

# Algebraic Statistics

**Benjamin C. W. Brown**

B.Brown@ed.ac.uk

~

27th October 2020

## 1. INDEPENDENCE MODELS

## 2. CLASSICAL ALGEBRAIC GEOMETRY

*planetmath.org*

A *statistical model* is usually parameterised by a function, called a *parametrisation*

$\Theta \rightarrow \mathcal{P}$ , given by  $\theta \mapsto P_\theta$ , so that  $\mathcal{P} = \{P_\theta : \theta \in \Theta\}$ ,

where  $\Theta$  is the *parameter space*.  $\Theta$  is usually a subset of  $\mathbb{R}^n$ .

*McCullagh, 2002*

This should be defined using category theory.

## Two-by-Two Contingency Tables

A contingency table contains counts obtained by cross-classifying observed cases according to two or more discrete criteria.

*Example*

TODO: Figure (Florida death sentences)

We ask whether the sentences were made independently of the defendant's race.

## Two-by-Two Contingency Tables

- ▶ Classify using two criteria with  $r$  and  $c$  levels, yields two random variables  $X$  and  $Y$ .
- ▶ Code outcomes as  $[r] := \{1, \dots, r\}$ , and  $[c] := \{1, \dots, c\}$ .

All information about  $X$  and  $Y$  is contained in the *joint probabilities*

$$p_{ij} = P(X = i; Y = j), \quad i \in [r], j \in [c].$$

- ▶ These in turn determine the *marginal probabilities*:

$$p_{i+} := \sum_{j=1}^c p_{ij} = P(X = i), \quad i \in [r],$$
$$p_{+j} := \sum_{i=1}^r p_{ij} = P(Y = j), \quad j \in [c].$$

*Definition*

Two random variables  $X$  and  $Y$  are *independent* if the joint probabilities factor as  $p_{ij} = p_{i+} \cdot p_{+j}$ , for all  $i \in [r]$  and  $j \in [c]$ . Denote independence of  $X$  and  $Y$  by  $X \perp\!\!\!\perp Y$ .

*Proposition*

Two random variables  $X$  and  $Y$  are independent if and only if the  $(r \times c)$ -matrix,  $p = (p_{ij})$ , has rank one.

For a  $(2 \times 2)$ -table, we thus have:

	$P(Y = 1)$	$P(Y = 2)$	
$P(X = 1)$	$p_{11}$	$p_{12}$	$\xrightarrow{X \perp\!\!\!\perp Y} p_{11}p_{22} = p_{12}p_{21}.$
$P(X = 2)$	$p_{21}$	$p_{22}$	

Suppose now we select  $n$  cases, giving rise to  $n$  independent pairs of discrete random variables:

$$\begin{pmatrix} X^{(1)} \\ Y^{(1)} \end{pmatrix}, \begin{pmatrix} X^{(2)} \\ Y^{(2)} \end{pmatrix}, \dots, \begin{pmatrix} X^{(n)} \\ Y^{(n)} \end{pmatrix},$$

all drawn from the same distribution, i.e.:

$$P(X^{(k)} = i; Y^{(k)} = j) = p_{ij}, \quad \text{for all } i \in [r], j \in [c], k \in [n].$$

Joint probability matrix  $p = (p_{ij})$  is an *unknown* element of the  $(rc - 1)$ -dimensional *probability simplex*,

$$\Delta_{rc-1} = \left\{ q \in \mathbb{R}^{r \times c} : q_{ij} \geq 0, \text{ for all } i, j, \text{ and } \sum_{i=1}^r \sum_{j=1}^c q_{ij} = 1 \right\}.$$

*Definitions*

A *statistical model*  $\mathcal{M}$  is a subset of  $\Delta_{rc-1}$ . It represents the set of all candidates for the unknown distribution  $p$ .

The *independence model* for  $X$  and  $Y$  is the set

$$\mathcal{M}_{X \perp\!\!\!\perp Y} := \{p \in \Delta_{rc-1} : \text{rank}(p) = 1\}.$$

$\mathcal{M}_{X \perp\!\!\!\perp Y}$  is the intersection of  $\Delta_{rc-1}$  and the set of all matrices  $p = (p_{ij})$  such that

$$p_{ij}p_{kl} - p_{il}p_{jk} = 0, \text{ for all } 1 \leq i < k \leq r, \text{ and } 1 \leq j < l \leq c.$$

These are examples of *Segre varieties* in algebraic geometry.



### *Projective Space*

Playing field is *n-dimensional projective space*,  $\mathbb{P}^n$ :

$$\mathbb{P}^n := \{(z_0, \dots, z_n) \in \mathbb{C}^n\} / (\mathbf{x} \sim \lambda \cdot \mathbf{y}), \quad \lambda \neq 0,$$

that is, its elements consists of *lines through the origin* in  $\mathbb{C}^n$ .

TODO: FIGURE

*Varieties* are the geometric studied in algebraic geometry, and are the *vanishing sets*<sup>1</sup> for a system of polynomials.

TODO: FIGURES

---

<sup>1</sup>from '*Verschwindungsmenge*'

*Segre varieties* come from  $\sigma : \mathbb{P}^n \times \mathbb{P}^m \rightarrow \mathbb{P}^{(n+1)(m+1)-1}$ , that sends  $([X], [Y])$  to the pairwise products of their components:

$$\sigma : ([X_1, \dots, X_{n+1}], [Y_1, \dots, Y_{m+1}]) \mapsto [\dots, X_i Y_j, \dots],$$

### Example

$$\sigma : \mathbb{P}^1 \times \mathbb{P}^1 \rightarrow \mathbb{P}^3, ([X_1, X_2], [Y_1, Y_2]) \mapsto [X_1 Y_1, X_1 Y_2, X_2 Y_1, X_2 Y_2].$$

$$\text{Set } [X_1 Y_1, X_1 Y_2, X_2 Y_1, X_2 Y_2] = [p_{11}, p_{12}, p_{21}, p_{22}],$$

$$\rightsquigarrow \det \begin{pmatrix} p_{11} & p_{12} \\ p_{21} & p_{22} \end{pmatrix} = 0 \iff \text{rank} \begin{pmatrix} p_{11} & p_{12} \\ p_{21} & p_{22} \end{pmatrix} \leq 1.$$

### *Rulings*

$\sigma(\mathbb{P}^1 \times \mathbb{P}^1) = \{[p_{11}, p_{12}; p_{21}, p_{22}] : \det(p_{ij}) = 0\}$  is an example of a *determinantal variety*.

This example has two families of lines inside of it; the images of  $\sigma([p_{11}, p_{12}] \times \{Q\})$  and  $\sigma(\{Q\} \times [p_{21}, p_{22}])$ , which are called *rulings of the surface*.

Let  $\Delta \subset \mathbb{R}^4$  be the tetrahedron with vertices given by the four basis vectors,  $A_i = e_i$ , and let a general point  $p = (p_{ij})$  inside of  $\Delta$  is represented by

$$p_{ij} = (p_{11}, p_{12}, p_{21}, p_{22}) = \begin{array}{|c|c|} \hline p_{11} & p_{12} \\ \hline p_{21} & p_{22} \\ \hline \end{array}$$