

# Computing risks based on family history in genetic diseases

**Alexandra Lefebvre**, Olivier Bouaziz, Grégory Nuel

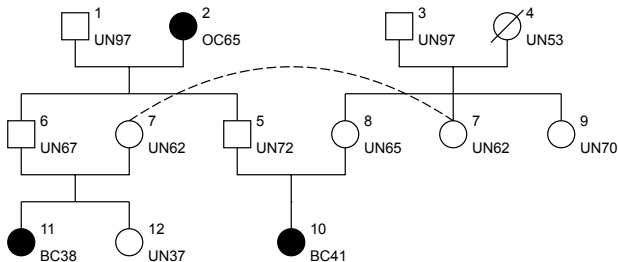
LPSM, UPMC, Sorbonne Université, Paris, France

SMPGD 2018  
Université de Montpellier

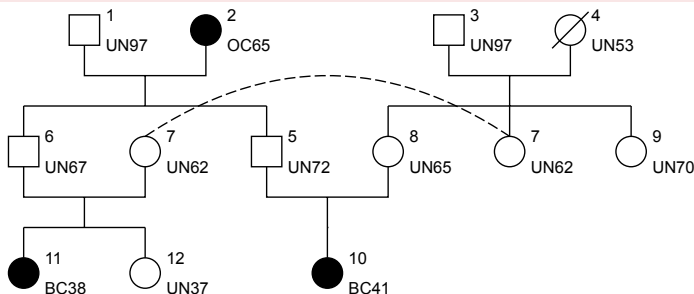


# Context of the breast cancer

- 1st cancer in women. 54,000 women in France each year
- Complex disease due to an accumulation of mutations (BRCA 1/2, PALB2, RAD51, etc.)
- **Inherited mutation** in 10 to 15% of the cases



# Family History (FH) : Pedigree + Sequencing (seq) + Survival (Y)



## Variant data (var)

- Known variant or VUS
- Database
- Functional tests
- Molecular dynamics

## Cancer pathology (patho)

- MSI status (Lynch)
- BRAF mutation (Lynch)
- ER, PR, HER2 status (Breast)
- invasive/in situ (Breast)

compute  $\Rightarrow$  posterior **variant status** **carrier risk** **tumoral risk**

$V$  : the set of variant status

$\text{var}$  : the set of variant data

$X$  : the set of genotypes

$\text{seq}$  : the set of sequencing data

$Y$  : the set of phenotypes (survival data)

$\text{patho}$  : the set of pathology reports

$$\left. \begin{array}{l} \text{seq} : \text{the set of sequencing data} \\ Y : \text{the set of phenotypes (survival data)} \end{array} \right\} \rightarrow \text{FH} = \{\text{seq}, Y\}$$

$$\begin{aligned} \mathbb{P}(\text{var}, V, X, \text{seq}, Y, \text{patho}) = & \prod_j \mathbb{P}(\text{var}_j | V_j) \mathbb{P}(V_j) \\ & \prod_i \mathbb{P}(X_i | X_{\text{pat}_i}, X_{\text{mat}_i}) \mathbb{P}(\text{seq}_i | X_i) \mathbb{P}(Y_i | X_i) \\ & \prod_{i, Y_i \neq \text{UN}} \mathbb{P}(\text{patho}_i | X_i) \end{aligned}$$

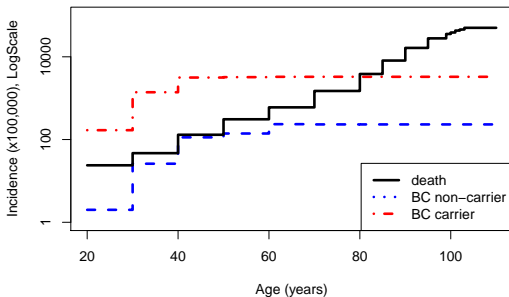
Our first objective : modelling efficiently the quantity

$$\prod_i \mathbb{P}(X_i | X_{\text{pat}_i}, X_{\text{mat}_i}) \mathbb{P}(Y_i | X_i)$$

Illustrated with a mendelian model and parameters taken from the breast cancer literature [Claus et al., 1991, Easton et al., 1993].

## The Claus-Easton model [Claus et al., 1991]

- Autosomal, biallelic, dominant, estimated allele frequency  $f=0.33\%$
- The hazard functions per genotype ( $\lambda_0$  and  $\lambda_1$ ) estimated from the densities in [Easton et al., 1993]:



$$\frac{\lambda_1(t)}{\lambda_0(t)} \in [20 - 80]$$

Objective : Implement the Claus-Easton model in a Bayesian network (sum/product algorithm) combined with survival data.

# Implementation : genotypes

$$\mathbb{P}(X, Y) = \prod_i \underbrace{\mathbb{P}(X_i | X_{\text{pat}_i}, X_{\text{mat}_i})}_{\text{genotype}} \underbrace{\mathbb{P}(Y_i | X_i)}_{\text{phenotypes}}$$

**Mode of inheritance** : 1 autosomal biallelic gene,  $f = 0.33\%$

Founders (**Hardy-Weinberg**) : 
$$\begin{cases} \mathbb{P}(X_i = 00) = (1 - f)^2 \\ \mathbb{P}(X_i = 10) = \mathbb{P}(X_i = 01) = f(1 - f) \\ \mathbb{P}(X_i = 11) = f^2 \end{cases}$$

Offsprings (**Mendel**): 
$$\begin{cases} \mathbb{P}(X_i = 00) = (1 - \Theta(X_{\text{pat}})) \times (1 - \Theta(X_{\text{mat}})) \\ \mathbb{P}(X_i = 10) = \Theta(X_{\text{pat}}) \times (1 - \Theta(X_{\text{mat}})) \\ \mathbb{P}(X_i = 01) = (1 - \Theta(X_{\text{pat}})) \times \Theta(X_{\text{mat}}) \\ \mathbb{P}(X_i = 11) = \Theta(X_{\text{pat}}) \times \Theta(X_{\text{mat}}) \end{cases}$$

with  $\Theta(00) = 0$ ,  $\Theta(10) = \Theta(01) = 0.5$ ,  $\Theta(11) = 1$

# Implementation : phenotypes

with  $T_i$ , the age at disease onset for the individual  $i$ .

$$Y_i \text{ (Survival data)} = \begin{cases} \{T_i > \tau_i\} & \text{if } i \text{ is censored (UN) at age } \tau_i \\ \{T_i = \tau_i\} & \text{if } i \text{ is affected (BC or OC) at age } \tau_i \end{cases}$$

## Dominant model of disease :

Hazard functions of  $T$  :  $\begin{cases} \lambda_0(t) & \text{for } X = 00 \text{ (non-carrier)} \\ \lambda_1(t) & \text{for } X \neq 00 \text{ (carrier)} \end{cases}$

Survival functions of  $T$  :  $\begin{cases} S_0(t) = \exp\left(-\int_0^t \lambda_0(t) dt\right) & \text{for non-carrier} \\ S_1(t) = \exp\left(-\int_0^t \lambda_1(t) dt\right) & \text{for carrier} \end{cases}$

## conditional probabilities

- For a **censored** individual at age  $\tau_i$

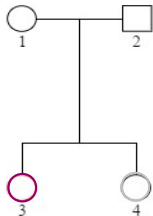
$$\mathbb{P}(Y_i|X_i) = \mathbb{P}(T_i > \tau_i|X_i) = \begin{cases} S_0(\tau_i) & \text{for non-carriers} \\ S_1(\tau_i) & \text{for carriers} \end{cases}$$

- For an **affected** individual at age  $\tau_i$

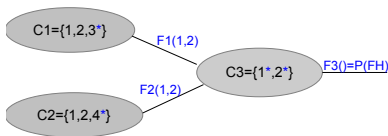
$$\mathbb{P}(Y_i|X_i) = \mathbb{P}(T_i = \tau_i|X_i) = \begin{cases} S_0(\tau_i)\lambda_0(\tau_i) & \text{for non-carriers} \\ S_1(\tau_i)\lambda_1(\tau_i) & \text{for carriers} \end{cases}$$

# Impl. in a Bayesian network <sup>1</sup> (sum-product algorithm)

$$\mathbb{P}(X, FH) = \prod_i \sum_{Y_i \in FH} \underbrace{\mathbb{P}(X_i | X_{pa_i}) \mathbb{P}(Y_i | X_i)}_{K_i(X_i, X_{pa_i})} \rightarrow \mathbb{P}(FH) = \sum_X \prod_i K_i(X_i, X_{pa_i})$$



With  $X \in \{00, 10, 01, 11\}^n \rightarrow 4^n$  configurations

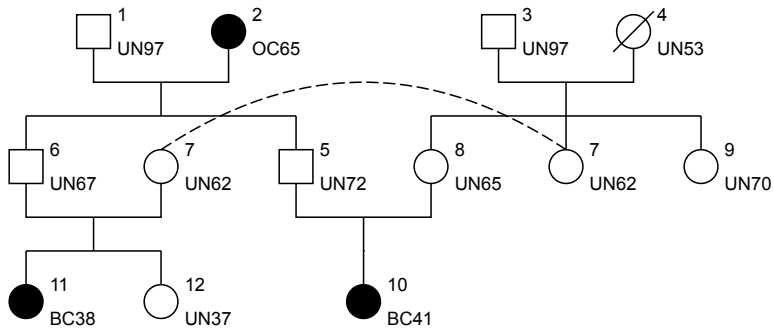


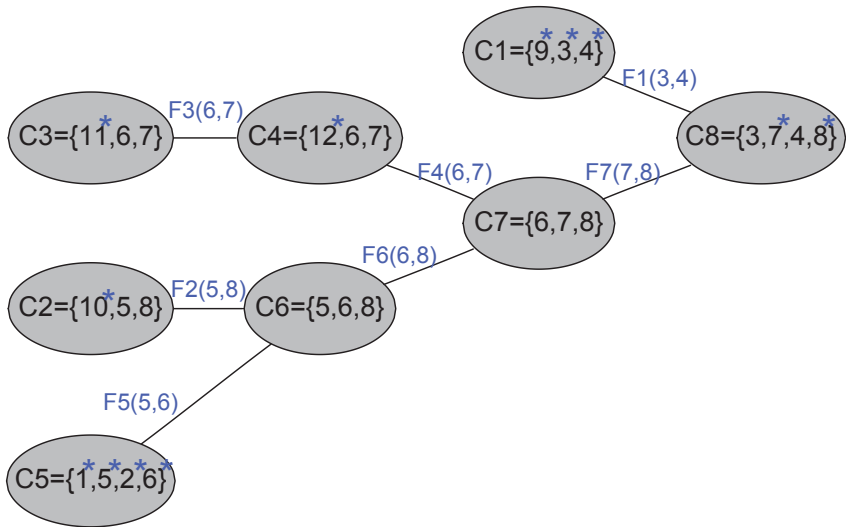
$$\mathbb{P}(FH) = \sum_{X_1} K_1(X_1) \sum_{X_2} K_2(X_2) \underbrace{\sum_{X_3} K_3(X_3, X_1, X_2)}_{F_1(X_1, X_2)} \underbrace{\sum_{X_4} K_4(X_4, X_1, X_2)}_{F_2(X_1, X_2)}$$

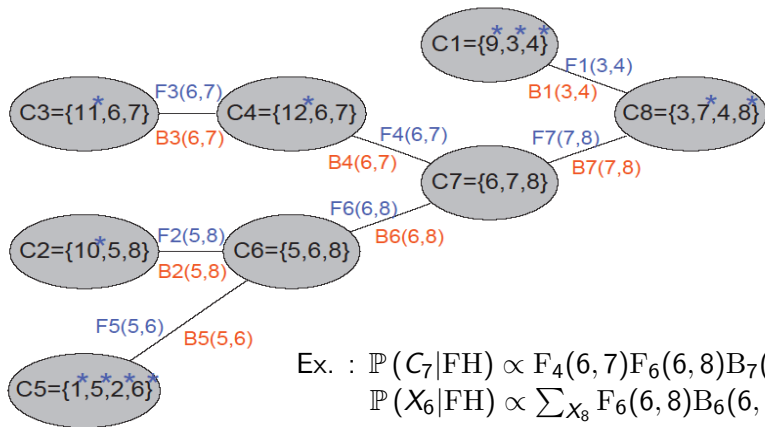
$$\mathbb{P}(FH) = F_3(\emptyset) = \sum_{X_1, X_2} K_1(X_1) K_2(X_2) F_1(X_1, X_2) F_2(X_1, X_2)$$

<sup>1</sup>[Koller and Friedman, 2009]



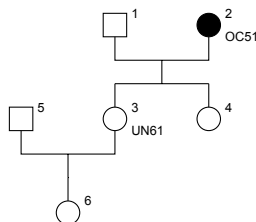
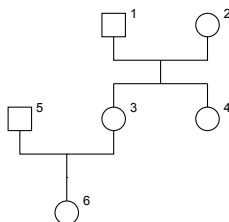
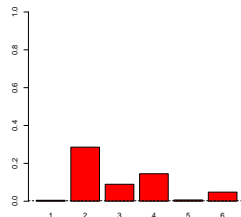
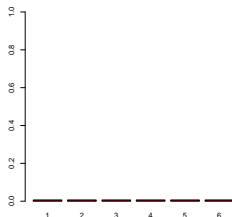




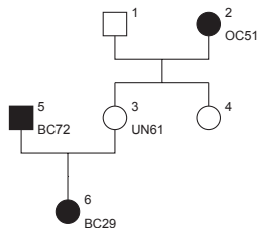
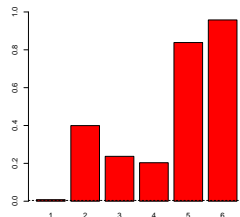
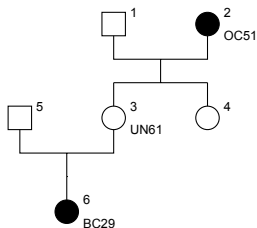
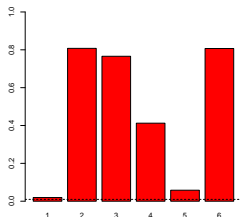


F & B computed once for any later marginal or joint distribution needed  
 Complexity  $\mathcal{O}(4^n) \rightarrow \mathcal{O}(n \times 4^k)$ , k:tree-width

# Results : Carrier risk, posterior marginal carrier distribution



# Results : Carrier risk, posterior marginal carrier distribution



# Implementation : Disease risk prediction

$$\pi(\tau) = \mathbb{P}(\text{carrier}|\text{FH})$$

- Breast cancer risk with no competing risk of death.

$$\mathbb{P}(T \leq t|\text{FH}) = 1 - S(t|\text{FH}) \text{ with}$$

$$S(t|\text{FH}) = \sum X_i \mathbb{P}(T > t, X_i|\text{FH}) = \pi(\tau) \frac{S_1(t)}{S_1(\tau)} + (1 - \pi(\tau)) \frac{S_0(t)}{S_0(\tau)}$$

- With competing risk of death.

$$T^* = \min(T_{\text{disease}}, T_{\text{death}})$$

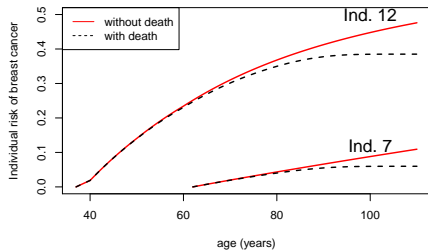
$$\lambda_{\text{both}}(t|\text{FH}) = \lambda_{\text{disease}}(t|\text{FH}) + \lambda_{\text{death}}(t)$$

$$\begin{aligned} \mathbb{P}(T \leq t|\text{FH}) &= \int_{\tau}^t S_{\text{both}}(u) \lambda_{\text{disease}}(u) du \\ &= \int_{\tau}^t \exp\left(-\int_{\tau}^u \lambda_{\text{both}}(v) dv\right) \lambda_{\text{disease}}(u) du \end{aligned}$$

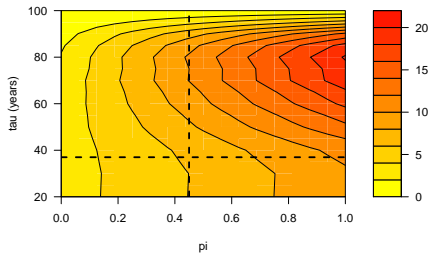
$$\lambda_{\text{disease}}(t|\text{FH}) = \pi(t|\text{FH}) \lambda_1(t) + (1 - \pi(t|\text{FH})) \lambda_0(t)$$

$$\pi(t|\text{FH}) = \frac{\pi(\tau) S_1(t)}{S(t|\text{FH}) S_1(\tau)}$$

# Results : Disease risk



Individual disease risk  
with vs without  
competing risk of death



Difference in % of disease risk at 100  
with vs without  
competing risk of death

- Adaptation to **Microsatellite Instability (MSI cancer / Lynch syndrome)** (A. Duval, Saint-Antoine Hospital)
- Taking into account the sequencing data, the pathology reports, variant data especially to help classify the **Variants of Uncertain Significance (VUS)**.
- **Parameters estimations**
- Complex distributions (number of carriers in the family) with **probability generating functions** (polynomials) → familial risk
- **Multi-state** and frailty survival models





Claus, E. B., Risch, N., and Thompson, W. D. (1991).

Genetic analysis of breast cancer in the cancer and steroid hormone study.

*American journal of human genetics*, 48(2):232.



Easton, D., Bishop, D., Ford, D., and Crockford, G. (1993).

Genetic linkage analysis in familial breast and ovarian cancer: results from 214 families. the breast cancer linkage consortium.

*American journal of human genetics*, 52(4):678.



Koller, D. and Friedman, N. (2009).

*Probabilistic graphical models: principles and techniques*.

MIT press.

# Thank you for your attention

