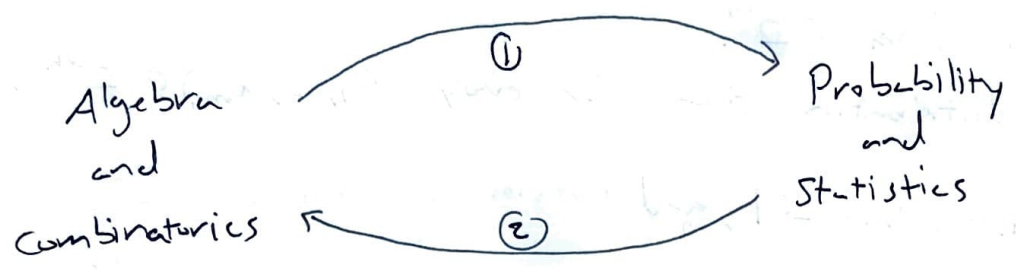


Lesson 1: What is Algebraic Statistics?

"Algebraic statistics" primarily describes the field of study that emerged in the 1990's centered on the connection between algebraic geometry and probability and statistics. Over the decades the "algebraic geometry" bit has become broadly interpreted to include commutative algebra and combinatorics.



① How can algebra and combinatorics be useful in probability and statistics?

② How can ideas from probability and statistics drive new research in algebra and combinatorics?

- Today, we will go through a first example, analyzing how algebra and statistics arise when thinking about statistical problems.
- We will try and state some of the "Big Ideas" that make up the popular research initiatives in the field. A key goal of this course is to give you a solid intro to these popular topics in algebraic statistics.

A First Example

Let X_1, \dots, X_n be random variables on the (finite) state space Ω .

Ex. let $\Omega = \{0, 1\}$. X_i is a function $X_i: \Omega \rightarrow \mathbb{R}$ such that there is a function $P: 2^\Omega \rightarrow [0, 1]$ such that for every measurable subset $B \subseteq \mathbb{R}$

$P(X_i^{-1}(B)) =$ the probability of the outcomes in $X_i^{-1}(B)$.

If $X_i = \text{id}$ then B contains either

• only 0	→	$P(X_i^{-1}(B)) = P(\{0\})$
• only 1	→	$= P(\{1\})$
• both 1 and 0	→	$= P(\{1 \text{ or } 0\}) = P(\Omega) = P(1) + P(0) = 1$
• neither	→	$= P(\text{neither}) = 0$.

- P is the probability distribution associated to \mathcal{X}_i .

• The distributions of each \mathcal{X}_i naturally induce a distribution on all of them considered together.

• The sequence $(\mathcal{X}_1, \dots, \mathcal{X}_m)$ has outcomes $\Omega^m = \{(1,0,1,1), \dots, (1,0), \dots\}$ and an associated probability distribution P such that

$$P(\mathcal{X}_1 = x_1, \dots, \mathcal{X}_m = x_m) = \text{probability of } \mathcal{X}_1 = x_1, \dots, \text{ and } \mathcal{X}_m = x_m$$

for every $(x_1, \dots, x_m) \in \Omega^m$.

• Since P a distribution then for every $(x_1, \dots, x_m) \in \Omega^m$:

(1) $0 \leq P(x_1, \dots, x_m) \leq 1$, and

(2) $\sum_{(x_1, \dots, x_m) \in \Omega^m} P(x_1, \dots, x_m) = 1$.

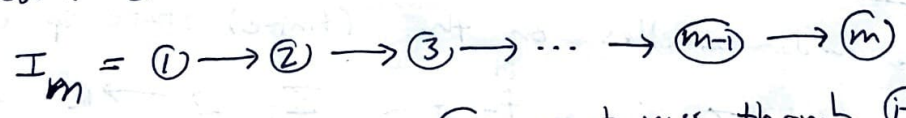
Definition 1 $(\mathcal{X}_1, \dots, \mathcal{X}_m)$ is called a Markov chain if for all $i=3, \dots, m$ and all $(x_1, \dots, x_i) \in \Omega^i$

$$P(\mathcal{X}_i = x_i \mid \mathcal{X}_1 = x_1, \dots, \mathcal{X}_{i-1} = x_{i-1}) = P(\mathcal{X}_i = x_i \mid \mathcal{X}_{i-1} = x_{i-1}).$$

That is, $\mathcal{X}_i \perp\!\!\!\perp \mathcal{X}_{\{1, \dots, i-2\}} \mid \mathcal{X}_{i-1}$.

- In words: the probability of the outcome $\mathcal{X}_i = x_i$ given every preceding outcome depends only on the immediately preceding outcome. In stats, we say, " \mathcal{X}_i is (conditionally) independent of $\mathcal{X}_1, \dots, \mathcal{X}_{i-2}$ given \mathcal{X}_{i-1} ."

- We represent a Markov chain with a directed graph



- info coming from 1, 2, 3, ..., i-2 must pass through i-1 to get to i...

Big Idea (Graphical Models) The combinatorial structure of graphs can be used to capture important information (e.g. conditional independence relations) that hold among random variables.

Definition 2 A statistical model is any subset of distributions that could be associated to a set of random variables.

- Ex. The model $M(I_m) := \left\{ \text{all distributions satisfying Def 1 for } \mathcal{X}_1, \dots, \mathcal{X}_m \right\}$.

- Like many popular statistical models, the underlying description of $M(I_m)$ is inherently algebraic.

- Fix $m=3$ and $\Omega = \{0, 1\}$.

- Any distribution P for $\mathcal{X}_1, \mathcal{X}_2, \mathcal{X}_3$ is a point in \mathbb{R}^8 :

$$P_{011} = P(\mathcal{X}_1=0, \mathcal{X}_2=1, \mathcal{X}_3=1) \in \mathbb{R}.$$

$$P = (P_{000}, P_{100}, P_{010}, P_{001}, P_{110}, P_{101}, P_{011}) \in \mathbb{R}^8.$$

- By (1) and (2)

$$P \in \Delta_{|\Omega|^m-1} = \left\{ p \in \mathbb{R}^8 \mid 0 \leq p_i < 1, \sum_{i \in \Omega^m} p_i = 1 \right\}$$

(The $(|\Omega|^m-1)$ -dim probability simplex)

$$\Rightarrow M(I_3) \subset \Delta_7 \subset \mathbb{R}^8.$$

- Let $\sum_{k \in \Omega} P_{ijk} =: P_{ij+}$. By def of conditional probability:

$$P(\mathcal{X}_3=k \mid \mathcal{X}_1=i, \mathcal{X}_2=j) = \frac{P_{ijk}}{\sum_{k \in \Omega} P_{ij+}}$$

- Def 1 $\Rightarrow \forall i, j, k \in \Omega, 13$,

$$\frac{P_{ijk}}{P_{ij+}} = \frac{P_{+jk}}{P_{+j+}}, \text{ or equivalently for all } i, i', j, k \in \Omega, 13, \frac{P_{ijk}}{P_{ij+}} = \frac{P_{i'jk}}{P_{i'j+}}.$$

- Clearly denominators yields a purely algebraic description of $M(I_3)$:

$$P \in M(I_3) \iff \left. \begin{aligned} (1) & P_{ijk} \geq 0 \quad \forall i, j, k \in \Omega, 13 \\ (2) & \sum_{i, j, k \in \Omega, 13} P_{ijk} = 1 \\ (3) & P_{000} P_{101} - P_{001} P_{100} = 0 \\ (4) & P_{010} P_{111} - P_{011} P_{110} = 0 \end{aligned} \right\} \star$$

- $M(I_3)$ is a semialgebraic set; i.e. a collection of points satisfying a system of polynomial equations and inequalities.

Big Idea Which popular statistical models are representable as "nice" semialgebraic sets? What does this "nice" structure capture about our model?

$M(I_m)$ is an example of a conditional independence model as it is defined by constraints of the form $X_A \perp\!\!\!\perp X_B \mid X_C$.

$\star \Rightarrow M(I_3) = \Delta_7 \cap \underbrace{V(p_{000}p_{101} - p_{001}p_{100}, p_{010}p_{111} - p_{011}p_{110})}_{\text{vanishing set of...}} \hookrightarrow \text{a conditional independence variety.}$

Big Idea What is the algebraic structure of a given conditional independence (ideal or) variety?

While many useful statistical models turn out to be conditional independence models, it is most common in statistics to define a model parametrically. The standard parametrization of $M(I_m)$ follows from the chain rule in probability.

The Chain Rule:

$$P(X_1=x_1, \dots, X_m=x_m) = \prod_{i \in [m]} P(X_i=x_i \mid X_1=x_1, \dots, X_{i-1}=x_{i-1})$$
$$= \prod_{i \in [m]} P(X_i=x_i \mid X_{i-1}=x_{i-1}) \leftarrow (\text{By Def } \textcircled{1})$$

For $M(I_3)$: $P(X_1=i, X_2=j, X_3=k) = \alpha_i \beta_{ji} \gamma_{kj}$ where $\alpha_i = P(X_1=i)$, $\beta_{ji} = P(X_2=j \mid X_1=i)$ and $\gamma_{kj} = P(X_3=k \mid X_2=j)$.

Big Idea Many statistical models have parametric representations in terms of polynomial functions of the parameters. What are their algebraic properties?

The algebraic structure of a statistical model translates into algebraic answers to statistical questions when we start to incorporate data.

Statistical Inference Questions:

- ① How can we fit a statistical model to given data?
- ② Does the model fit the given data well?

Starting (frequentist) Assumption S'pose there is some "true" unknown distribution P from which all our data are independently drawn samples.

- Example: For $M(I_3)$ we could have data

$$ID = \{000, 010, 110, 000, 101, 110, 100, 010, 110, 111, 000, 000, 010\}$$

13 independent and identically distributed samples (iid).

- iid assumption \Rightarrow probability of observing ID really only depends on the vector of counts : $\underline{u} = (u_{ijk} \mid i,j,k \in \{0,1\}) \in \mathbb{R}^8$, where

$u_{ijk} := \#$ of ijk 's in ID .

e.g. $\underline{u} = \left(\frac{4}{000}, \frac{0}{001}, \frac{3}{010}, \frac{1}{100}, \frac{0}{011}, \frac{1}{101}, \frac{3}{110}, \frac{1}{111} \right)$.

- The likelihood function records the probability of observing the data given the model (parameters):

$$L(P \mid \underline{u}) := P(ID) = \frac{n!}{\prod_{ijk} u_{ijk}!} \prod_{ijk \in \{0,1\}^3} P_{ijk}^{u_{ijk}}$$

$\binom{n}{\underline{u}}$ \leftarrow summing over all possible orderings of ID .

where $n = |ID|$. Here \underline{u} is given and P is unknown...

- To find \hat{P} that maximizes the likelihood of seeing ID (i.e. the maximum likelihood estimate (MLE)) we solve

$$\arg \max L(P \mid \underline{u}) \text{ subject to } P \in M(I_3), \text{ or}$$

$$\arg \max \prod_{(i,j,k) \in \{0,1\}^3} (\alpha_i \beta_{ji} \gamma_{kj})^{u_{ijk}} \text{ subject to } \alpha_i, \beta_{ji}, \gamma_{kj} \geq 0, \sum_i \alpha_i = 1, \sum_j \beta_{ji} = 1, \text{ and } \sum_k \gamma_{kj} = 1.$$

- Equivalently, we can solve the same problem for $\log L(P \mid \underline{u})$. To do so, we take the partial derivatives and set equal to 0.

\hat{P} is the solution to a system of polynomial equations. (the score (or critical) equations)

- For $M(I_3)$ this solution is unique!

$$\hat{P} = \left(\hat{p}_{ijk} = \frac{u_{ij+}}{u_{+++}} \cdot \frac{u_{+jk}}{u_{+++}} \mid ijk \in \{0,1\}^3 \right)$$

Big Idea For which statistical models do the score equations admit a unique solution (ML-degree 1)? When does this solution admit a closed-form, rational expression as above?

Testing Goodness-of-fit

• Our second main statistical inference question is whether we can reject the hypothesis:
 H_0 : the data-generating distribution is in $\mathcal{M}(I_3)$.

Question If we assume the data-generating distribution $P \in \mathcal{M}(I_3)$, among all data sets ID' that could have been generated, what proportion of such ID' are more likely to be observed than ID ?

→ If proportion large, reject H_0 .

↳ depends on unknown distribution $P \Rightarrow$ can't compute! $\ddot{=}$

Observation The value of $L(P|u)$ depends only on some lower-order functions of ID (sufficient statistics):

$$L(P|u) = P(u | \alpha_i, \beta_{ji}, \gamma_{kj}) = \binom{n}{u} \prod_{i,j,k \in \{0,1\}^3} (\alpha_i \beta_{ji} \gamma_{kj})^{u_{ijk}}$$
$$= \binom{n}{u} \prod_{i,j \in \{0,1\}^2} (\alpha_i \beta_{ji})^{u_{ij+}} \prod_{k,j \in \{0,1\}^2} (\gamma_{kj})^{u_{+jk}}$$

sufficient statistics = $\{u_{ij+}, u_{+jk}\}$.

Fisher's Question Among all ID' with the same sufficient statistics as ID , what proportion of such ID' are more likely to occur than ID ?

$$P(v | \alpha_i, \beta_{ji}, \gamma_{kj}, u_{ij+} = v_{ij+}, u_{+jk} = v_{+jk}) = \frac{\binom{n}{v} \prod_{ij} (\alpha_i \beta_{ji})^{v_{ij+}} \prod_{jk} \gamma_{kj}^{v_{+jk}}}{\sum_{v: \dots} \binom{n}{v} \prod_{ij} (\alpha_i \beta_{ji})^{v_{ij+}} \prod_{jk} \gamma_{kj}^{v_{+jk}}}$$
$$= \frac{\binom{n}{v}}{\sum_v \binom{n}{v}}$$

• Define the convex polytope

$$\mathcal{P}(u) := \{ v \in \mathbb{R}^8 \mid v_{ij+} = u_{ij+}, v_{+jk} = u_{+jk}, v_{ijk} \geq 0 \}.$$

• The proportion of interest (Fisher's p-value) is:

$$p = \frac{\sum_{v \in \mathcal{P}(u) \cap \mathbb{Z}^8} \binom{n}{v} \chi\left(\binom{n}{v} > \binom{n}{u}\right)}{\sum_{v \in \mathcal{P}(u) \cap \mathbb{Z}^8} \binom{n}{v}}$$

Big Idea Testing goodness-of-fit (via Fisher's exact test) amounts to enumerating lattice points in the convex polytope $\mathcal{P}(u)$.

Next time... Katherin will give us a review of some of the key algebraic concepts needed in this course. In the third lecture, I will do the same for probability and statistics. After that, the lectures will be focused on exploring the "big ideas" from today.

These lectures you will give: each person will be assigned a lecture and given associated reading. They will give their lecture to one of us the Monday before their Friday lecture.