

Statistical Analysis.

Karl Pearson (1904): "On the Theory of Contingency and Its Relation to Association and Normal Correlation." [Draper's Company Research Memoirs Biometric Series I.]

Contingency Table } Cross Tabulation
Cross Tab.

- ◊ A matrix $\in [0, 1]^{m \times n}$ displaying bi-variate (in general, multivariate) frequency distribution of the variables.

	Left-Handed	Right-Handed	
Male	2	3	5
Female	1	4	5
Marginal Total	3	7	(10) Grand Total

Q: Are female less likely to be left-handed than right-handed?

- ◊ Statistical Tests: { Pearson's chi-squared test
G-test
Fisher's Exact test
Barnard's test.

Normalized Table $\in [0, 1]^{2 \times 2}$

$$\bar{S}^2 = \frac{1}{m(m-1)} \sum_{i,j} S_{ij}^2 \left\{ \begin{array}{l} A = \begin{bmatrix} \frac{1}{5} & \frac{3}{10} \\ \frac{1}{10} & \frac{2}{5} \end{bmatrix} \\ A^T = \begin{bmatrix} \frac{1}{5} & \frac{1}{10} \\ \frac{3}{10} & \frac{2}{5} \end{bmatrix} \\ S = AA^T = \begin{bmatrix} \frac{13}{100} & \frac{14}{100} \\ \frac{14}{100} & \frac{17}{100} \end{bmatrix} \\ \bar{S}^2 = \frac{28/100}{2} = \frac{14}{100} \end{array} \right.$$

◊ Note, if there was a perfect correlation

$$A = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$$

and $AA^T = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$ and $\bar{s}^2 = \frac{2}{2} = 1$

◊ Note, also, if there was a perfect anticorrelation

$$A = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} \text{ then } AA^T = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} \quad \bar{s}^2 = \frac{0}{2} = 0.$$

DARWIN'S FINCHES:

Finch	Islands	Islands																
		A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P	Q
Large Ground Finch		0	0	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1
Medium Ground F.		1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1
Small Ground F.		1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1
Cactus Ground F.		1	1	1	0	1	1	1	1	1	1	1	1	1	1	1	1	1
Sharp Beaked G.F.		0	0	1	1	1	0	0	1	0	1	0	1	1	0	1	1	1
Large Cactus G.F.		0	0	0	0	0	0	0	0	0	0	1	0	1	0	0	0	0
Large Tree F.		0	0	1	1	1	1	1	1	1	1	0	0	1	0	1	1	0
Medium Tree F.		0	0	0	0	0	0	0	0	0	0	0	1	0	0	0	1	0
Small Tree F.		0	0	1	1	1	1	1	1	1	1	1	0	1	0	1	1	0
Vegetarian F.		0	0	1	1	1	1	1	1	1	1	1	0	1	0	0	0	0
Woodpecker F.		0	0	1	1	1	0	1	1	0	1	0	0	0	0	0	0	0
Mangrove F.		0	0	1	1	0	0	0	0	0	0	0	0	0	0	0	0	0
Warbler F.		1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1

13

Finch	Marginal Bounds (Row Sums)																	Island																
	A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P	Q																	
Large ground finch	0	0	1	1	1	1	1	1	1	1	0	1	1	1	1	1	1																	
Medium ground finch	1	1	1	1	1	1	1	1	1	1	0	1	0	1	1	0	1																	
Small ground finch	1	1	1	1	1	1	1	1	1	1	1	1	0	1	1	0	1																	
Sharp-beaked ground finch	0	0	1	1	1	0	1	1	0	1	1	1	1	0	1	1	0																	
Cactus ground finch	1	1	1	0	1	1	1	1	1	1	0	1	0	1	1	0	1																	
Large cactus ground finch	0	0	0	0	0	0	0	0	0	0	1	1	0	1	0	0	0																	
Large tree finch	0	0	1	1	1	1	1	1	1	0	0	1	0	1	0	0	0																	
Medium tree finch	0	0	0	0	0	0	0	0	0	0	0	1	0	0	0	0	0																	
Small tree finch	0	0	1	1	1	1	1	1	1	1	1	1	0	0	0	0	0																	
Vegetarian finch	0	0	1	1	1	1	1	1	1	1	0	1	0	0	0	0	0																	
Woodpecker finch	0	0	1	1	1	0	1	1	0	1	0	1	0	0	0	0	0																	
Mangrove finch	0	0	1	1	0	0	1	1	0	0	0	0	0	0	0	0	0																	
Warbler finch	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1																	
	15	13	14	10	12	2	10	1	10	11	6	2	17																					

NOTE: Island name code: A = Seymour, B = Baltra, C = Isabella, D = Fernandina, E = Santiago, F = Rábida, G = Pinzón, H = Santa Cruz, I = Santa Fe, J = San Cristóbal, K = Española, L = Floreana, M = Genoves N = Marchena, O = Pinta, P = Darwin, Q = Wolf.

Marginal Row ds. → 4 sums 11 10 10 8 9 10 8 9 10 8 3 10 4 7 9 3 3

(43)

Darwin's Data.

$$A = 13 \times 17 \text{ matrix over } \{0, 1\}$$

$$= \{0, 1\}^{13 \times 17}$$

Grand Total = 123

Row Sums

$$A(1, \cdot), A(2, \cdot), \dots, A(13, \cdot)$$

Column Sums

$$A(\cdot, 1), A(\cdot, 2), \dots, A(\cdot, 17)$$

Problem:

Data: Row Sums / Column Sums

$$A(i, \cdot) \quad i = 1, \dots, m$$

$$A(\cdot, j) \quad j = 1, \dots, n$$

$$\sum_i A(i, \cdot) = \sum_j A(\cdot, j) = T$$

Find:

(i) Compute k distinct 0-1 A matrices satisfying the marginals given by r 's Sums and Column Sums.

$$k = 1; 100; 10000.$$

There are

67, 149, 106, 137, 567, 626 distinct tables.

(ii) Compute one A matrix that maximizes \sum^2 .

$$A = (x_{ij})$$

(44)

Feasibility

$$\sum_{j=1}^n x_{ij} = A(i \cdot) = i^{\text{th}} \text{ row sum } \quad i=1, \dots, m$$

$$\sum_{i=1}^m x_{ij} = A(\cdot j) = j^{\text{th}} \text{ column sum } \quad j=1, \dots, n$$

$$x_{ij} \in \{0, 1\}$$

~~~~~

Optimization

Note: Feasibility is in  $P$ .

Based on Gale-Ryser Condition:

Using a trick called "Majorization"

$$\begin{aligned} \text{Note } s_{ij} &= \sum_k x_{ik} x_{kj} = \sum_k x_{ik} x_{jk} \\ &= \sum_k \mathbb{1}_{x_{ik} = x_{jk} = 1} \\ &= \sum_k \mathbb{1}_{x_{ik} = x_{jk} = 1} \end{aligned}$$

We wish to maximize  $\frac{1}{m(m-1)} \sum_{i \neq j} s_{ij}^2$

simplify to  $\sum_{i \neq j} s_{ij}^2 \rightarrow \sum_{i \neq j} \omega_{ij} s_{ij}$

$$\omega_{ij} = A_m A(i \cdot) A(j \cdot)$$

maximize  $\sum_{i \neq j} A(i \cdot) A(j \cdot) \sum_k \mathbb{1}_{x_{ik} = x_{jk} = 1}$

$$= \sum_{i \neq j} \sum_k A(i \cdot) A(j \cdot) \mathbb{1}_{x_{ik} = x_{jk} = 1}$$

subject to

$$\sum_{j=1}^n x_{ij} = A(i \cdot); \quad \sum_{i=1}^m x_{ij} = A(\cdot j);$$

$$x_{ij} \in \{0, 1\}.$$

## A General Problem:

(45)

$c_1, c_2, \dots, c_m$  are species.

$x_1, x_2, \dots, x_n$  are islands.

We say  $c_i$  and  $c_j$  are symbionts in islands  $X' \subseteq \{x_1, \dots, x_n\}$

iff  $x_{ik} \Leftrightarrow x_{jk} \quad \forall x_k \in X'$

Both  $c_i$  and  $c_j$  cohabit islands  $X'$

We say  $c_i$  and  $c_j$  are antibionts in islands  $X' \subseteq \{x_1, \dots, x_n\}$

iff  $\bar{x}_{ik} \vee \bar{x}_{jk} \quad \forall x_k \in X'$

$$\boxed{x_{ik} \Rightarrow \bar{x}_{jk} \wedge x_{jk} \Rightarrow \bar{x}_{ik}}$$

Neither  $c_i$  nor  $c_j$  cohabits any island in  $X'$

Row Sum  $A(i \cdot) = \#$  islands occupied by a species  $i$

Column Sum  $A(\cdot j) = \#$  species occupying island  $j$

Data:  $m, n;$

Symbiont relation  $\{(c_i, c_j, x'_{ij})\}$

Antibiont relation  $\{(c_i, c_j, x'_{ij})\}$

Row Sum

Column Sum

CONSTRAINTS

Desiderata: Find a matrix  $\in \{0, 1\}^{m \times n}$  satisfying all the constraints.

## NP Completeness

ONE-IN-THREE POSITIVE 3-SAT

For each clause  $C_i$ : Create two species,  $c_i, c_i'$

$(c_i, c_i', \{x_k \mid x_k \in C_i\}) \leftarrow$  Add to Antibiont reln.

Row Sum  $(c_i) = 1$ ; Row Sum  $(c_i') = 2$ .

If  $x_k \in C_i \cap C_j$

$(c_i, c_j, \{x_k\}) \leftarrow$  Add to Symbiont reln

Column Sum  $(x_k) = |\{C_i \mid x_k \in C_i\}|$

Example:

(16)

$$(x \vee y \vee z) \wedge (x \vee u \vee v) \wedge (w \vee u \vee v)$$

|        |            |            |       |       |       |            |   |
|--------|------------|------------|-------|-------|-------|------------|---|
|        | $x$        | $u$        | $v$   | $w$   | $y$   | $z$        |   |
| $c_1$  | 1          |            |       |       | 0     | 0          | } |
| $c_1'$ | 0          |            |       |       | 1     | 1          |   |
| $c_2$  | 1          | 0          |       |       |       | 0          |   |
| $c_2'$ | 0          | 1          |       |       |       | 1          |   |
| $c_3$  |            | 0          | 0     | 1     |       |            |   |
| $c_3'$ |            | 1          | 1     | 0     |       |            |   |
|        | 2          | 2          | 1     | 1     | 1     | 2          |   |
|        | ↓          | ↓          | ↓     | ↓     | ↓     | ↓          |   |
|        | $c_1, c_2$ | $c_2, c_3$ | $c_3$ | $c_3$ | $c_1$ | $c_1, c_2$ |   |

$$\text{Anti} \left\{ \begin{array}{l} (c_1, c_1', \{x, y, z\}) \\ (c_2, c_2', \{x, u, z\}) \\ (c_3, c_3', \{w, u, v\}) \end{array} \right.$$

$$\text{Sym} \left\{ \begin{array}{l} (c_1, c_2, \{x, z\}) \\ (c_2, c_3, \{u\}) \end{array} \right.$$



# LINEAR PROGRAMMING.

(47)

## Standard Form:

◊ A Linear Function to be maximized:

$$f(x_1, x_2) = c_1 x_1 + c_2 x_2 = [c_1 \ c_2] \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = C^T x$$

◊ PROBLEM CONSTRAINTS:

$$\begin{aligned} a_{11} x_1 + a_{12} x_2 &\leq b_1 \\ a_{21} x_1 + a_{22} x_2 &\leq b_2 \\ a_{31} x_1 + a_{32} x_2 &\leq b_3 \end{aligned} \Leftrightarrow \begin{bmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \\ a_{31} & a_{32} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \leq \begin{bmatrix} b_1 \\ b_2 \\ b_3 \end{bmatrix}$$

$$\Leftrightarrow Ax \leq b$$

◊ Non-negative Variables

$$\begin{aligned} x_1 &\geq 0 \\ x_2 &\geq 0 \end{aligned} \Leftrightarrow \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \geq \begin{bmatrix} 0 \\ 0 \end{bmatrix} \Leftrightarrow x \geq 0$$

## Matrix Form

$$\max \{ C^T x \mid Ax \leq b \wedge x \geq 0 \} \leftarrow \text{Primal}$$

## Symmetric Dual Form

$$\min \{ b^T y \mid A^T y \geq c \wedge y \geq 0 \} \leftarrow \text{Dual}$$

Strong Duality Thm: If primal has an optimal soln  $x^*$ , then the dual has an optimal soln  $y^*$ , and

$$C^T x^* = b^T y^*.$$

Algorithms: Classical: Dantzig's Simplex Algorithm (Exptime)

Fast in practice.

Criss-Cross Algorithm.

Modern: Ellipsoid Algorithm (Khachiyan) } Interior  
Projective Algorithm (Karmakar) } Point  
Path-following Algorithm }  $\in P$ .

## Integer Linear Programming (ILP):

(48)

If all of the unknown variables are required to be integers then the problem is called an ILP.

0-1 Integer Programming or Binary Integer Programming (BIP) is a special case of ILP, where variables are required to be 0 or 1.

## Mixed Integer Programming (MIP)

Only some of the unknown variables are required to be integers.

## 3-SAT.

For every variable  $x_i$  put

$$0 \leq x_i \leq 1$$

$$x_i \in \{0, 1\}$$

For every clause  $C_j = (x_{i_1} \vee x_{i_2} \vee \bar{x}_{i_3})$  put

$$x_{i_1} + x_{i_2} + (1 - x_{i_3}) > 0.$$

↓  
unnegated  
literals

↑  
Negated  
literals.

~~~~~


SEMIDEFINITE PROGRAMMING: (SDP)

(49)

$X = n \times n$ matrix

$X =$ Positive Semidefinite (psd) if

$$\forall v \in \mathbb{R}^n \quad v^T X v \geq 0$$

$X =$ Positive Definite (pd) if

$$\forall v \in \mathbb{R}^n \quad v^T X v > 0.$$

$X \succeq 0 \leftarrow X$ is symmetric & positive Definite.

$$A \bullet B = \sum_{i=1}^n \sum_{j=1}^n A_{ij} B_{ij} = \text{Tr}(A^T B)$$

SDP: minimize $C \bullet X$

$$\text{s.t.} \quad A_i \bullet X = b_i \quad i=1, \dots, m$$

$$X \succeq 0.$$

$$A_i = \begin{pmatrix} a_{i1} & & & \\ & a_{i2} & & \\ & & \dots & \\ 0 & & & a_{in} \end{pmatrix}$$

$$i=1, \dots, m \quad C = \begin{pmatrix} c_1 & & & \\ & c_2 & & \\ & & \dots & \\ 0 & & & c_n \end{pmatrix}$$

LP \equiv SDP:

minimize $C \bullet X$

$$\text{s.t.} \quad A_i X = b_i \quad i=1, \dots, m$$

$$X_{ij} = 0 \quad i=1, \dots, m; \quad j=i+1, \dots, n$$

$$X \succeq 0$$

$$X = \begin{pmatrix} x_1 & & & \\ & x_2 & & \\ & & \dots & \\ 0 & & & x_n \end{pmatrix}$$

Quadratically Constrained Quadratic Prog. (50)

QCQP.

$$\underset{x}{\text{minimize}} \quad x^T Q x + q_0^T x + c_0$$

$$\text{s.t.} \quad x^T Q_i x + q_i^T x + c_i \leq 0 \quad i=1, \dots, m.$$

Factor each Q_i :

$$Q_i = M_i^T M_i$$

Note

$$\begin{pmatrix} I & M_i x \\ x^T M_i^T & -c_i - q_i^T x \end{pmatrix} \succeq 0 \quad \Leftrightarrow \quad x^T Q_i x + q_i^T x + c_i \leq 0$$

$$\underset{x, \theta}{\text{minimize}} \quad \theta$$

s.t.

$$\begin{pmatrix} I & M_0 x \\ x^T M_0^T & -c_0 - q_0^T x + \theta \end{pmatrix} \succeq 0$$

$$\begin{pmatrix} I & M_i x \\ x^T M_i^T & -c_i - q_i^T x \end{pmatrix} \succeq 0 \quad i=1, \dots, m.$$

SDP can be solved in Polytime using interior point method.