

RECHERCHE DE RÉGIONS GÉNOMIQUES CONSERVÉES

MODÈLE MATHÉMATIQUE ET TEST STATISTIQUE

S. Grusea, E. Pardoux¹ V. Lopez Rascol, P. Pontarotti²

¹L.A.T.P., Centre de Mathématique et d'Informatique
Université de Provence

²Laboratoire de Phylogénomique EA 3781 Évolution Biologique
Université de Provence

ALPHY/GTGC 23-24 janvier 2006

1 INTRODUCTION

- Contexte biologique
- Significativité

2 MODÈLE MATHÉMATIQUE

- Formulation du problème
- Test statistique
- Solution par la méthode de Monte Carlo

3 VERS UNE SOLUTION PLUS MATHÉMATIQUE

- Modèle mathématique simplifié
- Exemple : le cas sans familles de gènes
- Comportement asymptotique
- Conclusion

RÉGIONS GÉNOMIQUES CONSERVÉES

Intérêt :

RÉGIONS GÉNOMIQUES CONSERVÉES

Intérêt : inférer

- des relations évolutives entre espèces

RÉGIONS GÉNOMIQUES CONSERVÉES

Intérêt : inférer

- des relations évolutives entre espèces
- des gènes sous une pression fonctionnelle sélective.

RÉGIONS GÉNOMIQUES CONSERVÉES

Intérêt : inférer

- des relations évolutives entre espèces
- des gènes sous une pression fonctionnelle sélective.

RÉGIONS GÉNOMIQUES CONSERVÉES

Intérêt : inférer

- des relations évolutives entre espèces
- des gènes sous une pression fonctionnelle sélective.

RÉGION CONSERVÉE (CLUSTER CONSERVÉ DE GÈNES) :

deux régions chromosomiques ayant un certain nombre de gènes orthologues en commun.

RÉGIONS GÉNOMIQUES CONSERVÉES

Intérêt : inférer

- des relations évolutives entre espèces
- des gènes sous une pression fonctionnelle sélective.

RÉGION CONSERVÉE (CLUSTER CONSERVÉ DE GÈNES) :

deux régions chromosomiques ayant un certain nombre de gènes orthologues en commun.

Recherche de régions conservées :

RÉGIONS GÉNOMIQUES CONSERVÉES

Intérêt : inférer

- des relations évolutives entre espèces
- des gènes sous une pression fonctionnelle sélective.

RÉGION CONSERVÉE (CLUSTER CONSERVÉ DE GÈNES) :

deux régions chromosomiques ayant un certain nombre de gènes orthologues en commun.

Recherche de régions conservées :

- Chez une espèce A : région chromosomique fixée (région de référence)

RÉGIONS GÉNOMIQUES CONSERVÉES

Intérêt : inférer

- des relations évolutives entre espèces
- des gènes sous une pression fonctionnelle sélective.

RÉGION CONSERVÉE (CLUSTER CONSERVÉ DE GÈNES) :

deux régions chromosomiques ayant un certain nombre de gènes orthologues en commun.

Recherche de régions conservées :

- Chez une espèce A : région chromosomique fixée (région de référence)
- On cherche chez une autre espèce B des régions génomiques similaires.

RECHERCHE DE RÉGIONS CONSERVÉES

- 1 Trouver les orthologues des gènes de la région de référence dans l'espèce B.

RECHERCHE DE RÉGIONS CONSERVÉES

- 1 Trouver les orthologues des gènes de la région de référence dans l'espèce B.
- 2 Chercher les régions avec une grande densité de gènes orthologues.

RECHERCHE DE RÉGIONS CONSERVÉES

- ① Trouver les orthologues des gènes de la région de référence dans l'espèce B.
- ② Chercher les régions avec une grande densité de gènes orthologues.
- ③ Tester si elles sont **significatives** (vraiment 'conservées').

RECHERCHE DE RÉGIONS CONSERVÉES

- ① Trouver les orthologues des gènes de la région de référence dans l'espèce B.
- ② Chercher les régions avec une grande densité de gènes orthologues.
- ③ Tester si elles sont **significatives** (vraiment 'conservées').

RECHERCHE DE RÉGIONS CONSERVÉES

- 1 Trouver les orthologues des gènes de la région de référence dans l'espèce B.
- 2 Chercher les régions avec une grande densité de gènes orthologues.
- 3 Tester si elles sont **significatives** (vraiment 'conservées').

RÉGION CONSERVÉE SIGNIFICATIVE :

C'est très peu probable qu'elle soit apparue par hasard.

RECHERCHE DE RÉGIONS CONSERVÉES

- 1 Trouver les orthologues des gènes de la région de référence dans l'espèce B.
- 2 Chercher les régions avec une grande densité de gènes orthologues.
- 3 Tester si elles sont **significatives** (vraiment 'conservées').

RÉGION CONSERVÉE SIGNIFICATIVE :

C'est très peu probable qu'elle soit apparue par hasard.

Test statistique

- adapté à l'approche de type région de référence

RECHERCHE DE RÉGIONS CONSERVÉES

- 1 Trouver les orthologues des gènes de la région de référence dans l'espèce B.
- 2 Chercher les régions avec une grande densité de gènes orthologues.
- 3 Tester si elles sont **significatives** (vraiment 'conservées').

RÉGION CONSERVÉE SIGNIFICATIVE :

C'est très peu probable qu'elle soit apparue par hasard.

Test statistique

- adapté à l'approche de type région de référence
- qui prend en compte l'existence de familles de gènes orthologues.

L'IDÉE DU TEST STATISTIQUE

! Région conservée entre les deux espèces.

L'IDÉE DU TEST STATISTIQUE

! Région conservée entre les deux espèces.

? Est-elle significative ?

L'IDÉE DU TEST STATISTIQUE

! Région conservée entre les deux espèces.

? Est-elle significative ?

Hypothèse H_0 = hypothèse de génome aléatoire (hasard) :

L'IDÉE DU TEST STATISTIQUE

! Région conservée entre les deux espèces.

? Est-elle significative ?

Hypothèse H_0 = hypothèse de génome aléatoire (hasard) :

- calculer, sous cette hypothèse, la probabilité de trouver une telle région dans le génome B

L'IDÉE DU TEST STATISTIQUE

! Région conservée entre les deux espèces.

? Est-elle significative ?

Hypothèse H_0 = hypothèse de génome aléatoire (hasard) :

- calculer, sous cette hypothèse, la probabilité de trouver une telle région dans le génome B
- si suffisamment petite

L'IDÉE DU TEST STATISTIQUE

! Région conservée entre les deux espèces.

? Est-elle significative ?

Hypothèse H_0 = hypothèse de génome aléatoire (hasard) :

- calculer, sous cette hypothèse, la probabilité de trouver une telle région dans le génome B
- si suffisamment petite

L'IDÉE DU TEST STATISTIQUE

! Région conservée entre les deux espèces.

? Est-elle significative ?

Hypothèse H_0 = hypothèse de génome aléatoire (hasard) :

- calculer, sous cette hypothèse, la probabilité de trouver une telle région dans le génome B
- si suffisamment petite
⇒ région conservée **significative** !

1 INTRODUCTION

- Contexte biologique
- Significativité

2 MODÈLE MATHÉMATIQUE

- Formulation du problème
- Test statistique
- Solution par la méthode de Monte Carlo

3 VERS UNE SOLUTION PLUS MATHÉMATIQUE

- Modèle mathématique simplifié
- Exemple : le cas sans familles de gènes
- Comportement asymptotique
- Conclusion

FORMULATION MATHÉMATIQUE

Un génome = séquence ordonnée de gènes, sans séparation en chromosomes.

FORMULATION MATHÉMATIQUE

Un génome = séquence ordonnée de gènes, sans séparation en chromosomes.

La longueur d'une région génomique = le nombre de gènes.

FORMULATION MATHÉMATIQUE

Un génome = séquence ordonnée de gènes, sans séparation en chromosomes.

La longueur d'une région génomique = le nombre de gènes.

LES DONNÉES :

- m gènes dans la région de référence du génome A qui ont des orthologues dans le génome B

FORMULATION MATHÉMATIQUE

Un génome = séquence ordonnée de gènes, sans séparation en chromosomes.

La longueur d'une région génomique = le nombre de gènes.

LES DONNÉES :

- m gènes dans la région de référence du génome A qui ont des orthologues dans le génome B
- $\phi_i, i = 1, \dots, m$ les tailles des familles de gènes orthologues

FORMULATION MATHÉMATIQUE

Un génome = séquence ordonnée de gènes, sans séparation en chromosomes.

La longueur d'une région génomique = le nombre de gènes.

LES DONNÉES :

- m gènes dans la région de référence du génome A qui ont des orthologues dans le génome B
- $\phi_i, i = 1, \dots, m$ les tailles des familles de gènes orthologues
- N la taille du génome de l'espèce B.

PRENDRE EN COMPTE LES FAMILLES DE GÈNES

Un dénombrement un peu différent :

NOMBRE D'ORTHOLOGUES DANS UNE RÉGION

Dans une certaine région de B :

n_i orthologues de la i - ème famille, $i = 1, \dots, m$

PRENDRE EN COMPTE LES FAMILLES DE GÈNES

Un dénombrement un peu différent :

NOMBRE D'ORTHOLOGUES DANS UNE RÉGION

Dans une certaine région de B :

n_i orthologues de la i - ème famille, $i = 1, \dots, m$

\implies le 'nombre' d'orthologues dans cette région :

$$h = \sum_{i=1}^m \frac{n_i}{\phi_i}.$$

LE TEST STATISTIQUE (I)

DÉFINITIONS :

Pour chaque h possible :

- L_h = la longueur de la plus petite région dans B contenant h orthologues.

LE TEST STATISTIQUE (I)

DÉFINITIONS :

Pour chaque h possible :

- L_h = la longueur de la plus petite région dans B contenant h orthologues.
- $r_\beta(h)$ = la longueur critique pour qu'un cluster avec h orthologues soit significatif au seuil β :

$$\mathbb{P}(L_h \leq r_\beta(h)) \simeq \beta.$$

LE TEST STATISTIQUE (II)

On cherche un cluster significatif.

LE TEST STATISTIQUE (II)

On cherche un cluster significatif.

On ne connaît pas a priori le nombre h d'orthologues !

LE TEST STATISTIQUE (II)

On cherche un cluster significatif.

On ne connaît pas a priori le nombre h d'orthologues !

LE TEST :

- 1 Pour β fixé, trouver pour tous les h possibles les longueurs critiques correspondantes $r_\beta(h)$:

$$\mathbb{P}(L_h \leq r_\beta(h)) \simeq \beta.$$

LE TEST STATISTIQUE (II)

On cherche un cluster significatif.

On ne connaît pas a priori le nombre h d'orthologues !

LE TEST :

- 1 Pour β fixé, trouver pour tous les h possibles les longueurs critiques correspondantes $r_\beta(h)$:

$$\mathbb{P}(L_h \leq r_\beta(h)) \simeq \beta.$$

- 2 Trouver le bon β tel que

$$\mathbb{P}(\text{il existe } h \text{ t.q. } L_h \leq r_\beta(h)) \simeq \alpha (0,01).$$

ESTIMATION PAR MONTE CARLO

La loi exacte de L_h :

ESTIMATION PAR MONTE CARLO

La loi exacte de L_h : difficile à obtenir.

ESTIMATION PAR MONTE CARLO

La loi exacte de L_h : difficile à obtenir.

LA MÉTHODE DE MONTE CARLO :

On approche une moyenne théorique par la moyenne empirique.

ESTIMATION PAR MONTE CARLO

La loi exacte de L_h : difficile à obtenir.

LA MÉTHODE DE MONTE CARLO :

On approche une moyenne théorique par la moyenne empirique.

On approche une probabilité par la fréquence empirique.

ESTIMATION PAR MONTE CARLO

La loi exacte de L_h : difficile à obtenir.

LA MÉTHODE DE MONTE CARLO :

On approche une moyenne théorique par la moyenne empirique.

On approche une probabilité par la fréquence empirique.

- Pour chaque h , on estime la fonction de répartition de L_h :

$$F(r) = \mathbb{P}(L_h \leq r) \simeq \frac{|\{i = 1, \dots, M : l_{h,i} \leq r\}|}{M}$$

$l_{h,i}, i = 1, \dots, M$ simulations indép. de L_h , avec M suffisamment grand.

ESTIMATION PAR MONTE CARLO

La loi exacte de L_h : difficile à obtenir.

LA MÉTHODE DE MONTE CARLO :

On approche une moyenne théorique par la moyenne empirique.

On approche une probabilité par la fréquence empirique.

- Pour chaque h , on estime la fonction de répartition de L_h :

$$F(r) = \mathbb{P}(L_h \leq r) \simeq \frac{|\{i = 1, \dots, M : l_{h,i} \leq r\}|}{M}$$

$l_{h,i}, i = 1, \dots, M$ simulations indép. de L_h , avec M suffisamment grand.

- Pour β fixé, on calcule $r_\beta(h)$ pour chaque h .

ESTIMATION PAR MONTE CARLO

La loi exacte de L_h : difficile à obtenir.

LA MÉTHODE DE MONTE CARLO :

On approche une moyenne théorique par la moyenne empirique.

On approche une probabilité par la fréquence empirique.

- Pour chaque h , on estime la fonction de répartition de L_h :

$$F(r) = \mathbb{P}(L_h \leq r) \simeq \frac{|\{i = 1, \dots, M : l_{h,i} \leq r\}|}{M}$$

$l_{h,i}, i = 1, \dots, M$ simulations indép. de L_h , avec M suffisamment grand.

- Pour β fixé, on calcule $r_\beta(h)$ pour chaque h .
- De façon analogue, on trouve le bon β .

1 INTRODUCTION

- Contexte biologique
- Significativité

2 MODÈLE MATHÉMATIQUE

- Formulation du problème
- Test statistique
- Solution par la méthode de Monte Carlo

3 VERS UNE SOLUTION PLUS MATHÉMATIQUE

- Modèle mathématique simplifié
- Exemple : le cas sans familles de gènes
- Comportement asymptotique
- Conclusion

MODÈLE MATHÉMATIQUE SIMPLIFIÉ

- Le génome $B =$ l'intervalle continu $[0,1]$

MODÈLE MATHÉMATIQUE SIMPLIFIÉ

- Le génome $B =$ l'intervalle continu $[0,1]$
- Sous l'hypothèse de génome aléatoire :
les positions dans B des gènes orthologues = v.a. indép. unif. distr.
sur $[0,1]$.

MODÈLE MATHÉMATIQUE SIMPLIFIÉ

- Le génome $B =$ l'intervalle continu $[0,1]$
- Sous l'hypothèse de génome aléatoire :
les positions dans B des gènes orthologues = v.a. indép. unif. distr.
sur $[0,1]$.
- On remplace h par $\frac{h}{m}$, donc $h \in [0, 1]$ aussi.

LE CAS SANS FAMILLES DE GÈNES ($\phi_i = 1, \forall i$)

- Mesure de comptage renormalisée

$$\mu_m(I) = \frac{1}{m} \sum_{i=1}^m 1_{\{X_i \in I\}} = \frac{|\{i : X_i \in I\}|}{m},$$

$X_i, i = 1, \dots, m$ des v.a. indép. unif. distr. sur $[0,1]$.

LE CAS SANS FAMILLES DE GÈNES ($\phi_i = 1, \forall i$)

- Mesure de comptage renormalisée

$$\mu_m(I) = \frac{1}{m} \sum_{i=1}^m 1_{\{X_i \in I\}} = \frac{|\{i : X_i \in I\}|}{m},$$

$X_i, i = 1, \dots, m$ des v.a. indép. unif. distr. sur $[0,1]$.

- Donc

$$L_h = \inf\{|I| : \mu_m(I) \geq h\}.$$

LE CAS SANS FAMILLES DE GÈNES ($\phi_i = 1, \forall i$)

- Mesure de comptage renormalisée

$$\mu_m(I) = \frac{1}{m} \sum_{i=1}^m 1_{\{X_i \in I\}} = \frac{|\{i : X_i \in I\}|}{m},$$

$X_i, i = 1, \dots, m$ des v.a. indép. unif. distr. sur $[0,1]$.

- Donc

$$L_h = \inf\{|I| : \mu_m(I) \geq h\}.$$

LE CAS SANS FAMILLES DE GÈNES ($\phi_i = 1, \forall i$)

- Mesure de comptage renormalisée

$$\mu_m(I) = \frac{1}{m} \sum_{i=1}^m 1_{\{X_i \in I\}} = \frac{|\{i : X_i \in I\}|}{m},$$

$X_i, i = 1, \dots, m$ des v.a. indép. unif. distr. sur $[0,1]$.

- Donc

$$L_h = \inf\{|I| : \mu_m(I) \geq h\}.$$

PROBLÈME

Étudier la loi de L_h .

LE CAS SANS FAMILLES DE GÈNES ($\phi_i = 1, \forall i$)

- Mesure de comptage renormalisée

$$\mu_m(I) = \frac{1}{m} \sum_{i=1}^m 1_{\{X_i \in I\}} = \frac{|\{i : X_i \in I\}|}{m},$$

$X_i, i = 1, \dots, m$ des v.a. indép. unif. distr. sur $[0,1]$.

- Donc

$$L_h = \inf\{|I| : \mu_m(I) \geq h\}.$$

PROBLÈME

Étudier la loi de L_h .

Même dans ce cas, difficile à obtenir.

LE COMPORTEMENT ASYMPTOTIQUE DE L_h ($m \rightarrow \infty$)

La fonction de répartition de μ_m :

$$F_m(t) = \frac{1}{m} \sum_{i=1}^m 1_{\{X_i \leq t\}}.$$

LE COMPORTEMENT ASYMPTOTIQUE DE $L_h (m \rightarrow \infty)$

La fonction de répartition de μ_m :

$$F_m(t) = \frac{1}{m} \sum_{i=1}^m 1_{\{X_i \leq t\}}.$$

- Loi des grands nombres :

$$F_m(t) \xrightarrow{m \rightarrow \infty} t, \text{ uniformément en } t \in [0, 1], \text{ p.s..}$$

LE COMPORTEMENT ASYMPTOTIQUE DE $L_h (m \rightarrow \infty)$

La fonction de répartition de μ_m :

$$F_m(t) = \frac{1}{m} \sum_{i=1}^m 1_{\{X_i \leq t\}}.$$

- Loi des grands nombres :

$$F_m(t) \xrightarrow{m \rightarrow \infty} t, \text{ uniformément en } t \in [0, 1], \text{ p.s..}$$

- Théorème Limite Centrale :

$$\sqrt{m}(F_m(t) - t) \xrightarrow{m \rightarrow \infty} B_t, \quad t \in [0, 1]$$

convergence en loi fonctionnelle vers le pont brownien.

APPROXIMATION GAUSSIENNE

On obtient :

$$\mathbb{P}(L_h \leq r) \gtrsim \mathbb{P}\left(\bigcup_i A_i\right),$$

où les A_i impliquent des v.a. i.i.d. gaussiennes.

LA FORMULE DE CHUNG - ERDÖS

$$\mathbb{P}\left(\bigcup_i A_i\right) \geq \frac{(\sum_i \mathbb{P}(A_i))^2}{\sum_i \mathbb{P}(A_i) + \sum_{i \neq j} \mathbb{P}(A_i \cap A_j)}.$$

CONCLUSION

- Pour β fixé, on calcule $r_\beta(h)$ pour tout h .

CONCLUSION

- Pour β fixé, on calcule $r_\beta(h)$ pour tout h .
- Algorithme stochastique (Robbins - Monro) pour trouver le β t.q.

$\mathbb{P}(\text{il existe } h \text{ t.q. } L_h \leq r_\beta(h)) \simeq \alpha(0,01)$.

CONCLUSION

- Pour β fixé, on calcule $r_\beta(h)$ pour tout h .
- Algorithme stochastique (Robbins - Monro) pour trouver le β t.q.

$\mathbb{P}(\text{il existe } h \text{ t.q. } L_h \leq r_\beta(h)) \simeq \alpha(0,01)$.

CONCLUSION

- Pour β fixé, on calcule $r_\beta(h)$ pour tout h .
- Algorithme stochastique (Robbins - Monro) pour trouver le β t.q.

$$\mathbb{P}(\text{il existe } h \text{ t.q. } L_h \leq r_\beta(h)) \simeq \alpha(0,01).$$

Généralisation pour le cas avec des familles de gènes.

CONCLUSION

- Pour β fixé, on calcule $r_\beta(h)$ pour tout h .
- Algorithme stochastique (Robbins - Monro) pour trouver le β t.q.

$$\mathbb{P}(\text{il existe } h \text{ t.q. } L_h \leq r_\beta(h)) \simeq \alpha (0,01).$$

Généralisation pour le cas avec des familles de gènes.

PERSPECTIVES :

- Prendre en compte l'ordre de gènes.

CONCLUSION

- Pour β fixé, on calcule $r_\beta(h)$ pour tout h .
- Algorithme stochastique (Robbins - Monro) pour trouver le β t.q.

$$\mathbb{P}(\text{il existe } h \text{ t.q. } L_h \leq r_\beta(h)) \simeq \alpha(0,01).$$

Généralisation pour le cas avec des familles de gènes.

PERSPECTIVES :

- Prendre en compte l'ordre de gènes.
- Reconstruire les génomes ancestraux.

CONCLUSION

- Pour β fixé, on calcule $r_\beta(h)$ pour tout h .
- Algorithme stochastique (Robbins - Monro) pour trouver le β t.q.

$$\mathbb{P}(\text{il existe } h \text{ t.q. } L_h \leq r_\beta(h)) \simeq \alpha(0,01).$$

Généralisation pour le cas avec des familles de gènes.

PERSPECTIVES :

- Prendre en compte l'ordre de gènes.
- Reconstruire les génomes ancestraux.
- ...