

MATH 425

Note Title

11/11/2011

11/11/2011

The Pearson χ^2 test. We draw N balls with return. The balls are of k colors color i occurs with probability p_i . We got a_i balls of color i .

$$\chi^2 = \frac{(a_1 - p_1 N)^2}{p_1 N} + \dots + \frac{(a_k - p_k N)^2}{p_k N}$$

Claim: This has the distribution of a sum of squares of $(k-1)$ independent standard Gaussian variables.

To prove the claim, note that we have a sum of k random variables, but they are not independent. (For example, if $k=2$, a_2 is determined by a_1 .)

How do we deal with sums of λ variables which are not independent? How to quantify the lack of independence. Recall we observed that if X and Y are independent then

$$E(X \cdot Y) = E(X) \cdot E(Y)$$

We define the covariance of X and Y as

$$\text{cov}(X, Y) = E(X \cdot Y) - E(X) \cdot E(Y).$$

Note that

$$\text{var}(X) = \text{cov}(X, X).$$

This is kind of analogous to the dot product

of vectors: $\vec{x} = (x_1, \dots, x_n) \in \mathbb{R}^n$, $\vec{y} = (y_1, \dots, y_n) \in \mathbb{R}^n$

$$\vec{x} \cdot \vec{y} = x_1 y_1 + \dots + x_n y_n.$$

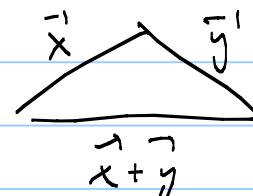
$$\|\vec{x}\|^2 = \vec{x} \cdot \vec{x}$$

$$\begin{aligned} \|\vec{x} + \vec{y}\|^2 - \|\vec{x}\|^2 - \|\vec{y}\|^2 &= (\vec{x} + \vec{y}) \cdot (\vec{x} + \vec{y}) - \vec{x} \cdot \vec{x} - \vec{y} \cdot \vec{y} \\ &= 2\vec{x} \cdot \vec{y} \end{aligned}$$

Similarly in probability,

$$(*) \quad \text{cov}(X, Y) = \frac{1}{2} (\text{var}(X+Y) - \text{var}(X) - \text{var}(Y))$$

What is the angle α between two vectors \vec{x}, \vec{y} ?

$$\cos \alpha = \frac{\vec{x} \cdot \vec{y}}{\|\vec{x}\| \|\vec{y}\|} \quad (\text{theorem of cosines})$$


The diagram shows a triangle formed by three vectors originating from the same point. One vector points upwards and to the left, labeled \vec{x} . Another vector points upwards and to the right, labeled \vec{y} . The third vector, which is the sum of \vec{x} and \vec{y} , points downwards and to the right, labeled $\vec{x} + \vec{y}$.

In probability, this is called the correlation coefficient

$$g(X, Y) = \frac{\text{cov}(X, Y)}{\sigma(X) \sigma(Y)} = \frac{\text{cov}(X, Y)}{\sqrt{\text{var}(X) \text{var}(Y)}}$$

In relevance to the X^2 , we want to look

at the variables

$$\frac{a_1 - p_1 N}{\sqrt{p_1 N}}, \dots, \frac{a_k - p_k N}{\sqrt{p_k N}}.$$

I want to compute all the variances and covariances of these variables.

a_i are binomial variables $X_{p_i, N}$

$$\text{var}(a_i) = N p_i (1-p_i)$$

$$\text{var}\left(\frac{a_i - p_i N}{\sqrt{p_i N}}\right) = \frac{N p_i (1-p_i)}{N p_i} = 1 - p_i$$

$$i \neq j \quad \text{Cov} \left(\frac{a_i - p_i N}{\sqrt{p_i N}}, \frac{a_j - p_j N}{\sqrt{p_j N}} \right) = \frac{1}{\sqrt{p_i N} \sqrt{p_j N}} \text{cov}(a_i, a_j)$$

probabilität $= \frac{1}{N \sqrt{p_i p_j}}$

binomial $\overset{\sim}{\sim} p_i + p_j, N$ trials

$$\text{cov}(a_i, a_j) = \frac{1}{2} (\text{var}(a_i + a_j) - \text{var}(a_i) - \text{var}(a_j))$$

$$= \frac{1}{2} ((p_i + p_j)(1 - p_i - p_j)N - p_i(1 - p_i)N - p_j((1 - p_j)N))$$

$$= -p_i p_j N$$

$$\text{cov} \left(\frac{a_i - p_i N}{\sqrt{p_i N}}, \frac{a_j - p_j N}{\sqrt{p_j N}} \right) = -\sqrt{p_i p_j}$$

if $i \neq j$

$$\text{cov} \left(\frac{a_i - p_i N}{\sqrt{p_i N}}, \frac{a_i - p_i N}{\sqrt{p_i N}} \right) = \text{var} \left(\frac{a_i - p_i N}{\sqrt{p_i N}} \right) = 1 - p_i$$

If we have n vectors $\vec{v}_1, \dots, \vec{v}_n$, it is useful to form the matrix

$$\begin{pmatrix} \vec{v}_1 \cdot \vec{v}_1 & \dots & \vec{v}_1 \cdot \vec{v}_n \\ \vdots & & \vdots \\ \vec{v}_n \cdot \vec{v}_1 & \dots & \vec{v}_n \cdot \vec{v}_n \end{pmatrix} \quad (\text{Gram matrix})$$

In our case, let us make the matrix of covariances of a_1, \dots, a_k : $\begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix}$

$$\begin{pmatrix} 1 - p_1 & -\sqrt{p_1 p_2} & \cdots & -\sqrt{p_1 p_k} \\ -\sqrt{p_2 p_1} & 1 - p_2 & \cdots & -\sqrt{p_2 p_k} \\ \vdots & & & \vdots \\ -\sqrt{p_k p_1} & -\sqrt{p_k p_2} & \cdots & 1 - p_k \end{pmatrix} = I - \left(\begin{pmatrix} \sqrt{p_1} & \cdots & \sqrt{p_k} \end{pmatrix} v_{ij} \right) =$$

$$= I - \underbrace{\left(\begin{pmatrix} \sqrt{p_1} & \cdots & \sqrt{p_k} \end{pmatrix} \cdot \begin{pmatrix} \sqrt{p_1} & \cdots & \sqrt{p_k} \end{pmatrix}^T \right)}_{\text{matrix product}} =$$

$$\begin{pmatrix} \sqrt{p_1} \\ \vdots \\ \sqrt{p_k} \end{pmatrix} \quad k \times 1 \quad k \times k$$

$$\|(\sqrt{p_1}, \dots, \sqrt{p_k})\| = \sqrt{\sqrt{p_1}^2 + \dots + \sqrt{p_k}^2} = \sqrt{p_1 + \dots + p_k} = \sqrt{1} = 1$$

$v = (\sqrt{p_1}, \dots, \sqrt{p_k})$. I can change coordinates in a way preserving the dot product such that v becomes $(1, 0, \dots, 0)$.

So the in those coordinates becomes

$$I - (1, 0, \dots, 0)^T \cdot (1, 0, \dots, 0) = I - \begin{pmatrix} 1 & 0 \\ 0 & 0 \end{pmatrix} =$$

$$= \begin{pmatrix} 0 & 0 \\ 0 & 1 & \dots & 1 \end{pmatrix}$$

here are my unrelated
k-1 independent
standard basisvans

(HW)

Find a χ^2 table on the internet.

①

Suppose a standard die is cast 10

times, we get 3 1's
7 2's
2 3's
2 4's
2 5's
4 6's.

Can you conclude the die is biased
with significance level 95%?

②

If we keep casting the die, how
many 6's in a row (from the first trial)

do we have to throw to conclude
the die is biased with significance level

95%?

[hint: degrees of freedom can be 1]

6 or. not a 6