Анализ номинативных

данных

Кто чаще обращается в службу знакомств: мужчины или женщины?

Можно ли утверждать, что водители-женщины чаще становятся участниками ДТП

Зависит ли количество аварий на производстве от дня недели?

(дорожно-транспортных происшествий)?

игрышей?

Можно ли утверждать, что выигрыши в игре распределены не случайно среди про-

	Распределение:							
	эмпирическое	теоретическое						
«3a»	30	25						
«Против»	20	25						
Сумма:	50	50						

Хи-квадрат

$$\chi_{9}^{2} = \sum_{i=1}^{P} \frac{(f_{9} - f_{T})^{2}}{f_{T}}$$

```
import org.apache.spark.ml.stat.ChiSquareTest
val data = Seq(
  (0.0, Vectors.dense(0.5, 10.0)),
  (0.0, Vectors.dense(1.5, 20.0)),
  (1.0, Vectors.dense(1.5, 30.0)),
  (0.0, Vectors.dense(3.5, 30.0)),
  (0.0, Vectors.dense(3.5, 40.0)),
  (1.0, Vectors.dense(3.5, 40.0))
val df = data.toDF("label", "features")
val chi = ChiSquareTest.test(df, "features", "label").head
println("pValues = " + chi.getAs[Vector](0))
println("degreesOfFreedom = " + chi.getSeq[Int](1).mkString("[", ",", "]"))
println("statistics = " + chi.getAs[Vector](2))
```

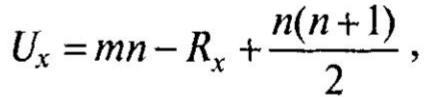
import org.apache.spark.ml.linalg.{Vector, Vectors}

Непараметрические методы

U - Манна - Уитни

Обозначим значения переменной для одной выборки X, а для другой выборки — Y и упорядочим значения обеих выборок по возрастанию.

Значения	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	19
Выборка	X	X	Y	X	X	X	Y	X	X	Y	X	Y	Y	Y	Y	Y



 $U_y = mn - R_y + \frac{m(m+1)}{2},$

 $U_x + U_y = mn$

<u>-</u>	Ранги <i>X</i>	$\frac{1}{1}$	2		4	5	6		8	9	10	11	1	10	* '		
<u>.</u> 5	Ранги У	\vdash		3				7			10		12	13	14	15	16

Шаг 2. Значения выборок ранжируются, и выписываются отдельно ранги для од-

Шаг 3. Вычисляются суммы рангов по $X(R_x)$ и по $Y(R_y)$: $R_x = 46$; $R_y = 90$.

Значения

той и другой выборке (строки 1 и 2).

ной и другой выборки (строки 3-5).

Ш а г 4. Вычисляются U_x и U_y по формуле 12.1:

$$\frac{1}{4} = \frac{1}{4} + \frac{1}$$

$$U_x = 8 \cdot 8 - 46 + \frac{8(8+1)}{2} = 54$$
, $U_y = 8 \cdot 8 - 90 + \frac{8(8+1)}{2} = 10$, $U_x + U_y = 64 = mn$.

Ш а г 5. Определяется p-уровень значимости: наименьшее из U сравнивается с табличным (приложение 9) для соответствующих объемов выборки m=8 и n=8. Значение $p \le 0.05$ (0.01), если вычисленное $U_{\text{змп}} \le U_{\text{габл}}$ В нашем случае наименьшим является $U_{\nu} = 10$, которое и принимается за эмпирическое значение критерия. Оно меньше критического для p=0.05 (U=13), но больше критического для p=0.01(U = 7). Следовательно, p < 0.05.

Т - Вилкоксон

Проверим гипотезу о различии значений показателя, измеренного дважды на одной и той же выборке («Условие 1» и «Условие 2»), на уровне $\alpha = 0.05$:

1	№ объекта:	1	2	3	4	5	6	7	8	9	10	11	12
2	Условие 1:	6	11	12	8	5	10	7	6	3	9	4	5
3	Условие 2:	14	5	8	10	14	7	12	13	11	10	15	16
4	Разность d_i :	-8	6	4	-2	<u>-9</u>	3	-5	- 7	-8	-1	-11	-11
5	Ранги d _i :	8,5	6	4	2	10	3	5	7	8,5	1	11,5	11,5
6	Ранги $d_i(+)$:		6	4			3						
7	Ранги <i>d</i> (-):	8,5			2	10		5	7	8,5	1	11,5	11,5

Ш а г 2. Ранжировать абсолютные значения разностей (строка 5). Ш а г 3. Выписать ранги положительных и отрицательных значений разностей (стро-

Ш а г 1. Подсчитать разности значений для каждого объекта выборки (строка 4).

ки 6 и 7). Ш а г 4. Подсчитать суммы рангов отдельно для положительных и отрицательных

ш а г 4. Подсчитать суммы рангов отдельно для положительных и отрицательных разностей: $T_1 = 13$; $T_2 = 65$. За эмпирическое значение критерия $T_{\text{эмп}}$ принимается меньшая сумма: $T_{\text{эмп}} = 13$.

Ш а г 5. Определяется p-уровень значимости: $T_{\text{эмп}}$ сравнивается с табличным (приложение 10) для соответствующего объема выборки. Значение $p \le 0.05$ (0,01), если вычисленное $T_{\text{эмп}} \le T_{\text{табл}}$ В нашем случае эмпирическое значение равно критическому значению для p = 0.05. Следовательно, p = 0.05.

Экспериментальные

планы

AB

ABn

$1 - (1 - \alpha)^m$

m=5, lpha=0,05 она равна pprox 22,6%

Поправка Бонферрони

$p_i < \alpha/m$

AABB

A-A+B-B+

Многорукие бандиты



Биномиальное распределение

$$f(k,n,p)=\Pr(k;n,p)=\Pr(X=k)=inom{n}{k}p^k(1-p)^{n-k}$$

Бернулли

$$f_a(y | \theta_a) = \theta_a^y (1 - \theta_a)^{1-y}$$

$$egin{array}{ll} heta_i & \sim & \operatorname{Beta}\left(lpha_i,eta_i
ight) \ y_i & \sim & \operatorname{Bernoulli}\left(heta_i
ight) \end{array}$$

модель легко интерпретируема, α — это количество успешных испытаний, а β — количество неуспешных испытаний; среднее значение будет $\frac{\alpha}{\alpha+\beta}$.