

Lecture 2

27.10.2019.

Recap: Discrete exponential families are toric varieties and are represented by data A, h where A is a $k \times r$ integer matrix with $\mathbf{1}$ in its row-span and h is a vector in $\mathbb{R}_{>0}^r$.

Example: X, Y binary. Independence model.

$$\mathcal{M}_{X \perp\!\!\!\perp Y} = \left\{ \begin{pmatrix} p_{00} & p_{01} \\ p_{10} & p_{11} \end{pmatrix} : p_{00}p_{11} - p_{01}p_{10} = 0, p_{ij} \geq 0, \sum p_{ij} = 1 \right\}$$

$$= \left\{ \begin{pmatrix} p_{00} & p_{01} \\ p_{10} & p_{11} \end{pmatrix} : \begin{pmatrix} p_{0+} \\ p_{1+} \end{pmatrix} \begin{pmatrix} p_{+0} & p_{+1} \end{pmatrix} \right\}$$

$$A = \begin{pmatrix} 1 & 1 & 0 & 0 \\ 0 & 0 & 1 & 1 \\ 1 & 0 & 1 & 0 \\ 0 & 1 & 0 & 1 \end{pmatrix} \quad h = (1, 1, 1, 1), \phi^{A, h}: \mathbb{R}^k \rightarrow \mathbb{R}^r$$

$$(\theta_1, \theta_2, \theta_3, \theta_4) \mapsto (\theta_1, \theta_3 - \theta_1, \theta_4, \theta_2 \theta_3 - \theta_2 \theta_4)$$

If $u = (u_{00}, u_{01}, u_{10}, u_{11}) \in \mathbb{N}^4$

$$Au = \begin{pmatrix} u_{0+} \\ u_{1+} \\ u_{+0} \\ u_{+1} \end{pmatrix} \text{ captures the marginal counts of the contingency table } \begin{array}{cc|c} u_{00} & u_{01} & u_{0+} \\ u_{10} & u_{11} & u_{1+} \\ \hline u_{+0} & u_{+1} & \end{array}$$

Statistical tests for $\overbrace{\text{DEF}}^{\text{DEF}}$ discrete exponential families.

General set up: Consider a DEF $\mathcal{M}_{A, h} \subseteq \Delta_{r-1}$ and suppose we collect data $X^{(1)}, \dots, X^{(n)} \in [r]$ which are indep. and identic. distributed according to a distribution $p \in \text{int}(\Delta_{r-1})$.

We would like to test

$$H_0: p \in \mathcal{M}_{A, h} \text{ vs. } p \notin \mathcal{M}_{A, h}$$

Example: Data 326 homicide indictments in Florida.

Binary classif: $X = \text{Race} = \begin{cases} B \\ W \end{cases}$, $Y = \text{Penalty} = \begin{cases} D \\ \cancel{X} \end{cases}$

Q: Were decisions of penalty made independent of race?

Penalty Race	D	X	Total
B	19	141	160
W	17	149	166
Total	36	290	326

Hypothesis testing problem: $H_0: p \in \mathcal{M}_{X \perp Y}$ vs. $H_1: p \notin \mathcal{M}_{X \perp Y}$.

Chi-squared test of indep: If H_0 is true $\Rightarrow p_{ij} = p_{i+} p_{+j}$

The expected number of occurrences of $\{X=i, Y=j\}$ is $n p_{i+} p_{+j}$

We can estimate the marginal probabilities by using

$$\hat{p}_{i+} = \frac{U_{i+}}{n} \quad \hat{p}_{+j} = \frac{U_{+j}}{n}, \quad U_{i+} = U_{i0} + U_{i1}$$
$$U_{+j} = U_{0j} + U_{1j}$$

\rightarrow We can estimate the number of counts $n p_{i+} p_{+j}$ by

$$\hat{U}_{ij} = n \hat{p}_{i+} \hat{p}_{+j}$$

The χ^2 -statistic
$$\chi^2(u) = \sum_{i=1}^r \sum_{j=1}^r \frac{(U_{ij} - \hat{u}_{ij})^2}{\hat{u}_{ij}}$$

If H_0 is true we expect a small value for $\chi^2(u)$.
We reject H_0 if $\chi^2(u)$ is too large. How large?

{ The probability that $\chi^2(u)$ takes a value greater than or equal to $\chi^2(u)$ provided that H_0 is true.
p-value

Goal: Compute $P(\chi^2(u) \geq \chi^2(u))$

Asymptotic approach: Suppose $n \rightarrow \infty$ and use

Prop. 9.1.1: If the joint distribution of X is given by a distribution $p \in \mathcal{M}_{A,n}$, then

$$\lim_{n \rightarrow \infty} P(\chi_n^2(v) \geq t) = P(\chi_{df}^2 \geq t) \text{ for all } t \geq 0,$$

where $df = r - 1 - \dim \mathcal{M}_{A,n}$ is the codim. or # of degrees of freedom of the model.

i.e. For a true distribution lying in the model, the Pearson χ^2 statistic converges in distrib. to a χ_{df}^2 .

Drawbacks: What if sample size is small and contingency table is sparse?

Exact Approach: Can we calculate this p-value some other way?

Denote by $[r]$ the outcome space. Assume $h=1$.

$$\text{Let } \mathcal{T}(n) = \left\{ u \in \mathbb{N}^r : \sum_{i \in r} u_i = n \right\}$$

Def 1.1.10: We call the vector Au the minimal sufficient statistics for the model \mathcal{M}_A , and the set of tables

$$\mathcal{F}(u) = \left\{ v \in \mathbb{N}^r : Av = Au \right\}$$

is called the fiber of the contingency table $u \in \mathcal{T}(n)$ w.r.t the model.

$$\mathbf{1} \in \text{rowspan} \Rightarrow xA = \mathbf{1}$$

$$x(Av) = x(Au)$$

$$\mathbf{1}v = \mathbf{1}u \Rightarrow \sum v_i = \sum u_i = n.$$

Prop 1: If $p \in \mathcal{M}_A$, $p(x) = \theta_1^{a_{1x}} \cdots \theta_k^{a_{kx}}$, $x \in [r]$, and $u \in \mathcal{T}(n)$, then

$$P(U=u) = \frac{n!}{\prod_{i \in r} u_i!} \theta^{Au}$$

and the conditional probability $P(U=u | AU=Au)$ does not depend on p .

$$\begin{pmatrix} a_{11} & a_{1x} & a_{1r} \\ \vdots & \vdots & \vdots \\ a_{k1} & a_{kx} & a_{kr} \end{pmatrix} \begin{pmatrix} u_1 \\ \vdots \\ u_r \end{pmatrix}$$

pf.

$$\begin{aligned} P(U=u) &= \frac{n!}{\prod_{i \in r} u_i!} \prod_{i \in r} p_i^{u_i} = \frac{n!}{\prod_{i \in r} u_i!} \prod_{i \in r} (\theta_1^{a_{i1}} \cdots \theta_k^{a_{ik}})^{u_i} \\ &= \frac{n!}{\prod_{i \in r} u_i!} \prod_{i \in r} \theta_1^{a_{i1} \cdot u_i} \cdots \theta_k^{a_{ik} \cdot u_i} = \frac{n!}{\prod_{i \in r} u_i!} \prod_{i \in r} \theta_1^{\sum_i a_{i1} \cdot u_i} \cdots \theta_k^{\sum_i a_{ik} \cdot u_i} \\ &= \frac{n!}{\prod_{i \in r} u_i!} \theta^{Au} \end{aligned}$$

Moreover,
$$P(U=u | AU=Au) = \frac{P(U=u)}{P(AU=Au)}$$

$$P(AU=Au) = \sum_{v \in \mathcal{F}(u)} \frac{n!}{\prod_{i \in \mathcal{R}} v_i} \theta^{Av} = n! \theta^{Au} \sum_{v \in \mathcal{F}(u)} \left(\prod_{i \in \mathcal{R}} v_i \right)^{-1}$$

$$\Rightarrow P(U=u | AU=Au) = \frac{1 / \left(\prod_{i \in \mathcal{R}} u_i! \right)}{\sum_{v \in \mathcal{F}(u)} 1 / \left(\prod_{i \in \mathcal{R}} v_i! \right)} \quad \square$$

- Based on Prop 1, we generalize Fisher's exact test by computing the p-value

$$P(X^2(U) \geq X^2(u) | AU=Au)$$

Here
$$X^2(U) = \sum_{i \in \mathcal{C}} \frac{(U_i - \hat{u}_i)^2}{\hat{u}_i}, \quad \hat{u}_i = n \hat{p}_i \quad \text{where}$$

Hence the p-value is:

\hat{p}_i is the MLE.
 \hat{p}_i is the \downarrow Lecture 4.
 same for all tables in the fiber.

$$\frac{\sum_{v \in \mathcal{F}(u)} 1_{X^2(v) \geq X^2(u)} / \left(\prod_{i \in \mathcal{R}} v_i! \right)}{\sum_{v \in \mathcal{F}(u)} 1 / \left(\prod_{i \in \mathcal{R}} v_i! \right)}$$

→ Exact computation of this quantity is prohibitive.
 Thus we sample from elements in the fiber.