

# 《数据科学与数据分析》实验课作业

## Week 11\_Exercise 5

191820019 陈文杰

**数据来源：**请下载课件中的“作业数据”。

**要求：**请对“作业数据”中的“专利清洗作业”数据进行清洗，实现以下功能，并回答下列问题。

**功能：**

(1) 清洗数据，利用申请号将申请表中的信息和授权表中的信息匹配起来，形成一张新的表，这张表要对应申请和授权两方面的信息。至于申请未授权，以及授权无申请信息的，酌情处理用于分析，无限制条件。

(2) **最少要实现**如课件第4章“南京专利数据清洗”案例中出现的所有统计工作，并绘制图表，可用于回答下列问题。

**问题：**

(1) 怎样评估江北地区的核心技术优势？

(2) 结合江北的部分产业规划，对江北地区的专利引进工作给出一定建议。**【简要作答，500字以内】**

(3) 数据清洗过程中遇到了哪些问题，你是如何解决的？谈一谈你对数据清洗工作的心得。**【最少200字】**

**格式与内容要求：**

(1) 回答问题时，请先使用 Word 作答，另存为 PDF 文件，**提交 PDF 版本**，上传到“教学立方”；**（只提交作业，不要提交数据！）**

(2) 回答问题**要有逻辑，条理清晰，注意格式**，各级标题要有区分。不标序号且写成一篇作文或日记的作业分数会较低。

# 一、数据清洗与图表绘制

## 1、数据清洗

```
library(dplyr)

library(xlsx)

library(readxl) #不需要 java 环境 read_xlsx

apply.data = read_xlsx("E:\\Desktop\\Week 11_数据整理作业数据 (1)\\专利清洗作业.xlsx",2)
authorize.data = read_xlsx("E:\\Desktop\\Week 11_数据整理作业数据 (1)\\专利清洗作业.xlsx",4)

#1 数据清洗

#1.1 缺失值处理

which(is.na(apply.data),arr.ind=T) #which 函数查看 df 总体是否有缺失值，arr.ind 返回坐标
sum(is.na(apply.data)) #求缺失值元素个数总和
sum(complete.cases(apply.data)) #求非空行数总和
table(is.na(apply.data$申请号)) #根据某列统计缺失的行数
table(complete.cases(apply.data$专利类型)) #True 代表非缺失
mean(is.na(apply.data$申请号)) #查看缺失值占比
#循环遍历，查询每列缺失值占比
for (i in colnames(apply.data)){
  print(table(is.na(apply.data[,i])))
  print(mean(is.na(apply.data[,i])))
}

apply.data <- select(apply.data,申请号:申请日)
authorize.data[!complete.cases(authorize.data),] #列出有缺失值的行
nrow(authorize.data[!complete.cases(authorize.data),]) #计算有缺失值的样本量

#1.2 数据表合并

# merge(x, y, by = intersect(names(x), names(y)),
#       by.x = by, by.y = by, all = FALSE, all.x = all, all.y = all,
#       sort = TRUE, suffixes = c(".x",".y"),
#       incomparables = NULL, ...)

data = merge(apply.data,authorize.data,by = "申请号")

str(data)

#1.2 数据去重

index <- duplicated(data$申请号)

data[!index,]

write.csv(data,"E:\\Desktop\\data.csv")
```

## 2、统计分析与图表绘制

因图表过多，详见本 word 文档附录部分。

二、简答题

(1) 怎样评估江北地区的核心技术优势？

行业名称	主分类号前三位	区属: x2Y				总计
		IV 高淳区	鼓楼区	建邺区	江北新区	
测量仪器	G01	14	125	66	359	564
一般的物理或化学的方法或装置	B01	5	38	18	138	199
发电、变电或配电 H02 131 22 6	H02	6	42	24	119	191
建筑物	E04	16	40	21	101	178
基本电路元件	H01	3	25	13	123	164
输送; 包装; 贮存; 搬运薄的或细丝	B65	11	19	23	107	160
机床; 其他类目中不包括的金属加工	B23	3	4	4	111	122
水、废水、污水或污泥的处理	C02	6	32	19	52	109
电通信技术	H04	3	14	13	46	76
家具; 家庭用的物品或设备; 咖啡磨;	A47	1	2	1	43	47
香料磨; 一般吸尘器						
医学或兽医学; 卫生学	A61		20	1	20	41
总计		68	361	203	1219	1851

上表是针对特定的几个行业，进行区属之间的专利授权统计所得的结果，据此对江北地区的核心技术优势进行评估，得出以下几点结论：

①江北地区技术优势明显。可以看到，江北新区在精密仪器、生化、电力、水利、通信等行业的专利授权量明显高于其他的几个区属，具有较为明显的技术优势。

②产业分布不合理。这些技术仍然是偏向服务于传统产业，在这些领域拥有核心技术优势不代表其在新材料、新能源以及互联网等领域有较强竞争力，而新兴产业又是当今经济的主要增长点和驱动力，所以江北新区在保持现有技术优势的同时，也应考虑加大对新兴产业的投入，促进专业的提质升级。

③样本分布的影响。值得注意的是，江北新区包含了浦口区、六合区和栖霞区八卦洲街道，但是与之进行对比的都是单独的分区，所以即使江北新区总体核心技术优势明显，下属分区技术优势仍有待考量，基于此应关注新区内部的统一、协调发展，掌握核心技术、增强区属竞争力。

(2) 结合江北的部分产业规划，对江北地区的专利引进工作给出一定建议。

①江北新区的产业规划

以往的江北地区主要承接江南的产业专业为主，容纳了许多夕阳产业，对江北地区的生态环境以及经济可持续发展造成了许多负面影响，因此在政府关于本轮打造江北新区的政策文件中，特意提出要着眼长远，对江北新区的产业进行高标准规划，实行清洁生产。

②专利引进工作建议

于此，江北地区的专利引进应着眼于包括互联网、新能源、新材料、通讯等

高新技术，为江北地区的产业提质升级注入生机，争取早日形成以新能源、信息技术、生物医药、汽车轨道交通以及节能环保在内的五大先进制造业，建立并完善现代产业体系。

**（3）数据清洗过程中遇到了哪些问题，你是如何解决的？谈一谈你对数据清洗工作的心得。**

#### **#数据清洗过程中的问题与解决方法**

##### ①数据不完整。

首先进行数据概览，并分析缺失值在行列上的集中分布，进而决定采用适合的方式进行处理（删除、填充...）。在本次实验中，缺失值主要存在于部分记录的部分属性中，缺失值较少，而且这些属性值大多无关后续的统计分析，故未对缺失值进行删除或填充。

##### ②数值不匹配。

以专利主题分类号为例，它大致分为两类，一类带字母，一类不带字母而是纯数字，难以进行合并分析。故通过人工查询，将纯数字的分类号转化为带字母的分类号，实现数值类型的统一，利于后续的专利类型统计分析。

##### ③数据重复。

数据集难免出现重复数据，因而需要时刻注意数据的去重。

##### ④数据无意义

在数据合并时如果未能选择合理的合并方式，可能得到笛卡尔积的数据统计表，里面的很多记录都无实际意义，因而选定特定的连接列、选用合适的合并方式至关重要。

#### **#数据清洗心得**

①数据清洗往往伴随数据采集进行，是完整数据分析流程的靠前期的工作，也是较为关键的部分，重要性不言而喻。

②数据清洗较为繁琐，一方面需要考虑数据的外况，看数据结构是否合理、数据缺失值是否有碍数据分析、数据是否有重复；另一方面需要结合现实，考察数据记录、数据值是否有意义，考察数据的内涵。

③在数据清洗阶段充分利用可视化图表，一方面提高清洗的效率和准度，另一方面可以加深对数据集的认知，为后续的模型搭建与研究分析奠定基础。

【附录：实验结果】

图 1：专利申请数区属分布条形图（1）

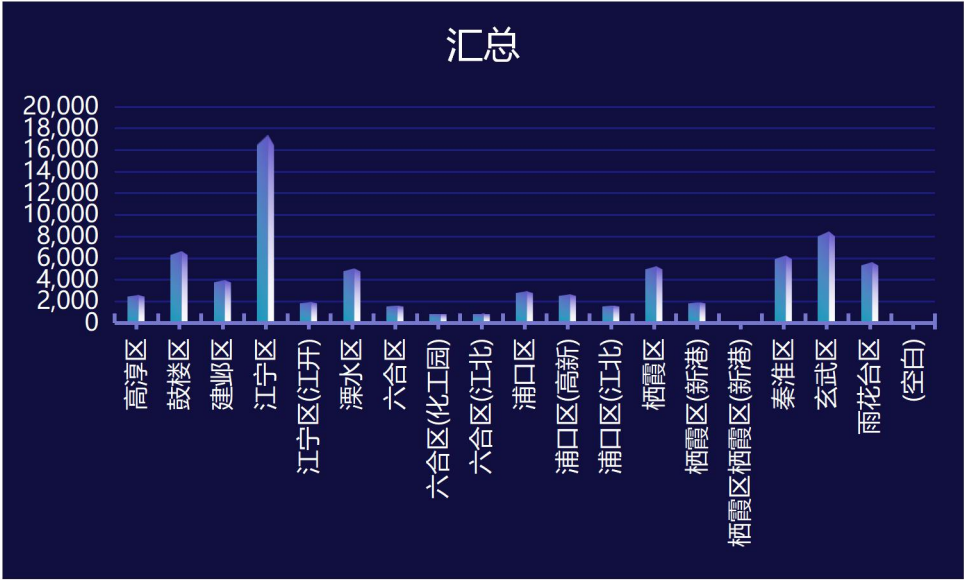


图 2：专利申请数区属分布条形图（降序）

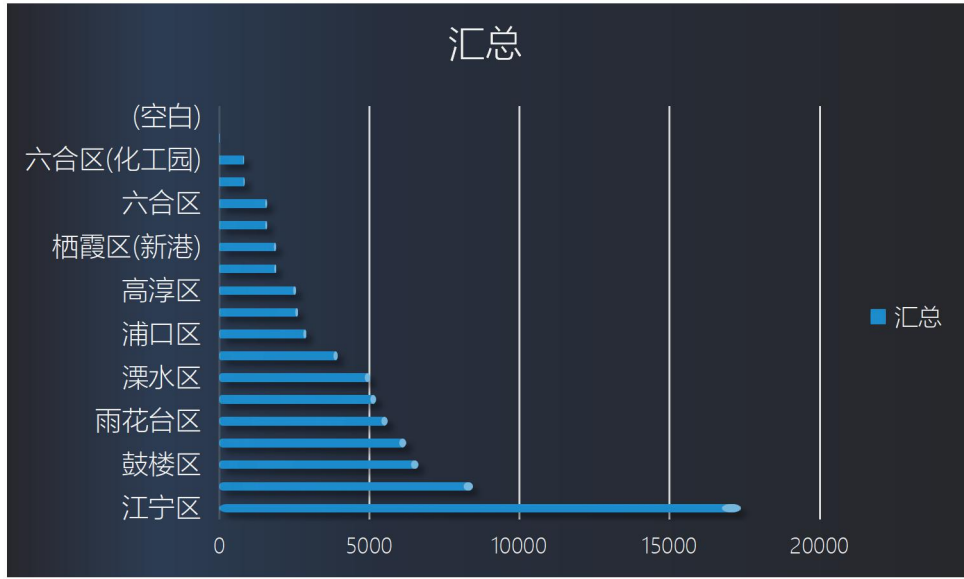


图 3：专利申请数区属分布饼状图（显示百分比，部分区属未组合）

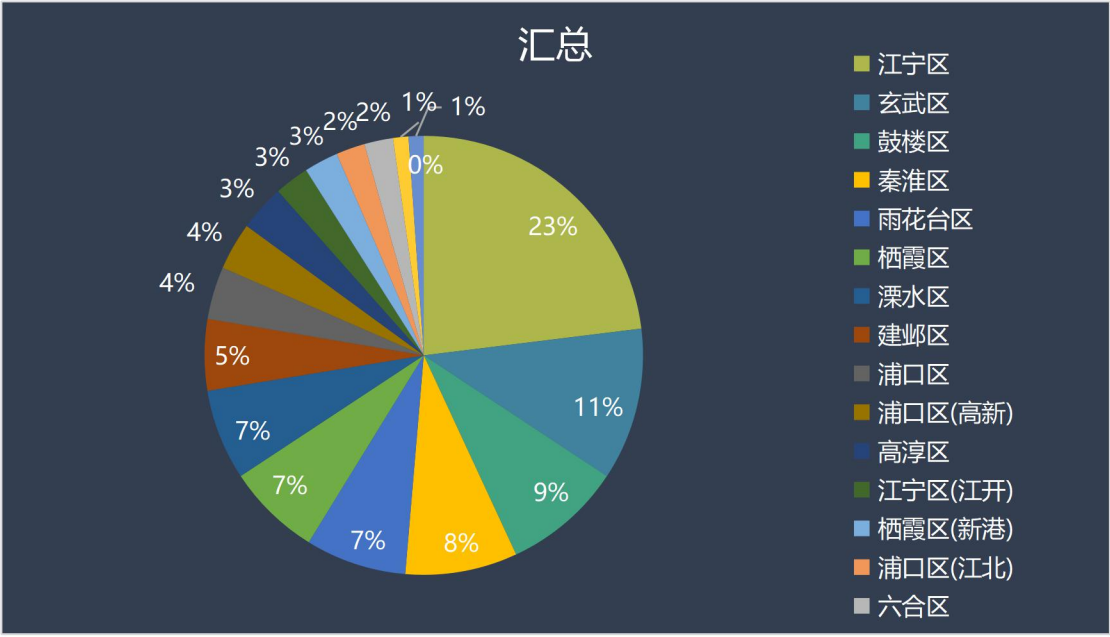


图 4：专利申请数区属分布饼状图（显示数量）

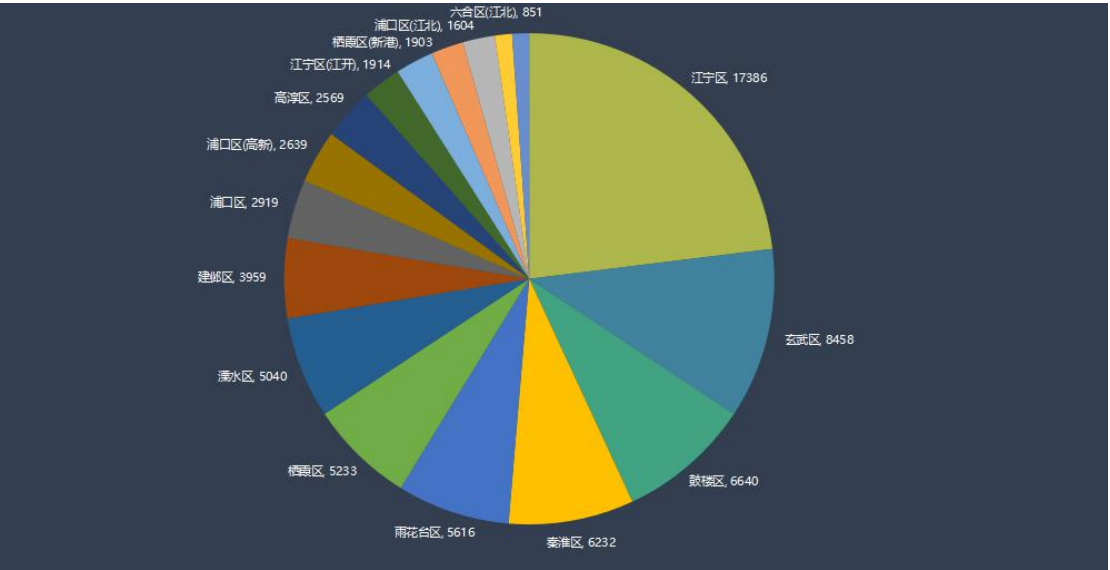


图 5：专利申请数区属分布饼状图（显示百分比）

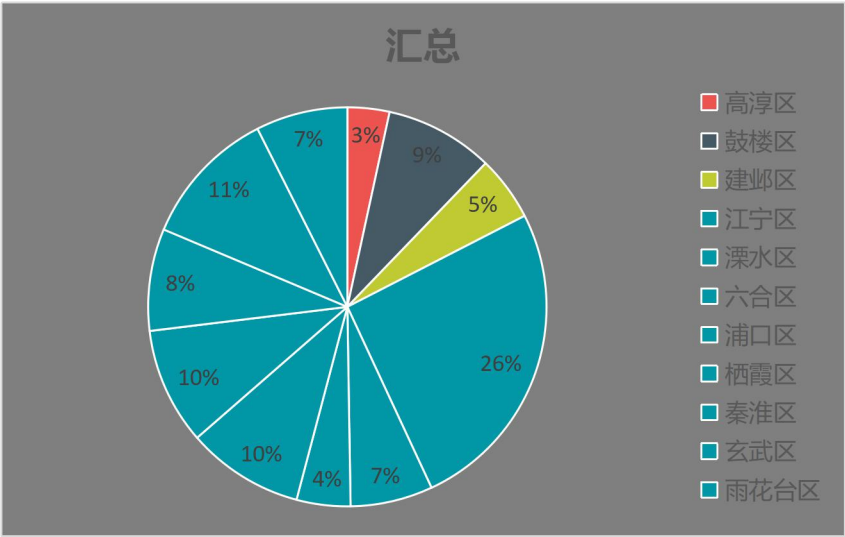


图 6：专利申请数区属分布条形图（部分区属汇总）

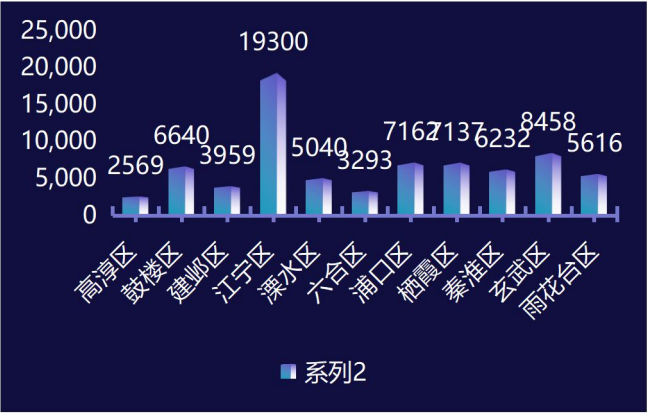
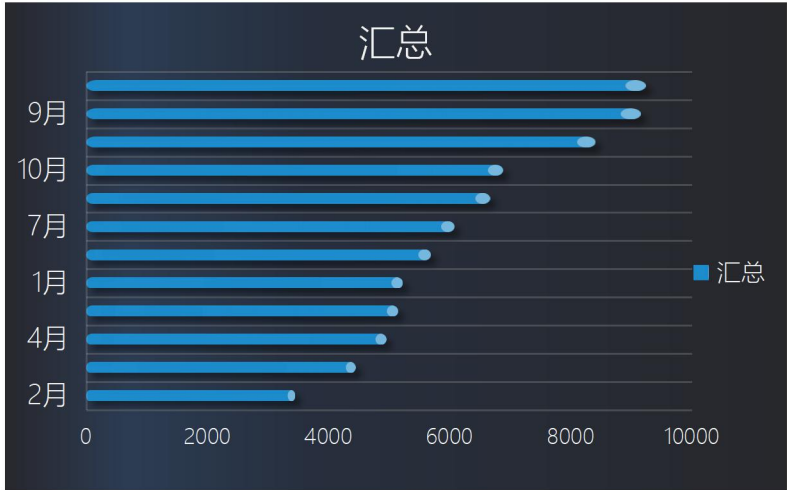


图 7：专利申请月份统计条形图（部分区属汇总）



计数项·区属	专利类型				
区属	发明	实用新型	外观设计	(空白)	总计
高淳区	1526	838	205		2569
鼓楼区	3404	2733	503		6640
建邺区	1388	1526	1045		3959
江宁区	7449	7140	2797		17386
江宁区(江开)	508	621	785		1914
溧水区	2131	2051	858		5040
六合区	514	938	149		1601
六合区(化工园)	301	401	49		841
六合区(江北)	460	326	65		851
浦口区	1348	1472	99		2919
浦口区(高新)	1152	1349	138		2639
浦口区(江北)	695	854	55		1604
栖霞区	2741	2171	321		5233
栖霞区(新港)	778	1036	90		1904
秦淮区	3806	2068	358		6232
玄武区	6080	1949	429		8458
雨花台区	2915	2168	533		5616
(空白)					
总计	37286	29641	8479		75406

表 1：专利类型统计表

计数项·区属	申请人类 型						
区属	大专院校	个人	机关团体	科研机构	企业	(空白)	总计
高淳区	12	212	5	7	2333		2569
鼓楼区	2026	1029	593	640	2352		6640
建邺区	343	345	23	51	3197		3959
江宁区	5502	3241	142	185	8316		17386
溧水区	197	400		74	4369		5040
六合区	22	315	31		1233		1601
六合区(化工园)	57	25	4	15	740		841
六合区(江北)	172	36			643		851
浦口区	701	924	6	4	1284		2919
浦口区(高新)	808	84			1747		2639
浦口区(江北)	637	274		27	666		1604
栖霞区	2179	629	31	141	2253		5233
栖霞区(新港)	1	38	1	29	1835		1904
秦淮区	1988	1599	227	290	2128		6232
玄武区	4055	1950	316	1108	1029		8458
雨花台区	17	762	315	268	4254		5616
总计	18717	11866	1694	2903	40226		75406

表 2：申请人类型统计表（数值）



计数项:申请人类型	申请人类型				
区属	大专院校	机关团体	科研机构	企业	总计
高淳区	0.47%	0.19%	0.27%	90.81%	100.00%
鼓楼区	30.51%	8.93%	9.64%	35.42%	100.00%
建邺区	8.66%	0.58%	1.29%	80.75%	100.00%
江宁区	31.65%	0.82%	1.06%	47.83%	100.00%
江宁区(江开)	0.00%	0.00%	3.34%	96.50%	100.00%
溧水区	3.91%	0.00%	1.47%	86.69%	100.00%
六合区	1.37%	1.94%	0.00%	77.01%	100.00%
六合区(化工园)	6.78%	0.48%	1.78%	87.99%	100.00%
六合区(江北)	20.21%	0.00%	0.00%	75.56%	100.00%
浦口区	24.02%	0.21%	0.14%	43.99%	100.00%
浦口区(高新)	30.62%	0.00%	0.00%	66.20%	100.00%
浦口区(江北)	39.71%	0.00%	1.68%	41.52%	100.00%
栖霞区	41.64%	0.59%	2.69%	43.05%	100.00%
栖霞区(新港)	0.05%	0.05%	1.52%	96.38%	100.00%
秦淮区	31.90%	3.64%	4.65%	34.15%	100.00%
玄武区	47.94%	3.74%	13.10%	12.17%	100.00%
雨花台区	0.30%	5.61%	4.77%	75.75%	100.00%
总计	24.82%	2.25%	3.85%	53.35%	100.00%

表 3：申请人类型统计表（横向百分比）

计数项:申请人类型	申请人类型					
区属	大专院校	个人	机关团体	科研机构	企业	总计
高淳区	0.06%	1.79%	0.30%	0.24%	5.80%	3.41%
鼓楼区	10.82%	8.67%	35.01%	22.05%	5.85%	8.81%
建邺区	1.83%	2.91%	1.36%	1.76%	7.95%	5.25%
江宁区	29.40%	27.31%	8.38%	6.37%	20.67%	23.06%
江宁区(江开)	0.00%	0.03%	0.00%	2.20%	4.59%	2.54%
溧水区	1.05%	3.37%	0.00%	2.55%	10.86%	6.68%
六合区	0.12%	2.65%	1.83%	0.00%	3.07%	2.12%
六合区(化工园)	0.30%	0.21%	0.24%	0.52%	1.84%	1.12%
六合区(江北)	0.92%	0.30%	0.00%	0.00%	1.60%	1.13%
浦口区	3.75%	7.79%	0.35%	0.14%	3.19%	3.87%
浦口区(高新)	4.32%	0.71%	0.00%	0.00%	4.34%	3.50%
浦口区(江北)	3.40%	2.31%	0.00%	0.93%	1.66%	2.13%
栖霞区	11.64%	5.30%	1.83%	4.86%	5.60%	6.94%
栖霞区(新港)	0.01%	0.32%	0.06%	1.00%	4.56%	2.52%
秦淮区	10.62%	13.48%	13.40%	9.99%	5.29%	8.26%
玄武区	21.66%	16.43%	18.65%	38.17%	2.56%	11.22%
雨花台区	0.09%	6.42%	18.60%	9.23%	10.58%	7.45%
总计	100.00%	100.00%	100.00%	100.00%	100.00%	100.00%

表 4：申请人类型统计表（纵向百分比）

**表 5：专利分类统计表（数值）**

表 6: 专利分类统计表 (纵向百分比)

表 7: 专利分类统计表 (横向百分比)

表 8: 专利分类按区属统计表 (1)

主分类号前三位		区属: x															
主分类号前三位		(1) 高淳区		高淳区	溧水区	江宁区	六合区	六合区(化工园)	浦口区	浦口区(高新)	鼓楼区	鼓楼区(新浦)	秦淮区	玄武区	雨花台区	(空白)	总计
01	14	120	66	287	23	89	24	44	23	83	21	80	30	22	1024		
02	6	42	24	158	8	18	21	31	39	20	11	39	19	12	447		
03	11	19	23	117	18	20	10	29	13	12	26	21	9	14	146		
04	5	38	18	90	10	31	36	23	18	19	13	10	10	10	339		
05	16	40	21	73	2	15	3	23	1	54	1	29	19	6	312		
06	3	23	13	89	14	11	1	36	21	29	18	16	7	9	306		
07	8	14	22	102	13	13	1	14	2	17	7	17	9	7	243		
08	3	4	4	11	18	26	1	27	15	18	14	7	6	6	202		
09	1	12	6	47	10	13	1	14	5	3	4	83	8	2	213		
10	3	14	13	55	3	4	6	8	21	11	2	16	14	40	207		
11	6	32	18	43	7	6	14	12	3	5	8	19	11	1	186		
12	1	27	17	26	22	3	1	3	1	9	4	18	8	1	186		
13	1	10	9	23	18	3	3	8	1	12	6	21	18	6	182		
14	2	13	13	53	2	2	6	17	13	2	18	13	13	6	182		
15	8	13	13	37	3	3	4	18	11	6	14	8	8	9	149		
16	3	5	6	24	7	3	1	7	8	9	8	9	11	7	150		
17	17	6	4	24	1	3	1	10	23	9	4	19	3	7	137		
18	6	13	10	28	12	13	1	7	20	12	1	6	3	5	117		
19	6	6	17	13	13	13	2	11	2	13	1	9	9	2	115		
20	2	1	18	50	7	7	1	4	3	4	6	3	3	3	106		
21	2	7	4	31	1	3	3	3	7	1	2	11	3	6	103		
22	2	9	8	41	1	3	2	6	10	6	4	4	3	3	103		
23	1	9	2	29	7	2	2	7	1	9	2	4	5	9	95		
24			1	40	18	14		4	3	1		22	12	2	94		
25	20			10	4			3	10	2	1	22	12	2	84		
26	2	1	1	11	1			6	6	9	8	6	7	4	81		

表 9：专利分类按区属统计表（2）

行业名称	计数项:主分类号前三位	区属: x2	高淳区	鼓楼区	建邺区	江北新区	总计
测量仪器	G01		14	125	66	359	564
一般的物理或化学的方法或装置	B01		5	38	18	138	199
发电、变电或配电 H02 131 22 6	H02		6	42	24	119	191
建筑物	E04		16	40	21	101	178
基本电器元件	H01		3	25	13	123	164
输送；包装；贮存；搬运薄的或细丝	B65		11	19	23	107	160
机床；其他类目中不包括的金属加工	B23		3	4	4	111	122
水、废水、污水或污泥的处理	C02		6	32	19	52	109
电通信技术	H04		3	14	13	46	76
家具；家庭用的物品或设备；咖啡磨；香料磨；一般吸尘器	A47		1	2	1	43	47
医学或兽医学；卫生学	A61			20	1	20	41
总计			68	361	203	1219	1851

表 10：专利行业归属统计表

区属: x2	计数项:主分类号前三位	计数项:主分类号前三位
高淳区	272	2.59%
鼓楼区	894	8.51%
建邺区	577	5.49%
江宁区	3389	32.26%
江宁区(江开)	6	0.06%
溧水区	659	6.27%
江北新区	2847	27.10%
六合区(江北)	3	0.03%
浦口区(江北)	1	0.01%
秦淮区	811	7.72%
玄武区	645	6.14%
雨花台区	402	3.83%
总计	10506	100.00%

表 11：特定产业（电通信技术）统计表（数值+百分比）