

Projets MAP311

# Modélisation bayésienne des genres

Xiang CHEN & Xin CHEN

26/06/2017

# Modélisation bayésienne des genres

## 1. Préliminaires.

1) On a  $f_{p,\theta}(x) = \frac{1}{B(a,b)} x^{a-1} (1-x)^{b-1} 1_{]0,1[}(x)$

Par la définition:

$$\mathbb{E}_\theta[P] = \int_R x f_{p,\theta}(x) dx = \int_0^1 \frac{1}{B(a,b)} x^a (1-x)^{b-1} dx = \frac{B(a+1,b)}{B(a,b)}$$

On rappelle  $\Gamma(n) = \int_0^{+\infty} t^{n-1} e^{-t} dt$

Par Intégration par partie

On a  $\Gamma(n) = \left[ \frac{1}{n} t^n e^{-t} \right]_0^{+\infty} + \int_0^{+\infty} \frac{1}{n} t^n e^{-t} dt = 0 + \int_0^{+\infty} \frac{1}{n} t^n e^{-t} dt = \frac{\Gamma(n+1)}{n}$

Comme  $B(a,b) = \frac{\Gamma(a)\Gamma(b)}{\Gamma(a+b)}$

$$\begin{aligned} \text{Alors } \mathbb{E}_\theta[P] &= \frac{B(a+1,b)}{B(a,b)} = \frac{\Gamma(a+1)\Gamma(b)}{\Gamma(a+b+1)} \times \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} \\ &= \frac{\Gamma(a+1)}{\Gamma(a)} \times \frac{\Gamma(a+b)}{\Gamma(a+b+1)} \\ &= \frac{a}{a+b} \end{aligned}$$

On a  $p_0 = \mathbb{E}_\theta[P] = \frac{a}{a+b}$

$$\begin{aligned} \text{De même, } \mathbb{E}_\theta[P^v(1-P)^w] &= \int_R x^v (1-x)^w f_{p,\theta}(x) dx \\ &= \int_0^1 \frac{1}{B(a,b)} x^{v+a-1} (1-x)^{w+b-1} dx \\ &= \frac{B(a+v,b+w)}{B(a,b)} \end{aligned}$$

2) On veut la loi conditionnelle de  $X_n$  sachant  $P$

Si  $P = p$

Alors  $\mathbb{P}_\theta(X_n = k | P = p) = \binom{n}{k} p^k (1-p)^{n-k}$

Donc  $\mathbb{P}_\theta(X_n = k | P) = \binom{n}{k} P^k (1-P)^{n-k}$

C'est-à-dire  $X_n | P \sim B(n, P)$

3) On prend  $g_1: \begin{matrix} \mathbb{R}^+ & \rightarrow & \mathbb{R}^+ \\ x & \mapsto & x \end{matrix}$  mesurable positive

Alors on a  $\mathbb{E}_\theta[\mathbb{E}_\theta[g_1(X_n) | P]] = \mathbb{E}_\theta[g_1(X_n)]$

Donc  $\mathbb{E}_\theta[X_n] = \mathbb{E}_\theta[\mathbb{E}_\theta[X_n | P]]$

$$= \mathbb{E}_\theta[nP] \text{ car } X_n|P \sim B(n, P)$$

$$= n\mathbb{E}_\theta[P]$$

$$= \frac{na}{a+b}$$

On prend  $g_2: \begin{matrix} \mathbb{R}^+ & \rightarrow & \mathbb{R}^+ \\ x & \mapsto & x \end{matrix}$  mesurable positive

$$\text{De même } \mathbb{E}_\theta[\mathbb{E}_\theta[g_2(X_n)|P]] = \mathbb{E}_\theta[g_2(X_n)]$$

$$\text{Donc } \mathbb{E}_\theta[X_n^2] = \mathbb{E}_\theta[\mathbb{E}_\theta[X_n^2|P]]$$

$$\text{Or } \mathbb{E}_\theta[X_n^2|P] = \text{Var}_\theta[X_n|P] + (\mathbb{E}_\theta[X_n|P])^2 = nP(1-P) + n^2P^2$$

$$\text{Donc } \mathbb{E}_\theta[X_n^2] = \mathbb{E}_\theta[\mathbb{E}_\theta[X_n^2|P]]$$

$$= \mathbb{E}_\theta[nP(1-P) + n^2P^2]$$

$$= n^2\mathbb{E}_\theta[P^2] + n\mathbb{E}_\theta[P(1-P)]$$

$$\text{Alors } \text{Var}_\theta[X_n] = \mathbb{E}_\theta[X_n^2] - \mathbb{E}_\theta[X_n]^2$$

$$= n^2\mathbb{E}_\theta[P^2] + n\mathbb{E}_\theta[P(1-P)] - n^2\mathbb{E}_\theta[P]^2$$

$$= n\mathbb{E}_\theta[P(1-P)] + n^2 \text{Var}_\theta[P]$$

$$\mathbb{E}_\theta[P^2] = \int_{\mathbb{R}} x^2 f_{p,\theta}(x) dx = \int_0^1 \frac{1}{B(a,b)} x^{a+1} (1-x)^{b-1} dx = \frac{B(a+2,b)}{B(a,b)}$$

$$= \frac{\Gamma(a+2)\Gamma(b)}{\Gamma(a+b+2)} \times \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)}$$

$$= \frac{\Gamma(a+2)}{\Gamma(a+1)} \times \frac{\Gamma(a+1)}{\Gamma(a)} \times \frac{\Gamma(a+b)}{\Gamma(a+b+1)} \times \frac{\Gamma(a+b+1)}{\Gamma(a+b+2)}$$

$$= \frac{a(a+1)}{(a+b)(a+b+1)}$$

$$= \frac{a+1}{a+b+1} \mathbb{E}_\theta[P]$$

$$\text{Var}_\theta[P] = \mathbb{E}_\theta[P^2] - \mathbb{E}_\theta[P]^2$$

$$\text{Alors } \text{Var}_\theta[X_n] = n\mathbb{E}_\theta[P] - n\mathbb{E}_\theta[P^2] + n^2 \text{Var}_\theta[P]$$

$$= n\mathbb{E}_\theta[P] - n \frac{a+1}{a+b+1} \mathbb{E}_\theta[P] + n^2 \frac{a+1}{a+b+1} \mathbb{E}_\theta[P] - n^2 \mathbb{E}_\theta[P]^2$$

$$\text{Comme } n^2 \frac{a+1}{a+b+1} \mathbb{E}_\theta[P] - n^2 \mathbb{E}_\theta[P]^2$$

$$= n^2 \mathbb{E}_\theta[P] \left( \frac{a+1}{a+b+1} - \frac{a}{a+b} \right)$$

$$= n^2 \frac{b}{(a+b)(a+b+1)} \mathbb{E}_\theta[P]$$

$$\begin{aligned}
\text{On a } \text{Var}_{\theta}[X_n] &= n\mathbb{E}_{\theta}[P] - n\frac{a+1}{a+b+1}\mathbb{E}_{\theta}[P] + n^2\frac{b}{(a+b)(a+b+1)}\mathbb{E}_{\theta}[P] \\
&= \mathbb{E}_{\theta}[P](n - n\frac{a+1}{a+b+1} + n^2\frac{b}{(a+b)(a+b+1)}) \\
&= \mathbb{E}_{\theta}[P](n\frac{b}{a+b+1} + n^2\frac{b}{(a+b)(a+b+1)}) \\
&= \mathbb{E}_{\theta}[P](n\frac{b}{a+b} \times \frac{a+b}{a+b+1} + n^2\frac{b}{a+b} \times \frac{1}{a+b+1}) \\
&= \mathbb{E}_{\theta}[P] \times n \times (1 - \frac{a}{a+b})(\frac{a+b}{a+b+1} + n\frac{1}{a+b+1}) \\
&= \mathbb{E}_{\theta}[P] \times (n - n\mathbb{E}_{\theta}[P]) \times \frac{a+b+n}{a+b+1} \\
&= \frac{1}{n}\mathbb{E}_{\theta}[X_n] \times (n - \mathbb{E}_{\theta}[X_n]) \times (1 + \frac{n-1}{a+b+1})
\end{aligned}$$

## 2. Estimation des paramètres

1) Si  $n = 1$

D'après la question précédente

$$\text{On a } \mathbb{E}_{\theta}[X_1] = 1 \times \frac{a}{a+b} = p_0$$

$$\text{Et } \text{Var}_{\theta}[X_n] = 1 \times \mathbb{E}_{\theta}[X_1] \times (1 - \mathbb{E}_{\theta}[X_1]) \times 1 = p_0(1 - p_0)$$

$$\mathbb{P}_{\theta}(X_1 = 1) = p_0$$

$$\mathbb{P}_{\theta}(X_1 = 0) = 1 - p_0$$

Donc la loi de  $X_1$  ne dépend que  $p_0$

De plus on ne peut que estimer la valeur de  $p_0 = \frac{a}{a+b}$

Pour une valeur de  $p_0$ , on peut trouver plusieurs couples (a,b) qui vérifient  $p_0 = \frac{a}{a+b}$

D'où c'est impossible d'estimer la paramètre  $\theta = (a, b)$

2) Si  $n \geq 2$

On veut construire  $\widehat{\theta}_k$

$$\text{à partir de } M_k = \frac{1}{K} \sum_{k=1}^K X_n^{(k)} \text{ et } V_k = \frac{1}{K} \sum_{k=1}^K (X_n^{(k)})^2 - M_k^2$$

$$\text{On note } \widehat{\theta}_k = (\widehat{a}_k, \widehat{b}_k)$$

$M_k$  est l'estimateur de moyenne

$V_k$  est l'estimateur de variance

$$\text{Comme } \begin{cases} \mathbb{E}_{\theta}[X_n] = \frac{na}{a+b} \\ \text{Var}_{\theta}[X_n] = \frac{1}{n}\mathbb{E}_{\theta}[X_n] \times (n - \mathbb{E}_{\theta}[X_n]) \times (1 + \frac{n-1}{a+b+1}) \end{cases}$$

$$\text{On a } \begin{cases} M_k = \frac{n\widehat{a}_k}{\widehat{a}_k + \widehat{b}_k} \\ V_k = \frac{M_k}{n} \times (n - M_k) \times \left(1 + \frac{n-1}{\widehat{a}_k + \widehat{b}_k + 1}\right) \end{cases}$$

D'après méthode des moments

On peut en déduire  $\widehat{a}_k$  et  $\widehat{b}_k$  en fonction de  $M_k$  et  $V_k$

On obtient

$$\begin{cases} \widehat{a}_k = \frac{nM_k(n-M_k)-nV_k}{nV_k-M_k(n-M_k)} \times \frac{M_k}{n} \\ \widehat{b}_k = \frac{nM_k(n-M_k)-nV_k}{nV_k-M_k(n-M_k)} \times \left(1 - \frac{M_k}{n}\right) \end{cases}$$

3) D'après la loi de Grande nombre :

$$M_K = \frac{x_n^{(1)} + \dots + x_n^{(K)}}{K} \xrightarrow[p.s. \ k \rightarrow +\infty]{} E_\theta [x_n^{(1)}] = \frac{na}{a+b}$$

$$\begin{aligned} V_K = \frac{1}{K} \sum_{k=1}^K \left(x_n^{(k)}\right)^2 - M_K^2 &\xrightarrow[p.s. \ k \rightarrow +\infty]{} E_\theta \left[\left(x_n^{(1)}\right)^2\right] - E_\theta^2 [x_n^{(1)}] = \text{Var}_\theta (x_n^{(1)}) \\ &= \frac{nab}{(a+b)^2} \left(1 + \frac{n-1}{a+b+1}\right) \end{aligned}$$

$$\begin{aligned} \text{Donc } \widehat{a}_k &\xrightarrow[p.s. \ k \rightarrow +\infty]{} \frac{n \times \frac{na}{a+b} \times \frac{nb}{a+b} - \frac{n^2 ab}{(a+b)^2} \left(1 + \frac{n-1}{a+b+1}\right)}{\frac{n^2 ab}{(a+b)^2} \left(1 + \frac{n-1}{a+b+1}\right) - \frac{na}{a+b} \times \frac{nb}{a+b}} \times \frac{na}{a+b} \times \frac{1}{n} \\ &= \frac{(n-1) \left(1 - \frac{1}{a+b+1}\right)}{1 + \frac{n-1}{a+b+1} - 1} \times \frac{a}{a+b} = a \end{aligned}$$

$$\text{De même } \widehat{b}_k \xrightarrow[p.s. \ k \rightarrow +\infty]{} (a+b) \times \left(1 - \frac{a}{a+b}\right) = b$$

$$\text{Alors } \widehat{\theta}_k = (\widehat{a}_k, \widehat{b}_k) \xrightarrow[k \rightarrow +\infty]{} (a, b) = \theta$$

4) On regroupe  $x_n^{(k)}$   $k=1,2,\dots,K$

$\forall r \in \llbracket 0, n \rrbracket$  on note l'ensemble  $K_r = \{x_n^{(k)} : x_n^{(k)} = r\}$  alors  $|K_r| = N_r$ ,  $|K_r|$  représente

cardinal, donc  $\sum_{k=1}^K x_n^{(k)} = \sum_{r=0}^n \sum_{k=1}^{|K_r|} (x_n^{(k)} = r) = \sum_{r=0}^n \sum_{k=1}^{N_r} r = \sum_{r=0}^n r N_r$

$$\sum_{k=1}^K \left(x_n^{(k)}\right)^2 = \sum_{r=0}^n \sum_{k=1}^{|K_r|} \left(x_n^{(k)} = r\right)^2 = \sum_{r=0}^n \sum_{k=1}^{N_r} r^2 = \sum_{r=0}^n r^2 N_r$$

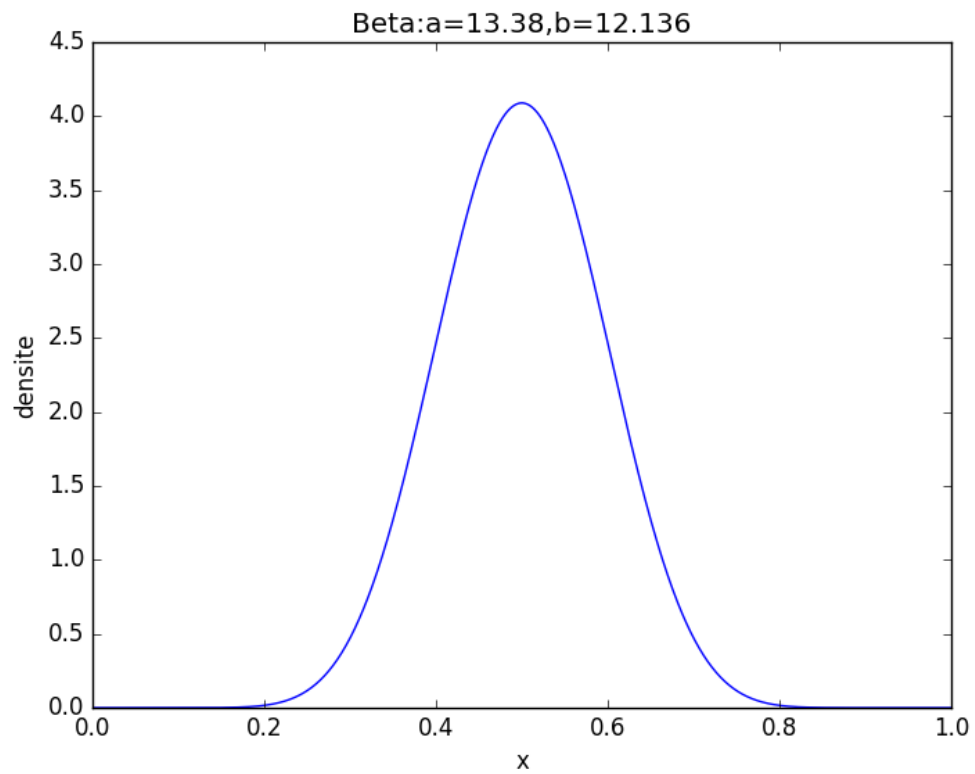
Valeur numérique :  $M_K = 2.621875$ ,  $V_K = 1.43514648437$ ,

$$\theta = (13,37995764, 12,13615216)$$

5) On utilise *scipy.stats.beta* pour tracer la densité et estime que  $p_1 = P_\theta (P > 1/2) =$

0.598253083758

Figure 1 la densité



6) On choisit la nombre de famille  $K=50, 200, 500$  et  $1000$ , et simule  $N=1000$  fois, tracer la densité, pour vérifier que  $\theta$  est asymptotiquement normale.

Figure 2 K=50

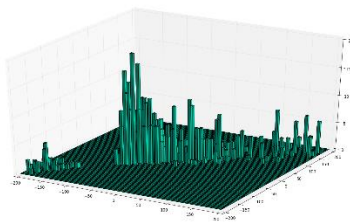


Figure 3 K=200

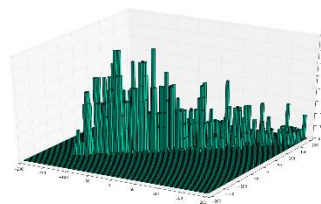


Figure 4 K=500

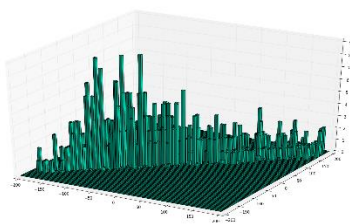
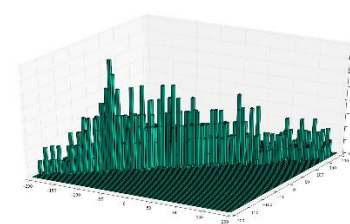


Figure 5 K=1000



Explication : pour obtenir la valeur moyenne  $M_k$  et la variante  $V_k$ , on choisit que la nombre de famille  $k$  soit 50, 200, 500 et 1000, dans ces  $k$  familles, tout d'abord, on stimulant la probabilité d'avoir un garçon dans une famille, et puis, on simule le nombre de garçon dans cette famille, donc dans tous les  $k$  familles, on obtient  $(x_5^{(1)}, \dots, x_5^{(k)})$ , en calculant  $M_k$  et  $V_k$ , or on a la relation entre  $(a, b)$  et  $(M_k, V_k)$ , on obtient  $(a, b)$  d'où  $\sqrt{K} (\widehat{\theta}_K - \theta)$ . On le répète  $N=1000$  fois, et puis on trace les résultats.

7) lorsque  $K=1000$ ,  $N=1000$ , on obtient la matrice de covariance est :

$$\begin{bmatrix} 635.214404365997 & 600.6343725306089 \\ 600.6343725306089 & 568.5224295398386 \end{bmatrix}$$

### 3. Estimation de la probabilité d'avoir un garçon

1) On veut calculer  $E_\theta[h(p)|x_n = k] = \frac{E_\theta[h(p)1_{\{x_n=k\}}]}{P_\theta(x_n=k)}$

En fait,  $P_\theta(x_n = k) = E_\theta[1_{\{x_n=k\}}] = E_\theta[E_\theta[1_{\{x_n=k\}}|P]] = E_\theta[P_\theta(x_n = k|P)]$

$$= E_\theta \left[ \binom{n}{k} P^k (1-P)^{n-k} \right]$$

$$E_\theta[h(p)1_{\{x_n=k\}}] = E_\theta \left[ E_\theta[h(P)1_{\{x_n=k\}}|P] \right] = E_\theta[h(p)P_\theta(x_n = k|P)]$$

$$= E_\theta[h(P) \binom{n}{k} P^k (1-P)^{n-k}]$$

$$\text{D'où } E_\theta[h(p)|x_n = k] = \frac{E_\theta[h(P) \binom{n}{k} P^k (1-P)^{n-k}]}{E_\theta[\binom{n}{k} P^k (1-P)^{n-k}]} = \frac{E_\theta[h(P) P^k (1-P)^{n-k}]}{E_\theta[P^k (1-P)^{n-k}]}$$

2) D'après les résultats de question 1 de la partie 1, on a :

$$E_\theta[P^v(1-P)^w] = \frac{B(v+a, w+b)}{B(a, b)}$$

On prend  $h: x \mapsto x$  bornée sur  $[0, 1]$

$$\text{Donc } E_\theta[P|x_n = k] = \frac{E_\theta[P^{k+1}(1-P)^{n-k}]}{E_\theta[P^k(1-P)^{n-k}]} = \frac{B(k+a+1, n-k+b)}{B(k+a, n-k+b)} = \frac{\Gamma(k+a+1)\Gamma(n-k+b)}{\Gamma(k+a+1+n-k+b)} \times$$

$$\frac{\Gamma(k+a+n-k+b)}{\Gamma(k+a)\Gamma(n-k+b)} = \frac{k+a}{a+b+n}$$

$$\text{De même, } E_\theta[P^2|x_n = k] = \frac{B(k+a+2, n-k+b)}{B(k+a, n-k+b)} = \frac{(k+a+1)(k+a)}{(a+b+n)(a+b+n+1)}$$

$$\text{Donc } \text{Var}[P|x_n = k] = E_\theta[P^2|x_n = k] - E_\theta^2[P|x_n = k] = \frac{k+a}{a+b+n} \times \frac{(a+k)(b+n-k)}{(a+b+n)(a+b+n+1)}$$

$$\text{Alors pour } P|x_n = k, E_\theta[P|x_n = k] = \frac{a+k}{a+b+n},$$

$$\text{Var}[P|x_n = k] = \frac{(a+k)^2(a+n-k)}{(a+b+n)^2(a+b+n+1)}$$

Pour  $P \quad E_{\theta}[P] = \frac{a}{a+b}, Var_{\theta}[P] = \frac{a^2 b}{(a+b)^2(a+b+1)}$

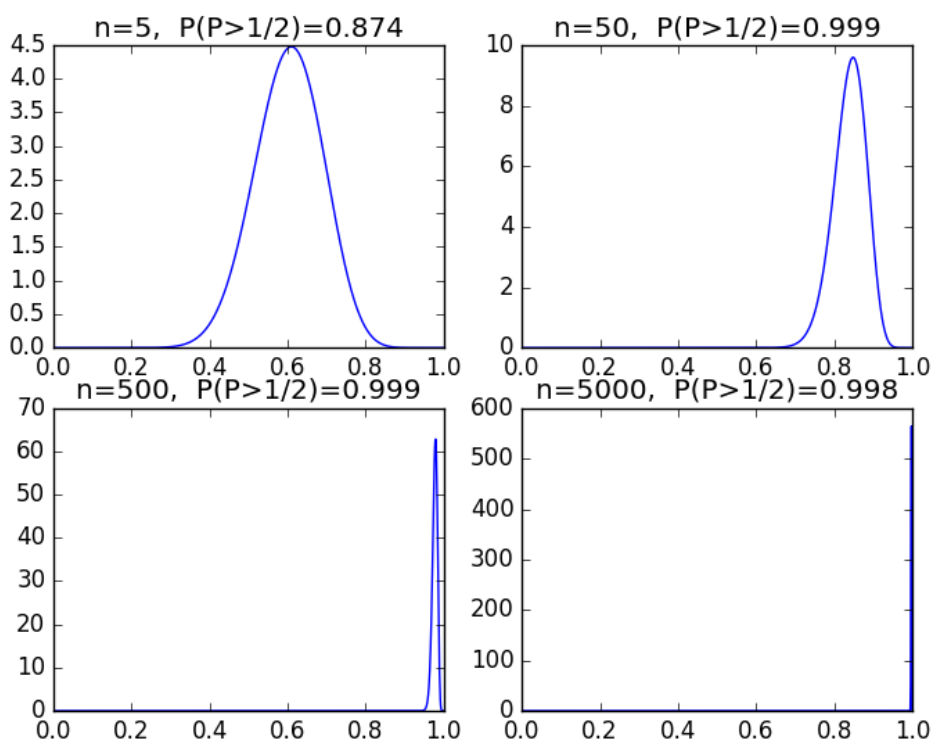
On fait analogue  $a + k \leftrightarrow a, b + n - k \leftrightarrow b$

Alors  $P|x_n = k$  suit une même loi que  $P$  avec la paramètre  $(a+k, b+n-k)$

Donc  $f_{P|x_n=k,\theta}(x) = \frac{1}{B(a+k, b+n-k)} x^{a+k-1} (1-x)^{b+n-k-1} \times 1_{]0,1[}(x)$

$$f_{P|x_n}(x) = \frac{1}{B(a+x_n, b+n-x_n)} x^{a+x_n-1} (1-x)^{b+n-x_n-1} 1_{]0,1[}(x)$$

3) Ci-dessous c'est la figure de densité pour  $n=5, 50, 500$  et  $5000$



4) Merci à Python, on a calculé  $P_{\theta}(P > 1/2|x_n = n)$  affiché ci-dessus.

Si les  $n$  précédents sont des garçons, alors le sexe du  $n+1$ -ème enfant est :

*garçon, si  $n$  est grand  
on ne peut pas prévoir, si  $n$  n'est pas grand*



Liste des programmes :

2-5.py pour la question 2.5

2-6.py pour la question 2.6

2-7.py pour la question 2.7

3-3et3-4.py pour la question 3.3 et la question 3.4