

# Algebraic Statistics

**Benjamin C. W. Brown**

B.Brown@ed.ac.uk

~

28th October 2020

## 1. INDEPENDENCE MODELS

## 2. CLASSICAL ALGEBRAIC GEOMETRY

## 3. MIXTURE MODELS & SECANT VARIETIES

## 4. SUMMARY

*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:

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

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*.

## Manifold of Independence

- ▶ 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  $\Delta$  be represented by

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

- ▶ Fienberg and Gilbert have shown that the two rulings are given by

$$\begin{array}{|c|c|} \hline st & s(1-t) \\ \hline t(1-s) & (1-s)(1-t) \\ \hline \end{array} \quad (0 \leq s, t \leq 1),$$

which is a hyperbolic paraboloid inside of  $\Delta$ .

- ▶ They call this the *manifold of independence*; any point on this surface has independent row and column marginal totals.

TODO: FIGURE.

- ▶ Suppose  $\mathcal{P} \subset \Delta_{r-1}$  is a model for a random variable  $X$  with state space  $[r]$ .
- ▶ Moreover, assume that there is a *hidden* or *latent* random variable  $Y$  with state space  $[s]$ , and for each  $j \in [s]$ , the conditional distribution of  $X$  given  $Y = j$  is  $p^{(j)} \in \mathcal{P}$ .
- ▶ The hidden variable  $Y$  also has some probability distribution  $\pi \in \Delta_{s-1}$ .

So the joint distribution of  $Y$  and  $X$  is given by the formula

$$P(Y = j; X = i) = \pi_j \cdot p_i^{(j)}.$$

## Mixture Models

- ▶ But as  $Y$  is hidden, we can only observe the marginal distribution of  $X$ , that is

$$P(X = i) = \sum_{j=1}^s \pi_j \cdot p_i^{(j)}.$$

- ▶ In other words, the marginal distribution of  $X$  is the convex combination of the  $s$  distributions  $p^{(1)}, \dots, p^{(s)}$ , with weights given by  $\pi$ .

*Definition*

Let  $\mathcal{P} \subset \Delta_{r-1}$  be a statistical model. The  $s$ -th mixture model is

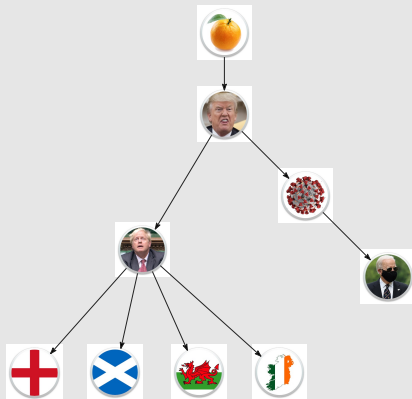
$$\text{Mixt}^s(\mathcal{P}) := \left\{ \sum_{j=1}^s \pi_j \cdot p^{(j)} : \pi \in \Delta_{s-1}, p^{(j)} \in \mathcal{P}, \text{ for all } j \right\}.$$



- ▶ Mixture models provide ways to build complex models out of simpler ones.
- ▶ Basic assumption is that the underlying population to be modelled can be split into  $s$  disjoint sub-populations.
- ▶ Restricted to each sub-population, the observable  $X$  follows a probability distribution from the simple model  $\mathcal{P}$ .
- ▶ After marginalisation though, the structure becomes significantly more complex as it is now a convex combination of these simple distributions.

## Phylogenetic Trees

- ▶ Introduce *phylogenetic trees*; describe the descent of species from a common ancestor:

*Example Cartoon*

- ▶ Sequence of DNA molecules in a genome is represented as a sequence of letters from the four letter alphabet  $\Sigma = \{A, C, G, T\}$ .
- ▶ At a particular site of the genome, any of the four nucleotides  $X \in \Sigma$ , say, might be observed.
- ▶ Based on a particular ancestral nucleotide of  $X$ , might expect evolution to occur in a way that the state of the current nucleotide is independent of one another:

$$A \mapsto X, C \mapsto X, G \mapsto X, \text{ or } T \mapsto X.$$

- ▶ So for each ancestral nucleotide, we have an independence model; its distribution is determined by a point of  $\Delta_3 = \Delta_{4-1}$ .

## Example

- ▶  $X$  is a hidden variable though; could have been any one of A, C, G, or T.
- ▶ For *exactly one choice* of  $X$ , we had the distribution  $\Delta_3$ ; need to consider *all choices* of ancestral nucleotide.
- ▶ Hence we get the mixture model:  $\text{Mixt}^4(\Delta_3) =$

$$\left\{ \sum_{i=1}^4 \text{Prob}(X = i) \cdot \lambda_i : i \in \Sigma, \lambda_i \in [0, 1], \lambda_1 + \dots + \lambda_4 = 1 \right\}.$$

- ▶ **Q?:** What is the analogue for mixture models in algebraic statistics?

- ▶ **A!:** Secant<sup>2</sup> varieties!

### Definitions

- ▶ Consider two varieties  $V, W \subseteq \mathbb{R}^k$ . The *join* of  $V$  and  $W$  is the variety

$$\mathcal{J}(V, W) := \{\lambda v + (1 - \lambda)w : v \in V, w \in W, \lambda \in [0, 1]\}.$$

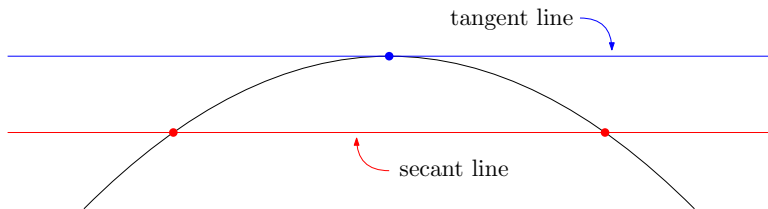
- ▶ If  $V = W$ , then this is the *secant variety* of  $V$ , denoted  $\text{Sec}^2(V) = \mathcal{J}(V, V)$ . The *s-th higher secant variety* is:

$$\text{Sec}^1(V) := V, \quad \text{Sec}^s(V) := \mathcal{J}(\text{Sec}^{s-1}(V), V).$$

---

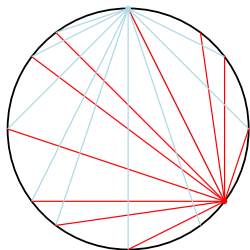
<sup>2</sup>from *secare*, “to cut” in Latin; c.f. *tangō*, “to touch”.

## Secant Varieties

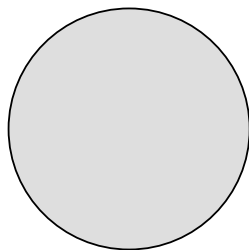


$$S^1 = V(x^2 + y^2 - 1)$$

$$\text{Sec}^2(S^1)$$

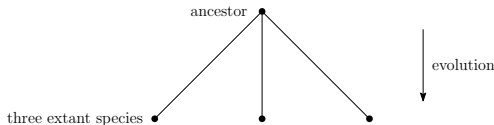


etc.  
→



## More Complicated Phylogenetic Trees

- ▶ Last example only had one extant species;  $X \xrightarrow{?} A, C, G, T$ .
- ▶ What if we had three extant species, coming from one ancestor?



- ▶ Now we have to consider:  $\text{Sec}^4(\mathbb{P}^3 \times \mathbb{P}^3 \times \mathbb{P}^3)$ ; or equivalently  $\text{Mixt}^4(\Delta_3 \times \Delta_3 \times \Delta_3)$ .
- ▶ Finding the minimal set of polynomials defining  $\text{Sec}^4(\mathbb{P}^3 \times \mathbb{P}^3 \times \mathbb{P}^3)$  once gave rise to a very important application of algebraic statistics...

*The Salmon Problem**Statement*

Determine the ideal<sup>3</sup> defining  $\text{Sec}^4(\mathbb{P}^3 \times \mathbb{P}^3 \times \mathbb{P}^3)$ .

*Prize*

Personally caught, and smoked just for you, copper river salmon from Alaska.

*Current Status*

Solved.

- ▶ At an IMA workshop in 2007, Elizabeth Allman offered this prize to whomever solved the above problem.
- ▶ It was solved in 2010 by Shmuel Friedland.

---

<sup>3</sup>read this as “set of defining polynomials”.



Why  $\text{Sec}^4(\mathbb{P}^3 \times \mathbb{P}^3 \times \mathbb{P}^3)$  again?

- ▶ Three independent variables (nucleotides in extant species)  $\rightsquigarrow$  three factors in product;
- ▶ Each independently assumes one value from  $\Sigma = \{\text{A, C, G, T}\} \rightsquigarrow$  distribution is a point in  $\mathbb{P}^3 = \mathbb{P}^{4-1}$ ;
- ▶ The ancestral nucleotide is unknown, but could assume any of the four values in  $\Sigma \rightsquigarrow$  mix four such independence models;
- ▶ The model for the three observed nucleotides is therefore

$$\text{Sec}^4(\mathbb{P}^3 \times \mathbb{P}^3 \times \mathbb{P}^3), \quad \text{c.f.}, \quad \text{Mixt}^4(\Delta_3 \times \Delta_3 \times \Delta_3).$$

## An Opportunity for a Stupid Joke

*Henri Poincaré*

*“[L]a mathématique est l’art de donner le même nom à des choses différentes.”*



Figure: Henri Poincaré, 1887.



Figure: Poincaré, colourised.

The solution to the salmon conjecture is equivalent to:

*Some familiar:*

- ▶ the mixture of four models for three independent variables;
- ▶ the fourth secant variety of the Segre variety  $\mathbb{P}^3 \times \mathbb{P}^3 \times \mathbb{P}^3$ .

*And some less familiar:*

- ▶ the set of  $(4 \times 4 \times 4)$ -tables of tensor rank  $\leq 4$ ;
- ▶ the naive Bayes model with four classes;
- ▶ the conditional independence model  $[X_1 \perp\!\!\!\perp X_2 \perp\!\!\!\perp X_3 | Y]$ ;
- ▶ the general Markov model for the phylogenetic tree,  $K_{1,3}$ ;
- ▶ superposition of four pure states in a quantum system.

## Applications

We finish by mentioning that algebraic statistics has at least a few important applications:

- ▶ It can win you salmon;
- ▶ It can win you 100 Swiss francs<sup>4</sup> (CHF 100  $\sim$  £85);
- ▶ One gets to learn lots of polysyllabic words;
- ▶ It can provide an individual with a topic for an (excellent) colloquium talk;
- ▶ Algebraists & statisticians *could* talk to one other (not that they *would* want to).

---

<sup>4</sup>Not mentioned in this talk.

**Questions?**