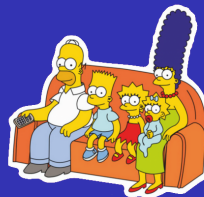


Bayesian Networks,
Mendelian Genetics,
and ...

THE
SIMPSONS



G. Nuel
May 5, 2020
Chaire MMB
SU, Paris, France



Blood Feud

In “*Blood Feud*” (S02E22):

- The rich Mr Burns badly needs O- blood for a transfusion
- Desperate Mr Smithers asks nuclear power plant’s employees to give blood for the great man



Homer Simpson sees a golden opportunity:

- Mr Burns being immensely rich, he will probably reward the one who saves him
- Homer asks Marge about his own blood group, but he’s A+ (Homer: “D’oh !”)

What are the chances to find a O- blood donor in the Simpson family ?

The Blood Group/Rhesus system

Lisa asks Pr Frink about blood group genetics:

- **ABO gene** with three alleles: two codominant A, B, and one recessive O
 $\Rightarrow p_O = 0.60, p_A = 0.30, p_B = 0.10$
- **RHD gene** with two alleles: dominant D (positive Rh) and recessive d (negative Rh)
 $\Rightarrow q_D = 0.60, q_d = 0.40$
- This leads to a total of 8 blood phenotypes: A+, B+, AB+, O+, A-, B-, AB-, O-



With ABO/RHD independence we get:



	ABO	OO	OA	OB	AA	AB	BB
RHD		0.36	0.36	0.12	0.09	0.06	0.01
DD	0.36	O+	A+	B+	A+	AB+	B+
Dd	0.48	O+	A+	B+	A+	AB+	B+
dd	0.16	O-	A-	B-	A-	AB-	B-

The Blood Group/Rhesus system

Lisa asks Pr Frink about blood group genetics:

- **ABO gene** with three alleles: two codominant A, B, and one recessive O
 $\Rightarrow p_O = 0.60, p_A = 0.30, p_B = 0.10$
- **RHD gene** with two alleles: dominant D (positive Rh) and recessive d (negative Rh)
 $\Rightarrow q_D = 0.60, q_d = 0.40$
- This leads to a total of 8 blood phenotypes: A+, B+, AB+, O+, A-, B-, AB-, O-



With ABO/RHD independence we get:



	group	O	A	B	AB
Rh		0.36	0.45	0.13	0.06
+	0.84	0.3024	0.378	0.1092	0.0504
-	0.16	0.0576	0.072	0.0208	0.0096

Nuclear Family



The problem to solve:

- 1: Homer, 2: Marge, 3: Bart, 4: Lisa, 5: Maggie
- X_i (resp. Y_i) genotype (resp. phenotype) or ind. i
- we need to compute

$$\pi = \mathbb{P}(\exists i, Y_i = O \mid Y_1 = A+)$$

Model 1: independent genotypes:

$$\mathbb{P}(X, Y) = \prod_{i=1}^5 \mathbb{P}(X_i) \mathbb{P}(Y_i \mid X_i)$$

- $\mathbb{P}(Y_1 = O \mid Y_1 = A+) = 0$, for $i \neq 1$, $\mathbb{P}(Y_i = O) = 0.0576$
- $\pi = 1 - (1 - 0.0576)^4 \simeq 0.2112$, **easy, right !?**

Lisa: “*But genotypes are not independent !*”

Homer: “*D’oh !*”

Nuclear Family

The problem to solve:



- 1: Homer, 2: Marge, 3: Bart, 4: Lisa, 5: Maggie
- X_i (resp. Y_i) genotype (resp. phenotype) or ind. i
- we need to compute

$$\pi = \mathbb{P}(\exists i, Y_i = O \mid Y_1 = A+)$$

Model 2: Mendelian transmission of alleles:

$$\mathbb{P}(X, Y) = \mathbb{P}(X_1)\mathbb{P}(X_2) \prod_{i=3}^5 \mathbb{P}(X_i \mid X_1, X_2) \times \prod_{i=1}^5 \mathbb{P}(Y_i \mid X_i)$$

- $\mathbb{P}(X_1 = OADd \mid Y_1 = A+) = \frac{0.36}{0.36+0.09} \times \frac{0.48}{0.36+0.48} \simeq 0.4571$
- $\mathbb{P}(X_2 = O^A_B Dd) = 0.48 \times 0.48 = 0.2304$
- $\mathbb{P}(X_2 = O^A_B dd) = 0.48 \times 0.16 = 0.0768$
- $\mathbb{P}(X_2 = OO Dd) = 0.36 \times 0.48 = 0.1728$
- $\mathbb{P}(X_2 = OO dd) = 0.36 \times 0.16 = 0.0576$

Nuclear Family

The problem to solve:



- 1: Homer, 2: Marge, 3: Bart, 4: Lisa, 5: Maggie
- X_i (resp. Y_i) genotype (resp. phenotype) or ind. i
- we need to compute

$$\pi = \mathbb{P}(\exists i, Y_i = O- | Y_1 = A+)$$

Model 2: Mendelian transmission of alleles:

$$\mathbb{P}(X, Y) = \mathbb{P}(X_1)\mathbb{P}(X_2) \prod_{i=3}^5 \mathbb{P}(X_i | X_1, X_2) \times \prod_{i=1}^5 \mathbb{P}(Y_i | X_i)$$

$ev = \{Y_1 = A+\}$ and N is the number of $O-$ in the nuclear family

- $\mathbb{P}(X_1 = OADD, AADd, AADD | ev) = 0.5429 \Rightarrow N \sim \mathcal{B}(1, 0.0567)$
- $\mathbb{P}(X_1 = OADd, X_2 \text{ not carrier} | ev) \simeq 0.2114 \Rightarrow N = 0$
- $\mathbb{P}(X_1 = OADd, X_2 = O_B^A Dd | ev) \simeq 0.1053 \Rightarrow N \sim \mathcal{B}(3, \frac{1}{16})$
- $\mathbb{P}(X_1 = OADd, X_2 = O_B^A dd, OODd | ev) \simeq 0.1141 \Rightarrow N \sim \mathcal{B}(3, \frac{1}{8})$
- $\mathbb{P}(X_1 = OADd, X_2 = OOdd | ev) \simeq 0.0263 \Rightarrow N \sim 1 + \mathcal{B}(3, \frac{1}{4})$

Nuclear Family

The problem to solve:



- 1: Homer, 2: Marge, 3: Bart, 4: Lisa, 5: Maggie
- X_i (resp. Y_i) genotype (resp. phenotype) or ind. i
- we need to compute

$$\pi = \mathbb{P}(\exists i, Y_i = O \mid Y_1 = A+)$$

Model 2: Mendelian transmission of alleles:

$$\mathbb{P}(X, Y) = \mathbb{P}(X_1)\mathbb{P}(X_2) \prod_{i=3}^5 \mathbb{P}(X_i \mid X_1, X_2) \times \prod_{i=1}^5 \mathbb{P}(Y_i \mid X_i)$$

$ev = \{Y_1 = A+\}$ and N is the number of O- in the nuclear family

$$N \mid ev \sim 0.5429 \times \mathcal{B}(1, 0.0567) + 0.2114 \times \mathbf{1}_{N=0} + 0.1053 \times \mathcal{B}(3, 1/16) \\ + 0.1141 \times \mathcal{B}(3, 1/8) + 0.0263 \times (1 + \mathcal{B}(3, 1/4))$$

$$\sum_{n \geq 0} \mathbb{P}(N = n \mid ev) z^n \simeq 0.8867 + 0.0920z + 0.0169z^2 + 0.0039z^3 + 0.0004z^4$$

$$\pi = 0.0920 + 0.0169 + 0.0039 + 0.0004 = 0.1132$$

What does it mean ?



Homer: “So $\pi = 0.1132$ or whatever, is that good ?”

- $\mathbb{P}(\text{at least one O-} | \text{Homer is A+}) = 11.32\%$
- $\mathbb{P}(\text{not any O-} | \text{Homer is A+}) = 88.68\%$

Homer: “Huh ...”

9 to 10 chances of not being able to help Mr Burns

Homer: “D’oh !”

Frink: “But I found a blood test in the criminal record of Bart ...”

Homer: “... and ?”

Frink: “**Bart is O- !!**”

Homer: “Woohoo !”



Weak-D Bart

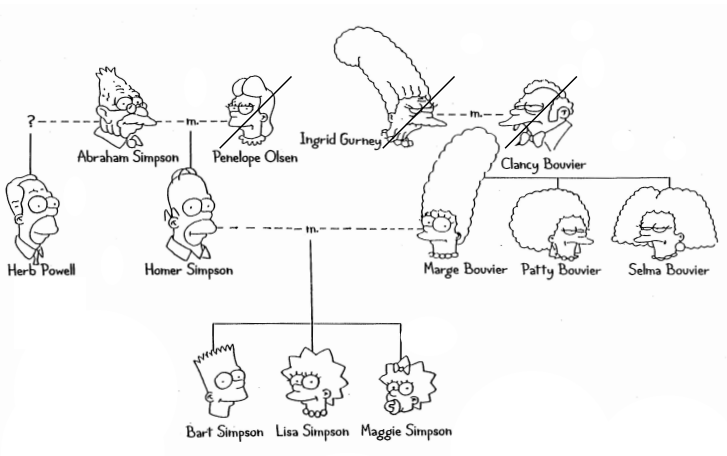


- Frink: “*Alas, Bart has the weak-D phenotype. His blood might kill Mr Burns !*”
- RHD gene, 3 alleles: D, d, and w (weak-D, pseudo neg Rh)
 $\Rightarrow q_D = 0.60, q_d = 0.39, q_w = 0.01$
- Only **OOdd** is compatible with Mr Burns !

RHD	ABO	OO	OA	OB	AA	AB	BB
		0.36	0.36	0.12	0.09	0.06	0.01
DD	0.3600	O+	A+	B+	A+	AB+	B+
Dd	0.4680	O+	A+	B+	A+	AB+	B+
Dw	0.0120	O+	A+	B+	A+	AB+	B+
dd	0.1521	O-	A-	B-	A-	AB-	B-
dw	0.0078	Ow	Aw	Bw	Aw	ABw	Bw
ww	0.0001	Ow	Aw	Bw	Aw	ABw	Bw

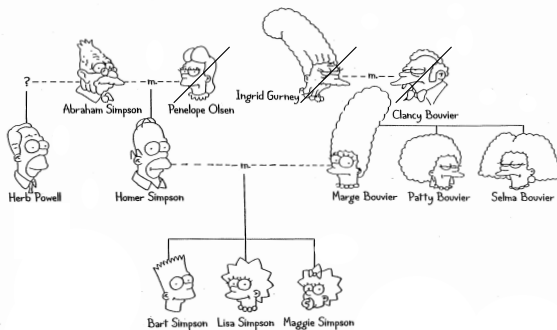
Lisa: “*We need to consider the **extended Simpson family !***”

Simpsons' Pedigree and Bayesian Network



- 1: Herb's mother, 2: Abraham, 3: Penelope, 4: Ingrid, 5: Clancy,
6: Herb, 7: Homer, 8: Marge, 9: Patty, 10: Selma,
11: Bart, 12: Lisa, 13: Maggie

Simpsons' Pedigree and Bayesian Network



- 1:** Herb's mother, **2:** Abraham, **3:** Penelope, **4:** Ingrid, **5:** Clancy,
6: Herb, **7:** Homer, **8:** Marge, **9:** Patty, **10:** Selma,
11: Bart, **12:** Lisa, **13:** Maggie

$$\mathbb{P}(X) = \mathbb{P}(X_1)\mathbb{P}(X_2)\mathbb{P}(X_3)\mathbb{P}(X_4)\mathbb{P}(X_5)$$

$$\mathbb{P}(X_6 | X_{1,2})\mathbb{P}(X_7 | X_{2,3})\mathbb{P}(X_8 | X_{4,5})\mathbb{P}(X_9 | X_{4,5})\mathbb{P}(X_{10} | X_{4,5})$$

$$\mathbb{P}(X_{11} | X_{7,8})\mathbb{P}(X_{12} | X_{7,8})\mathbb{P}(X_{13} | X_{7,8})$$

Simpsons' Pedigree and Bayesian Network

$$\mathbb{P}(X) = \mathbb{P}(X_1)\mathbb{P}(X_2)\mathbb{P}(X_3)\mathbb{P}(X_4)\mathbb{P}(X_5)$$

$$\mathbb{P}(X_6 | X_{1,2})\mathbb{P}(X_7 | X_{2,3})\mathbb{P}(X_8 | X_{4,6})\mathbb{P}(X_9 | X_{4,6})\mathbb{P}(X_{10} | X_{4,6})$$

$$\mathbb{P}(X_{11} | X_{7,8})\mathbb{P}(X_{12} | X_{7,8})\mathbb{P}(X_{13} | X_{7,8})$$

$$X_i \in \mathcal{G} = \{O, A, B\}^2 \times \{D, d, w\}^2 \quad |\mathcal{G}| = 3^2 \times 3^2 = 81$$

$$ev = \{\text{Homer A+ and Bart Ow}\} \quad \mathbb{P}(X|ev) = \frac{\mathbb{P}(X, ev)}{\sum_{X'} \mathbb{P}(X', ev)}$$



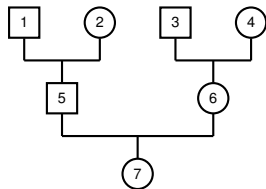
- $X = (X_1, X_2, \dots, X_{13})$ is the family genotype
- in order to compute $\mathbb{P}(ev) = \sum_{X'} \mathbb{P}(X', ev)$
- we *just* have to sum over 81^{13} configurations

$$81^{13} = 6\,461\,081\,889\,226\,672\,446\,898\,176$$

\Rightarrow simply impossible !

Local computations in a simple pedigree

Idea: we consider a smaller (but similar) family, ev (evidence) still represents the available information.



- for *founders* (1, 2, 3, 4) i :

$$\varphi_i(\mathbf{X}_i) = \mathbb{P}(\mathbf{X}_i \cap ev)$$

- for *offsprings* (5, 6, 7) k with parents i, j :

$$\varphi_j(\mathbf{X}_i, \mathbf{X}_j, \mathbf{X}_k) = \mathbb{P}(\mathbf{X}_k \cap ev \mid \mathbf{X}_i, \mathbf{X}_j)$$

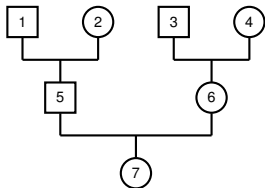


$$\mathbb{P}(ev) = \sum_{\mathbf{X}_1} \sum_{\mathbf{X}_2} \sum_{\mathbf{X}_3} \sum_{\mathbf{X}_4} \sum_{\mathbf{X}_5} \sum_{\mathbf{X}_6} \sum_{\mathbf{X}_7} \varphi_1(\mathbf{X}_1) \varphi_2(\mathbf{X}_2) \varphi_3(\mathbf{X}_3) \varphi_4(\mathbf{X}_4) \\ \varphi_5(\mathbf{X}_1, \mathbf{X}_2, \mathbf{X}_5) \varphi_6(\mathbf{X}_3, \mathbf{X}_4, \mathbf{X}_6) \varphi_7(\mathbf{X}_5, \mathbf{X}_6, \mathbf{X}_7)$$

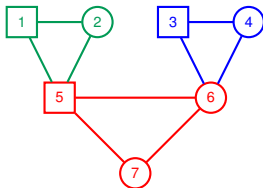
$\Rightarrow 81^7 = 22\,876\,792\,454\,961$ still too large !!

Local computations in a simple pedigree

Pedigree

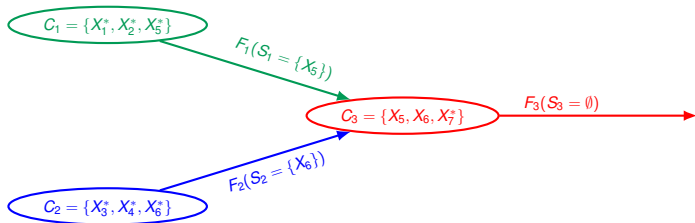


Clique decomposition



$$\mathbb{P}(\text{ev}) = \sum_{X_5} \sum_{X_6} \sum_{X_7} \left\{ \overbrace{\left(\sum_{X_1} \sum_{X_2} \varphi_1(X_1) \varphi_2(X_2) \varphi_5(X_1, X_2, X_5) \right)}^{F_1(X_5)} \right. \\
 \left. \overbrace{\left(\sum_{X_3} \sum_{X_4} \varphi_3(X_3) \varphi_4(X_4) \varphi_6(X_3, X_4, X_6) \right)}^{F_2(X_6)} \varphi_7(X_5, X_6, X_7) \right\}$$

Local computations in a simple pedigree



$$F_j(S_j) = \sum_{C_j \setminus S_j} \left(\prod_{i \in \text{from}_j} F_i(S_i) \right) \times \prod_{X_U \in C_j^*} \varphi_U(X_{\text{pa}_U}, X_U) \quad F_3(\emptyset) = \mathbb{P}(\text{ev})$$

Complexity:

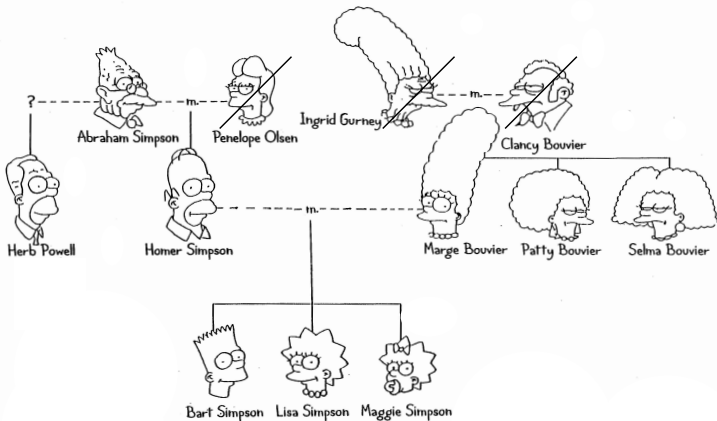
- from $81^7 = 22\,876\,792\,454\,961$
- to $3 \times 81^3 = 1\,594\,323$

Lisa: “*Much better !*”

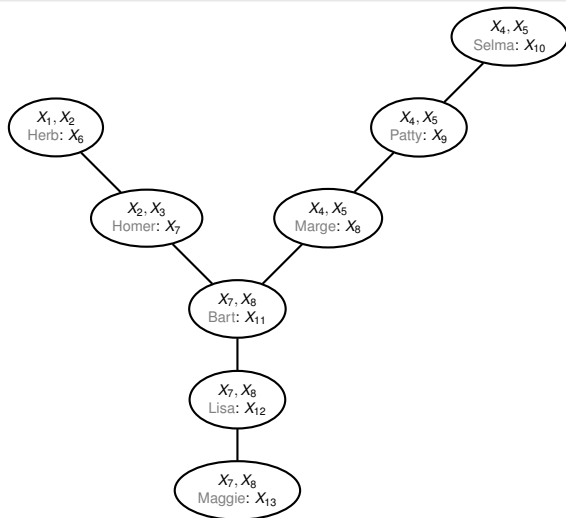
Homer: “*Woohoo !*”



Clique decomposition for the Simpsons

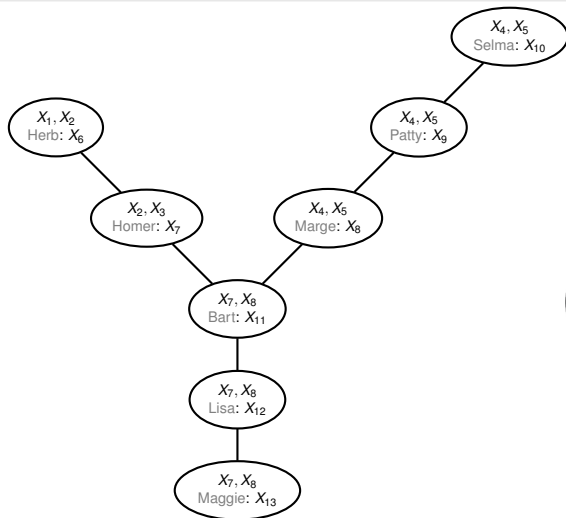


Clique decomposition for the Simpsons



- from $81^{13} = 6\,461\,081\,889\,226\,672\,446\,898\,176$
- to $8 \times 81^3 = 4\,251\,528$

Clique decomposition for the Simpsons



- from $81^{13} = 6\,461\,081\,889\,226\,672\,446\,898\,176$
- to $8 \times 81^3 = 4\,251\,528$

Posterior Distribution

After the [sum-product algorithm](#) we can combine forward/backward quantities with the potentials φ_i in order to derive:

- the marginal distribution $\mathbb{P}(C_j|ev)$ of each clique;
- the marginal distribution $\mathbb{P}(S_j|ev)$ of each separator;
- the full posterior $\mathbb{P}(X|ev)$ as a heterogenous Markov chain.

Introducing for all ind. i the *donor compatibility* random variable:

$$C_i = \mathbf{1}_{X_i=O\text{Odd}} = \begin{cases} 1 & \text{if } i \text{ compatible donor i.e. } X_i = O\text{Odd} \\ 0 & \text{if } i \text{ not compatible donor i.e. } X_i \neq O\text{Odd} \end{cases}$$

$$\mathbb{P}(C_i = 1|ev)$$

Herb's mum	Abraham	Penelope	Ingrid	Clancy
5.5%	2.4%	2.4%	16.2%	16.2%
Herb	Homer	Marge	Patty	Selma
3.8%	0.0%	13.5%	10.9%	10.9%
Bart	Lisa	Maggie		
0.0%	1.5%	1.5%		

Probability Generating Function

Idea: introduce a dummy variable z in the sum-product s.t.

$$\sum_{n \geq 0} \mathbb{P}(N = n, \text{ev}) z^n = \sum_X \prod_i \varphi_i(X_{\text{pa}_i}, X_i) z^{1_{X_i = \text{Odd}}}$$

complexity is just d times the previous one (d max degree).

$$\begin{aligned} \sum_{n \geq 0} \mathbb{P}(N = n | \text{ev}) z^n &= 0.528 + 0.282z + 0.0895z^2 + 0.0404z^3 \\ &+ 0.0329z^4 + 0.0250z^5 + 0.00189z^6 + 0.000148z^7 + 3.76 \times 10^{-7}z^8 \end{aligned}$$

$$\mathbb{P}(C_i = 1 | \text{ev})$$

Herb's mum	Abraham	Penelope	Ingrid	Clancy
5.5%	2.4%	2.4%	16.2%	16.2%
Herb	Homer	Marge	Patty	Selma
3.8%	0.0%	13.5%	10.9%	10.9%
Bart	Lisa	Maggie		
0.0%	1.5%	1.5%		

Probability Generating Function

Idea: introduce a dummy variable z in the sum-product s.t.

$$\sum_{n \geq 0} \mathbb{P}(N = n, \text{ev}) z^n = \sum_X \prod_i \varphi_i(X_{\text{pa}_i}, X_i) z^{1_{X_i=0\text{Odd}}}$$

complexity is just d times the previous one (d max degree).

$$\begin{aligned} \sum_{n \geq 0} \mathbb{P}(N = n | \text{ev}) z^n &= 0.528 + 0.282z + 0.0895z^2 + 0.0404z^3 \\ &+ 0.0329z^4 + 0.0250z^5 + 0.00189z^6 + 0.000148z^7 + 3.76 \times 10^{-7}z^8 \end{aligned}$$

$$\mathbb{P}(C_i = 1 | N = 0, \text{ev})$$

Herb's mum	Abraham	Penelope	Ingrid	Clancy
0.0%	0.0%	0.0%	0.0%	0.0%
Herb	Homer	Marge	Patty	Selma
0.0%	0.0%	0.0%	0.0%	0.0%
Bart	Lisa	Maggie		
0.0%	0.0%	0.0%		

Probability Generating Function

Idea: introduce a dummy variable z in the sum-product s.t.

$$\sum_{n \geq 0} \mathbb{P}(N = n, \text{ev}) z^n = \sum_X \prod_i \varphi_i(X_{\text{pa}_i}, X_i) z^{1_{X_i = \text{Odd}}}$$

complexity is just d times the previous one (d max degree).

$$\begin{aligned} \sum_{n \geq 0} \mathbb{P}(N = n | \text{ev}) z^n &= 0.528 + 0.282z + 0.0895z^2 + 0.0404z^3 \\ &+ 0.0329z^4 + 0.0250z^5 + 0.00189z^6 + 0.000148z^7 + 3.76 \times 10^{-7}z^8 \end{aligned}$$

$$\mathbb{P}(C_i = 1 | N = 1, \text{ev})$$

Herb's mum	Abraham	Penelope	Ingrid	Clancy
9.4%	5.4%	6.5%	28.3%	28.3%
Herb	Homer	Marge	Patty	Selma
5.8%	0.0%	11.8%	0.0%	0.0%
Bart	Lisa	Maggie		
0.0%	2.3%	2.3%		

Probability Generating Function

Idea: introduce a dummy variable z in the sum-product s.t.

$$\sum_{n \geq 0} \mathbb{P}(N = n, \text{ev}) z^n = \sum_X \prod_i \varphi_i(X_{\text{pa}_i}, X_i) z^{1_{X_i = \text{Odd}}}$$

complexity is just d times the previous one (d max degree).

$$\begin{aligned} \sum_{n \geq 0} \mathbb{P}(N = n | \text{ev}) z^n &= 0.528 + 0.282z + 0.0895z^2 + 0.0404z^3 \\ &+ 0.0329z^4 + 0.0250z^5 + 0.00189z^6 + 0.000148z^7 + 3.76 \times 10^{-7}z^8 \end{aligned}$$

$$\mathbb{P}(C_i = 1 | N = 2, \text{ev})$$

Herb's mum	Abraham	Penelope	Ingrid	Clancy
18.6%	6.5%	4.5%	31.3%	31.3%
Herb	Homer	Marge	Patty	Selma
15.5%	0.0%	40.0%	19.8%	19.8%
Bart	Lisa	Maggie		
0.0%	6.3%	6.3%		

Probability Generating Function

Idea: introduce a dummy variable z in the sum-product s.t.

$$\sum_{n \geq 0} \mathbb{P}(N = n, \text{ev}) z^n = \sum_X \prod_i \varphi_i(X_{\text{pa}_i}, X_i) z^{1_{X_i = \text{Odd}}}$$

complexity is just d times the previous one (d max degree).

$$\begin{aligned} \sum_{n \geq 0} \mathbb{P}(N = n | \text{ev}) z^n &= 0.528 + 0.282z + 0.0895z^2 + 0.0404z^3 \\ &+ 0.0329z^4 + 0.0250z^5 + 0.00189z^6 + 0.000148z^7 + 3.76 \times 10^{-7}z^8 \end{aligned}$$

$$\mathbb{P}(C_i = 1 | N = 3, \text{ev})$$

Herb's mum	Abraham	Penelope	Ingrid	Clancy
14.6%	6.0%	3.0%	33.2%	33.2%
Herb	Homer	Marge	Patty	Selma
11.8%	0.0%	27.9%	80.2%	80.2%
Bart	Lisa	Maggie		
0.0%	4.9%	4.9%		

Probability Generating Function

Idea: introduce a dummy variable z in the sum-product s.t.

$$\sum_{n \geq 0} \mathbb{P}(N = n, \text{ev}) z^n = \sum_X \prod_i \varphi_i(X_{\text{pa}_i}, X_i) z^{1_{X_i = \text{Odd}}}$$

complexity is just d times the previous one (d max degree).

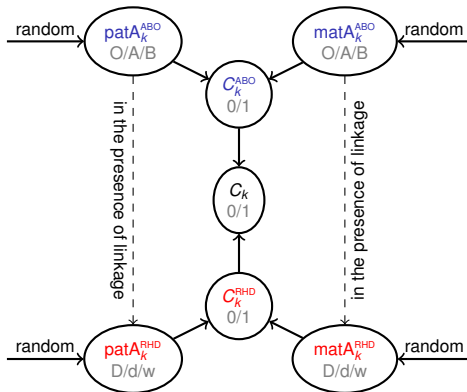
$$\begin{aligned} \sum_{n \geq 0} \mathbb{P}(N = n | \text{ev}) z^n &= 0.528 + 0.282z + 0.0895z^2 + 0.0404z^3 \\ &+ 0.0329z^4 + 0.0250z^5 + 0.00189z^6 + 0.000148z^7 + 3.76 \times 10^{-7}z^8 \end{aligned}$$

$$\mathbb{P}(C_i = 1 | N = 4, \text{ev})$$

Herb's mum	Abraham	Penelope	Ingrid	Clancy
6.9%	2.1%	1.5%	48.2%	48.2%
Herb	Homer	Marge	Patty	Selma
4.8%	0.0%	86.5%	97.8%	97.8%
Bart	Lisa	Maggie		
0.0%	3.2%	3.2%		

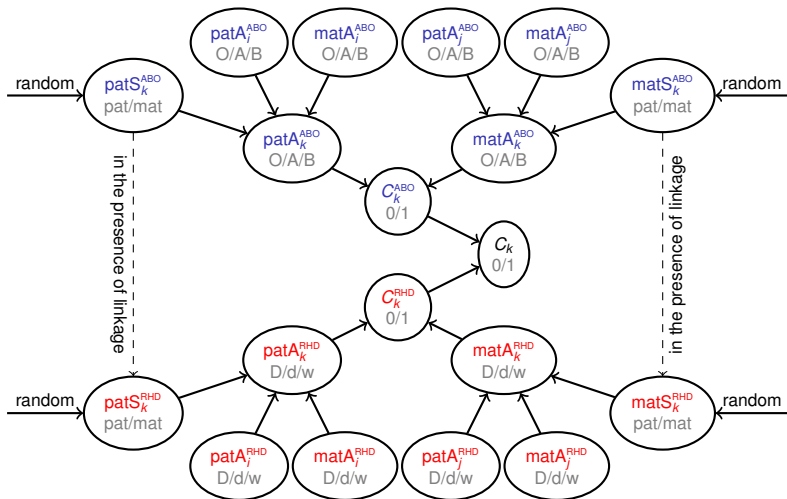
Extended Pedigree: Small Variables

For a *founder* i , instead of $X_i \in \mathcal{G}$ we have:



Extended Pedigree: Small Variables

For a *offspring* k (with father i and mother j),
instead of $X_k \in \mathcal{G} | X_i, X_j$ we have:



Extended Pedigree: Small Variables

Recall on complexity:

- naive $81^{13} = 6\,461\,081\,889\,226\,672\,446\,898\,176$
- genotypes $8 \times 81^3 = 4\,251\,528$

Small variables with the three heuristics:

- **min-neighbors**: the smallest clique
⇒ 61154 61649 89051
- **min-fill**: the clique with minimum fill-in
⇒ 85205 92333 92360
- **weighted min-fill**: the clique with minimum weighted fill-in
⇒ 57530 43841 43112



Many human diseases:

- Cancers:
 - Breast and Ovarian: *Institut Curie*
 - MSI Cancer and Lynch Syndrome: *Saint-Antoine*
 - Gliomas: *La Pitié-Salpêtrière*
- Rare Genetic Diseases:
 - Neuropathy Amyloid Hereditary: *Henri Mondor*
 - Pulmonary Arterial HT: *Marie Lannelongue*
 - Huntington Disease: *Hôpital Saint-Anne*
- Common Disease with Genetic Factors:
 - Alzheimer Disease: *CHU Rouen*
 - Diabetes, autism, cardio-vascular, obesity, . . .



And other applications: linkage analysis, genetic epidemiology, agronomy, recreative genetics, . . .

Aknowledgments

Funding:



and Special Thanks:

- Homer, Marge, Bart, Lisa, and Maggie Simpson
- Matt Groening

