.907	.000					
CGPA	.135	.009	.580	15.183	.000	.882	.563	.296	.260	3.840
GREScore	.003	.000	.257	7.441	.000	.810	.317	.145	.318	3.145
LetterofRecommendation Strength	.021	.004	.140	5.544	.000	.645	.242	.108	.594	1.685
4 (常量)	-1.502	.090		-16.648	.000					
CGPA	.125	.009	.534	13.227	.000	.882	.511	.255	.229	4.368
GREScore	.002	.000	.195	5.011	.000	.810	.220	.097	.245	4.081
LetterofRecommendation Strength	.021	.004	.135	5.372	.000	.645	.235	.104	.591	1.692
TOEFLScore	.003	.001	.125	3.315	.001	.792	.147	.064	.264	3.794

a. 因变量: ChanceofAdmit

VIF 值一般大于等于 10 时代表自变量之间有强烈的多重共线性。这里自变量之间的多重共线性并不强，而且系数均为正数，符合直觉。

使用向前选择的方法逐步回归，得到以下的模型结果：

输入 / 移去的变量^a

模型	输入的变量	移去的变量	方法
1	CGPA	.	向前 (准则: F-to-enter 的概率 <= .050)
2	GREScore	.	向前 (准则: F-to-enter 的概率 <= .050)
3	LetterofRecommendationStrength	.	向前 (准则: F-to-enter 的概率 <= .050)
4	TOEFLScore	.	向前 (准则: F-to-enter 的概率 <= .050)

a. 因变量: ChanceofAdmit

模型汇总^e

模型	R	R 方	调整 R 方	标准估计的误差	Durbin-Watson
1	.882 ^a	.779	.778	.06647	.847
2	.894 ^b	.800	.799	.06331	
3	.901 ^c	.811	.810	.06150	
4	.903 ^d	.815	.814	.06089	

a. 预测变量: (常量), CGPA。

b. 预测变量: (常量), CGPA, GREScore。

c. 预测变量: (常量), CGPA, GREScore, LetterofRecommendationStrength。

d. 预测变量: (常量), CGPA, GREScore, LetterofRecommendationStrength, TOEFLScore。

e. 因变量: ChanceofAdmit

Anova^a

模型		平方和	df	均方	F	Sig.
1	回归	7.740	1	7.740	1751.850	.000 ^b
	残差	2.200	498	.004		
	总计	9.940	499			
2	回归	7.948	2	3.974	991.448	.000 ^c
	残差	1.992	497	.004		
	总计	9.940	499			
3	回归	8.064	3	2.688	710.750	.000 ^d
	残差	1.876	496	.004		
	总计	9.940	499			
4	回归	8.105	4	2.026	546.547	.000 ^e
	残差	1.835	495	.004		
	总计	9.940	499			

a. 因变量: ChanceofAdmit

b. 预测变量: (常量), CGPA。

c. 预测变量: (常量), CGPA, GREScore。

d. 预测变量: (常量), CGPA, GREScore, LetterofRecommendationStrength。

e. 预测变量: (常量), CGPA, GREScore, LetterofRecommendationStrength, TOEFLScore。

模型通过前向选择了四个变量，排除了个人陈述的得分，这说明了个人陈述

对申请结果的影响相较于其他指标影响不大，可能的解释是：个人陈述与其他自变量有信息量上的重复，比如某位同学在个人陈述中突出了自己的英语能力，而这点从 GRE 与托福成绩中也可以看出来。从显著性分析来看，模型明显优于空模型（至少比不做分析要强），并且所有被选择的自变量共可以解释约 81% 的因变量的变化，因此所有被选择的指标是影响申请成功概率的主要因素，该线性模型是有参考价值的。下面观察各个自变量的系数情况：从未标准化系数来看，托福分数与 GRE 分数的系数很低，然而这并不能代表他们的重要程度低，因为托福和 GRE 成绩在数据范围上远大于其他指标。我们在预测申请成功概率的时候使用未标准化的系数，而观察各个自变量对因变量的影响时应该依据标准化后的系数。经过标准化后，回归系数从大到小依次是在校成绩，GRE 成绩，推荐信得分以及托福成绩。在相关性分析中，我们得出的相关性强度排名从高到低分别是在校成绩，GRE 成绩，托福成绩，推荐信得分。通过多元回归分析与相关性分析，我们交叉印证了在校成绩与 GRE 成绩对于提升申请成功率的强作用，然而在托福成绩与推荐信得分哪个更能影响申请成功率这个问题上相关性分析与多元线性回归给出了相反的结果，我觉得这可能是数据量不够以及在做多元线性分析时“满足正态分布”的前提假设不一定成立的缘故。

四、结论与展望：

本分析报告的主要成果是探究了学生申请指标的分布规律以及对研究生申请成功率的影响。在分布上，约 51%-60% 的学生拥有至少一段科研经历，在五个学校档次中，位于中间三个档次的学生占总体比例最高。托福成绩，推荐信得分与个人陈述得分呈左偏分布，而 GRE 成绩，在校成绩以及申请成功的概率均呈

正态分布。针对托福成绩，分数高于 105 分的同学高出小于等于 105 分的同学平均约 16%-19% 的申请成功概率，而推荐信得分高于 3.5 的同学对小于等于 3.5 分的同学平均高出 11%-16% 的申请成功概率。相关分析结果显示数据集中的所有指标均和申请成功概率有着显著的线性关系，相关系数排名从高到低依次是在校成绩，GRE 成绩，托福成绩，学校档次，个人陈述，推荐信以及科研经历。进一步对排名第一的在校成绩完成方差分析后发现，每提升一个档次的学校成绩对申请成功概率的提升都是有统计学上的显著意义的。经过多元线性回归后，我们发展了一套可以根据学生指标预测申请成功概率的拟合模型，根据标准化后的系数，在校成绩与 GRE 成绩对申请成功概率的影响最突出。

本报告对期望出国的学生的指导意义在于：在精力有限的条件下，应将大部分精力放在提升在校成绩以及英语标准化考试的成绩上，同时应该根据均值与比例参数区间估计的结果明确自己在某个指标上可能存在的劣势并加以重视，着手改善。

本报告的优点在于主题选择切中了广大学生的痛点，并且从数据分析的角度挖掘价值，佐证了市场上一些关于留学准备的流行观点。报告略显不足之处在于数据量较小，这会导致正态性检验时的偏差以及参数估计时的误差。在未来，随着数据的完善与补充，该报告所采用的研究方法会具有更高的可靠性和参考价值。

第6章到第11章节部分作业:

第6章:

1. 试简述 χ^2 分布, t 分布与 F 分布之间的关系.
卡方分布由正态分布随机变量的平方和计算得来, $\chi^2 = \sum_{i=1}^n x_i^2$
 t 分布建立在卡方与正态分布的基础上: $T = \frac{\bar{x}}{\sqrt{s/n}}$
而 F 分布也建立在 χ^2 分布的基础上: $F = \frac{\chi^2/m}{\chi^2/n}$

2. 中心极限定理的意义?

中心极限定理指出大量随机变量近似服从于正态分布. 为数理统计学与误差分析的基础定理.

第7章

1. 解释 95% 的置信区间.

该置信区间以 95% 的概率覆盖参数的真值.

2. $Z_{\alpha/2}$ 的含义是什么.

是指边际误差. 其前提条件是统计量服从正态分布.

3. 简述样本量与置信水平, 总体方差, 估计误差的关系.

置信水平越大, 样本量 (所需) 越大. 样本量越大, 总体方差越大, 估计误差越大, 所需样本量越少.

第8章

1. 两类错误之间存在什么数量关系?

往往控制了一类错误之后另一类错误发生的概率会上升.

2. p 值如何解释?

p 值代表“发生比该数据更极端的情况”的概率. 更小, 越说明数据的出现不是偶然, 而是有统计学意义的.

第9章

1. 简述计算 χ^2 统计量的步骤.

1° 计算每个单元格内的期望频数. (计算时可能用到比例的计算)

2° 计算公式: $\chi^2 = \sum \frac{(f_o - f_e)^2}{f_e}$

2. 简述 ϕ 系数, C 系数, V 系数各自的特点.

ϕ 系数主要用于 2×2 列联表的相关性测量. 数值越小, 相关性越弱.

C 系数主要用于大于 2×2 列联表的相关性测量. $C=0$ 时变量互相独立.

V 系数在 $r=2$ 或 $C=2$ 时就是 ϕ 系数.

第10章

1. 方差分析中有哪些基本假定.

1° 数据间互相独立

2° 服从正态分布

3° 方差齐性.

2. 组内误差和组间误差的含义.

组内误差反应了随机因素对观测数据的影响.

组间误差反应了不同处理对观测数据造成的误差.

3. 什么是交互作用?

因子与因子之间放在一起时产生对因变量的额外的效果.

第11章

1. 为什么要对相关系数进行检验.

因为要杜绝虚假相关的现象. 该现象可能由样本量过小 ~~造成~~ 或者随机误差造成.

2. 在回归分析中, F 检验与 t 检验分别有什么用?

t 检验是对单个变量系数的显著性检验.

F 检验是对整体模型是否有意义的显著性检验.