

Estimation de paramètres par **Statistic Model-Checking** pour un modèle de croissance de forêt

Gilles & Guillaume
LS2N – Équipe VELO

16 juin 2022



- ▶ Les données de Paracou

- ▶ Les données de Paracou
- ▶ Le modèle de dynamique forestière de Kuznetsov & Aponina

- ▶ Les données de Paracou
- ▶ Le modèle de dynamique forestière de Kuznetsov & Aponina
- ▶ Implémentation du Statistic Model-Checking

Les données de Paracou

Paracou ? Qu'est-ce que c'est ? C'est où ?

1/ Qu'est-ce que c'est ?

Paracou ? Qu'est-ce que c'est ? C'est où ?

1/ Qu'est-ce que c'est ? ~→ Un site de recherche géré par le CIRAD (Géraldine Derroire).



Paracou ? Qu'est-ce que c'est ? C'est où ?

1/ Qu'est-ce que c'est ? ↪ Un site de recherche géré par le CIRAD (Géraldine Derroire).



2/ C'est où ?

Paracou ? Qu'est-ce que c'est ? C'est où ?

1/ Qu'est-ce que c'est ? → Un site de recherche géré par le CIRAD (Géraldine Derroire).



2/ C'est où ? → Quelque part en Guyane...

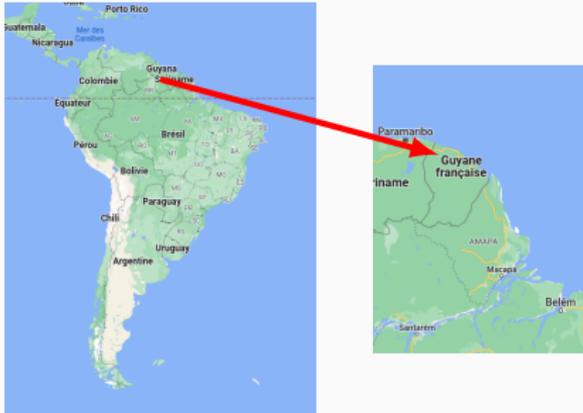


Paracou ? Qu'est-ce que c'est ? C'est où ?

1/ Qu'est-ce que c'est ? ↗ Un site de recherche géré par le CIRAD (Géraldine Derroire).



2/ C'est où ? ↗ Quelque part en Guyane...

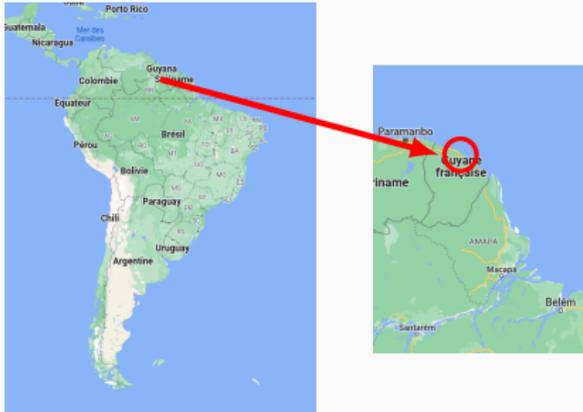


Paracou ? Qu'est-ce que c'est ? C'est où ?

1/ Qu'est-ce que c'est ? ↗ Un site de recherche géré par le CIRAD (Géraldine Derroire).



2/ C'est où ? ↗ Quelque part en Guyane...

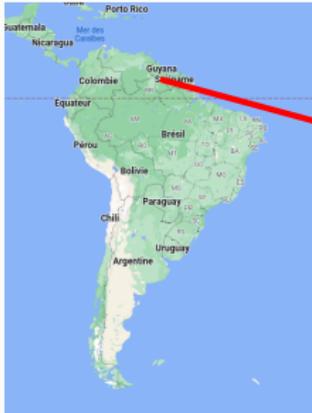


Paracou ? Qu'est-ce que c'est ? C'est où ?

1/ Qu'est-ce que c'est ? → Un site de recherche géré par le CIRAD (Géraldine Derroire).



2/ C'est où ? → Quelque part en Guyane...

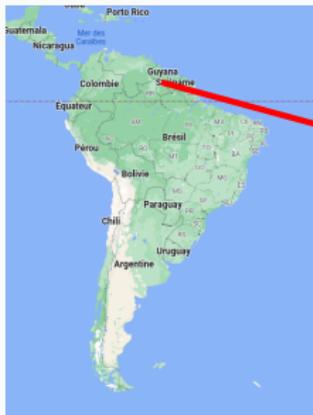


Paracou ? Qu'est-ce que c'est ? C'est où ?

1/ Qu'est-ce que c'est ? ↗ Un site de recherche géré par le CIRAD (Géraldine Derroire).



2/ C'est où ? ↗ Quelque part en Guyane...



Paracou ? Qu'est-ce que c'est ? C'est où ?

1/ Qu'est-ce que c'est ? ↗ Un site de recherche géré par le CIRAD (Géraldine Derroire).



2/ C'est où ? ↗ Quelque part en Guyane...



Présentation des données

- ▶ Depuis environ 40 ans, des données sont relevées sur une quinzaine de parcelles du site de Paracou.

Présentation des données

- ▶ Depuis environ 40 ans, des données sont relevées sur une quinzaine de parcelles du site de Paracou.
- ▶ Surface des parcelles : environ 6 hectares.

Présentation des données

- ▶ Depuis environ 40 ans, des données sont relevées sur une quinzaine de parcelles du site de Paracou.
- ▶ Surface des parcelles : environ 6 hectares.
- ▶ Certaines parcelles sont exploitées, d'autres ne le sont pas.

Présentation des données

- ▶ Depuis environ 40 ans, des données sont relevées sur une quinzaine de parcelles du site de Paracou.
- ▶ Surface des parcelles : environ 6 hectares.
- ▶ Certaines parcelles sont exploitées, d'autres ne le sont pas.
- ▶ Chaque arbre est répertorié (vivant/mort, essence, taille...)

Présentation des données

- ▶ Depuis environ 40 ans, des données sont relevées sur une quinzaine de parcelles du site de Paracou.
- ▶ Surface des parcelles : environ 6 hectares.
- ▶ Certaines parcelles sont exploitées, d'autres ne le sont pas.
- ▶ Chaque arbre est répertorié (vivant/mort, essence, taille...)



Présentation des données

- ▶ Depuis environ 40 ans, des données sont relevées sur une quinzaine de parcelles du site de Paracou.
- ▶ Surface des parcelles : environ 6 hectares.
- ▶ Certaines parcelles sont exploitées, d'autres ne le sont pas.
- ▶ Chaque arbre est répertorié (vivant/mort, essence, taille...)

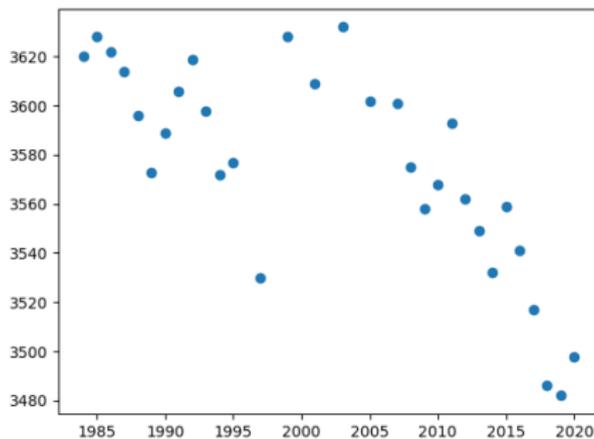


- ▶ Exemple : nombre d'arbres sur la parcelle 6 en fonction du temps, depuis 1983.

Visualisation des données

- ▶ Exemple : nombre d'arbres sur la parcelle 6 en fonction du temps, depuis 1983.

Parcelle 6 : ($\#$ arbres)(t)



► Questions

▶ Questions

~→ *Peut-on modéliser la croissance de la forêt ?*

► Questions

~→ *Peut-on modéliser la croissance de la forêt ?*

~→ *Si oui, comment ?*

► Questions

~→ *Peut-on modéliser la croissance de la forêt ?*

~→ *Si oui, comment ?*

~→ *Si on y arrive, quelles informations en tirer ?*

Le modèle de dynamique forestière de Kuznetsov & Aponina

Présentation du modèle

Le modèle d'Antonovsky et Korzukhin (1990) est donné par :

$$\begin{cases} \dot{u} = \rho v - \gamma(v)u - fu, \\ \dot{v} = fu - hv. \end{cases}$$

Présentation du modèle

Le modèle d'Antonovsky et Korzukhin (1990) est donné par :

$$\begin{cases} \dot{u} = \rho v - \gamma(v)u - fu, \\ \dot{v} = fu - hv. \end{cases}$$

u	\rightsquigarrow	arbres jeunes,
v	\rightsquigarrow	arbres âgés,
ρv	\rightsquigarrow	natalité,
fu	\rightsquigarrow	vieillissement,
hv	\rightsquigarrow	mortalité,
$\gamma(v)u$	\rightsquigarrow	compétition,
$\gamma(v)$	$=$	$a(v - b)^2 + c.$

Présentation du modèle

Le modèle d'Antonovsky et Korzukhin (1990) est donné par :

$$\begin{cases} \dot{u} = \rho v - \gamma(v)u - fu, \\ \dot{v} = fu - hv. \end{cases}$$

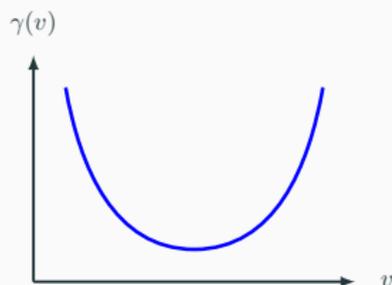


Figure 1 – Allure de $\gamma(v)$.

Trois régimes de paramètres

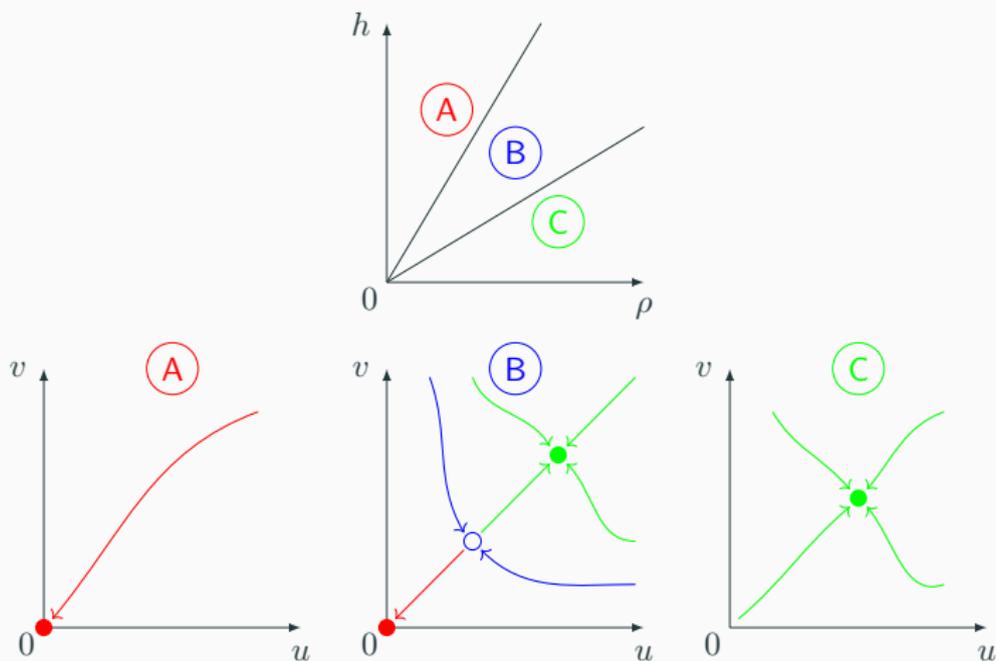


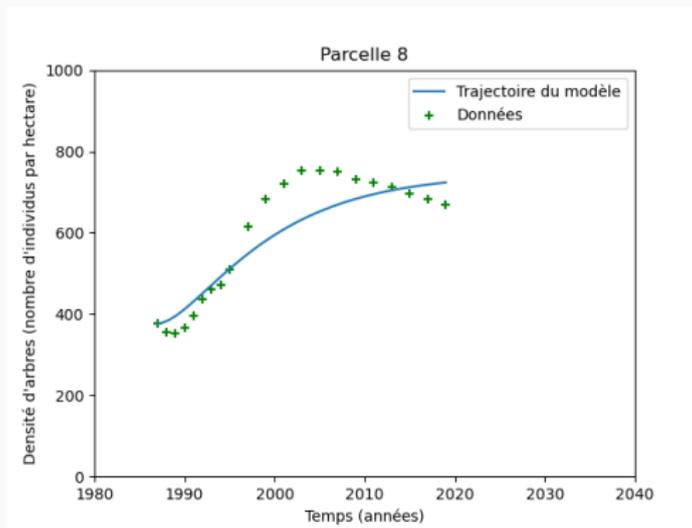
Figure 2 – Trois dynamiques possibles pour le modèle d'Antonovskiy et Korzukhin.

Modification du terme de vieillissement

- ▶ Les premières simulations du modèle montrent que sa dynamique reproduit mal l'évolution de certaines parcelles.

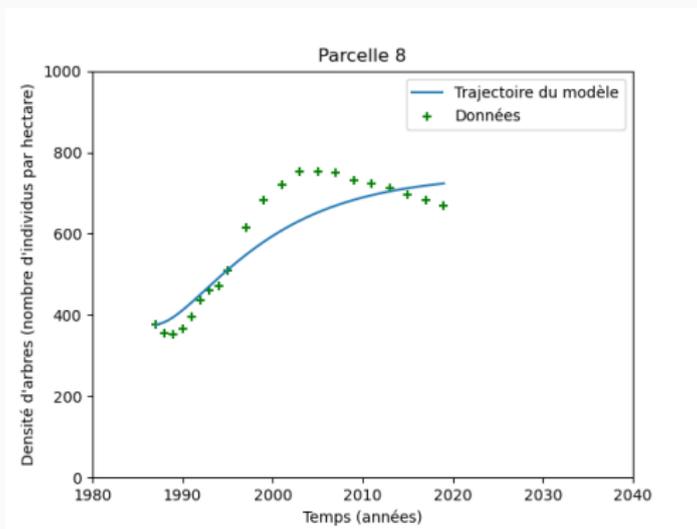
Modification du terme de vieillissement

- ▶ Les premières simulations du modèle montrent que sa dynamique reproduit mal l'évolution de certaines parcelles.



Modification du terme de vieillissement

- ▶ Les premières simulations du modèle montrent que sa dynamique reproduit mal l'évolution de certaines parcelles.



- ▶ Intuition : le terme de vieillissement $\pm fu$ n'est pas satisfaisant.

Modèle avec vieillissement non linéaire

- ▶ On pressent que le vieillissement dépend de la densité d'arbres.

Modèle avec vieillissement non linéaire

- ▶ On pressent que le vieillissement dépend de la densité d'arbres.
- ▶ On remplace donc $\pm fu$ par

$$fu \times v \times (T_{max} - (u + v)).$$

Modèle avec vieillissement non linéaire

- ▶ On pressent que le vieillissement dépend de la densité d'arbres.
- ▶ On remplace donc $\pm fu$ par

$$fu \times v \times (T_{max} - (u + v)).$$

- ▶ $\times v \rightsquigarrow$ le vieillissement diminue si v diminue.

Modèle avec vieillissement non linéaire

- ▶ On pressent que le vieillissement dépend de la densité d'arbres.
- ▶ On remplace donc $\pm fu$ par

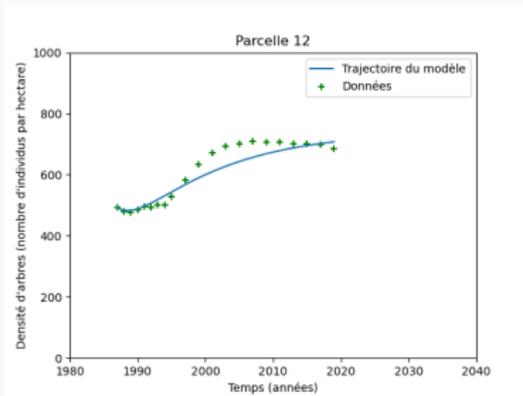
$$fu \times v \times (T_{max} - (u + v)).$$

- ▶ $\times v \rightsquigarrow$ le vieillissement diminue si v diminue.
- ▶ $\times (T_{max} - (u + v)) \rightsquigarrow$ le vieillissement est pondéré par la capacité maximale.

Est-ce que ça marche ?

► Résultats sur la parcelle 12.

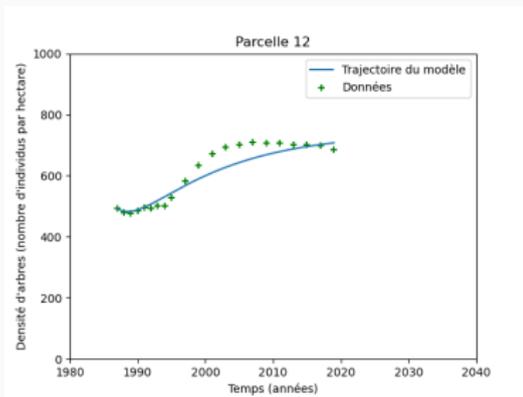
Vieillessement linéaire



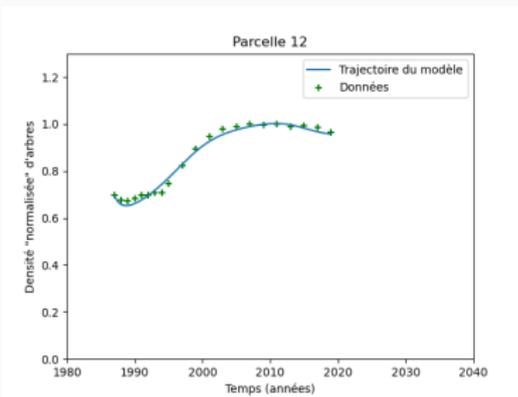
Est-ce que ça marche ?

► Résultats sur la parcelle 12.

Viellissement linéaire



Viellissement non linéaire



Parcelle-XX.csv : 85 - 125k entrées. Données manquantes (zone inaccessible certaines années).

Traitement préliminaire des données - Correction

Parcelle-XX.csv : 85 - 125k entrées. Données manquantes (zone inaccessible certaines années).

Si on se contentait de compter, une interpolation aurait pu suffire, mais les données sont plus riches. Donc pas de supposition sur les données (variations = incertitude).

Traitement préliminaire des données - Correction

Parcelle-XX.csv : 85 - 125k entrées. Données manquantes (zone inaccessible certaines années).

Si on se contentait de compter, une interpolation aurait pu suffire, mais les données sont plus riches. Donc pas de supposition sur les données (variations = incertitude).

Seule supposition raisonnable : Un arbre est vivant entre deux dates où il est recensé en vie.

- 13 - 1000 lignes par parcelle, (peu significatif à part sur la 15).
- Ne suffit pas à lisser complètement (arbres nés ou morts sur une date sans recensement)

Pour étudier la reprise après bucheronnage :

- parcelles 3, 4, 8 et 10
- sur la même période (1988 - 2019)
- exprimées en densité d'arbres /ha
- normalisé ($T_{\max} = 1 = 800$ arbres/ha)

Le modèle fit les données (et à quel point)

Le modèle fit les données (et à quel point)

Qualitatif

- nombre d'outliers (seuil de distance = tunnel) inférieur à une valeur

Le modèle fit les données (et à quel point)

Qualitatif

- nombre d'outliers (seuil de distance = tunnel) inférieur à une valeur

Quantitatif

- distance moyenne simu-data point à point (en % ...)
- si non fit : % de la donnée atteinte

Le modèle fit les données (et à quel point)

Qualitatif

- nombre d'outliers (seuil de distance = tunnel) inférieur à une valeur

Quantitatif

- distance moyenne simu-data point à point (en % ...)
- si non fit : % de la donnée atteinte

Meilleure distance \neq Minimum d'outliers

Objectifs

- Trouver jeu de paramètres satisfaisant le plus possible la propriété
- Cartographier l'espace de paramètres (minimums locaux éloignés, corrélations. . .)

Objectifs

- Trouver jeu de paramètres satisfaisant le plus possible la propriété
- Cartographier l'espace de paramètres (minimums locaux éloignés, corrélations. . .)

La contrainte de temps impose des stratégies opposées

- Propriété forte pour évacuer rapidement les zones de l'espace où elle n'est pas vérifiée.
- Propriété plus faible pour conserver suffisamment d'information pour chaque zone

Objectifs

- Trouver jeu de paramètres satisfaisant le plus possible la propriété
- Cartographier l'espace de paramètres (minimums locaux éloignés, corrélations. . .)

La contrainte de temps impose des stratégies opposées

- Propriété forte pour évacuer rapidement les zones de l'espace où elle n'est pas vérifiée.
- Propriété plus faible pour conserver suffisamment d'information pour chaque zone

Variabilité de la distance au sein d'une cellule : estimer le besoin de subdiviser pour

- chercher un meilleur optimum
- obtenir une moyenne plus représentative par cellule

Objectifs

- Trouver jeu de paramètres satisfaisant le plus possible la propriété
- Cartographier l'espace de paramètres (minimums locaux éloignés, corrélations. . .)

La contrainte de temps impose des stratégies opposées

- Propriété forte pour évacuer rapidement les zones de l'espace où elle n'est pas vérifiée.
- Propriété plus faible pour conserver suffisamment d'information pour chaque zone

Variabilité de la distance au sein d'une cellule : estimer le besoin de subdiviser pour

- chercher un meilleur optimum
- obtenir une moyenne plus représentative par cellule

Compromis : conserver le meilleur jeu de paramètres en plus de moyenne et var pour chaque cellule

Discrétisation.

- Nombre de cellules $\prod_{i=1}^n step_i$
- Nombre de simulations (evt. minoré par maxSPRT)
- Nombre de pas d'intégration (evt minoré par un cut distance ou nb outliers)
- Ordre de la fonction d'intégration (rk4 ici)

1. écrire le modèle en C++
2. Définir l'espace de paramètres (range et nb steps), configurer SPRT
3. fournir le ou les tunnels de données (données, width, nb_outliers, dt...)
4. build, run, wait...
5. Analyser/Filtrer les résultats
6. Reconstruire les simus et convertir dans un format plus exploitable pour les graphiques

Modèle C++

```
state_type modele(const double x_young, const double y_old,
                 const double a, const double b, const double c,
                 const double birth_rate, const double aging_rate, const double mortality)
{
    double d_young, d_old;
    const double aging(aging_rate * x_young * y_old * (1 - x_young - y_old));
    d_young = birth_rate * y_old - (a * (y_old - b) * (y_old - b) + c) * x_young - aging;
    d_old = aging - mortality * y_old;
    return {d_young, d_old};
}
```

Modèle C++

```
state_type modele(const double x_young, const double y_old,
                 const double a, const double b, const double c,
                 const double birth_rate, const double aging_rate, const double mortality)
{
    double d_young, d_old;
    const double aging(aging_rate * x_young * y_old * (1 - x_young - y_old));
    d_young = birth_rate * y_old - (a * (y_old - b) * (y_old - b) + c) * x_young - aging;
    d_old = aging - mortality * y_old;
    return {d_young, d_old};
}

constexpr state_type zero()
{
    return {0.0, 0.0};
}

// evaluate data in order to compare it to the tunnel
// for more flexibility, should return an eval_type
double eval_state(const state_type x) { return x[0] + x[1]; }

state_type init_state(const std::vector<double> &data, const params_type &)
{
    return {data[0] * 0.7, data[0] * 0.3}; // how is it determined ? maybe using params
}
```

Paramètres C++

```
ParamRanges<params_size> config({
    Param("a", 0., 1.0, 100),
    Param("b", 0, 1.0, 10),
    Param("c", 0, 1.0, 10),
    Param("f", 0, 1.0, 10),
    Param("r", 0, 1.0, 10),
    Param("h", 0, 1.0, 50),
});

SPRT sprt(sprt_alpha, sprt_beta, sprt_gamma, sprt_theta, sprt_nbsim);
```

Run

```
 1 [||||| 2.6%] 17 [|||||100.0%] 33 [|||||100.0%] 49 [|||||100.0%]
 2 [|||||100.0%] 18 [|||||100.0%] 34 [|||||100.0%] 50 [|||||100.0%]
 3 [|||||100.0%] 19 [|||||100.0%] 35 [|||||100.0%] 51 [|||||100.0%]
 4 [|||||100.0%] 20 [|||||100.0%] 36 [|||||100.0%] 52 [|||||100.0%]
 5 [|||||100.0%] 21 [|||||100.0%] 37 [|||||100.0%] 53 [|||||100.0%]
 6 [|||||100.0%] 22 [|||||100.0%] 38 [|||||100.0%] 54 [|||||100.0%]
 7 [|||||100.0%] 23 [|||||100.0%] 39 [|||||100.0%] 55 [|||||100.0%]
 8 [|||||100.0%] 24 [|||||100.0%] 40 [|||||100.0%] 56 [|||||100.0%]
 9 [|||||100.0%] 25 [|||||100.0%] 41 [|||||100.0%] 57 [|||||100.0%]
10 [|||||100.0%] 26 [|||||100.0%] 42 [|||||100.0%] 58 [|||||100.0%]
11 [|||||100.0%] 27 [|||||100.0%] 43 [|||||100.0%] 59 [|||||100.0%]
12 [|||||100.0%] 28 [|||||100.0%] 44 [|||||100.0%] 60 [|||||100.0%]
13 [|||||100.0%] 29 [|||||100.0%] 45 [|||||100.0%] 61 [|||||100.0%]
14 [|||||100.0%] 30 [|||||100.0%] 46 [|||||100.0%] 62 [|||||100.0%]
15 [|||||100.0%] 31 [|||||100.0%] 47 [|||||100.0%] 63 [|||||100.0%]
16 [|||||100.0%] 32 [|||||100.0%] 48 [|||||100.0%] 64 [|||||100.0%]
Mem[||||| 2.06G/62.6G] Tasks: 204, 458 thr; 64 running
Swp[||||| 39.8M/128G] Load average: 16.12 4.00 1.34
Uptime: 71 days, 18:47:41
```

PID	USER	CPU	PRI	NI	VIRT	RES	SHR	S	CPUN	MEM%	TIME+	Command
0422	ardourel-	5	20	0	14060	1824	1680	R	100.	0.0	0:09.23	a.out 3 63 ifile res1_test 63 ifile
0423	ardourel-	6	20	0	14060	1832	1688	R	100.	0.0	0:09.23	a.out 4 63 ifile res1_test 63 ifile
0424	ardourel-	7	20	0	14060	1868	1724	R	100.	0.0	0:09.23	a.out 5 63 ifile res1_test 63 ifile
0429	ardourel-	12	20	0	14060	1828	1684	R	100.	0.0	0:09.22	a.out 10 63 ifile res1_test 63 ifile
0434	ardourel-	17	20	0	14060	1904	1764	R	100.	0.0	0:09.22	a.out 15 63 ifile res1_test 63 ifile

!help F2Setup F3Search F4Filter F5Tree F6SortB F7Nice F8Kill F9Quit F10Quit

Data mining ;)

Analyser les résultats, Extraire les plus pertinents (nb outliers, distances)

Actuellement : grep, cut, wc...

Reconstruire les simus

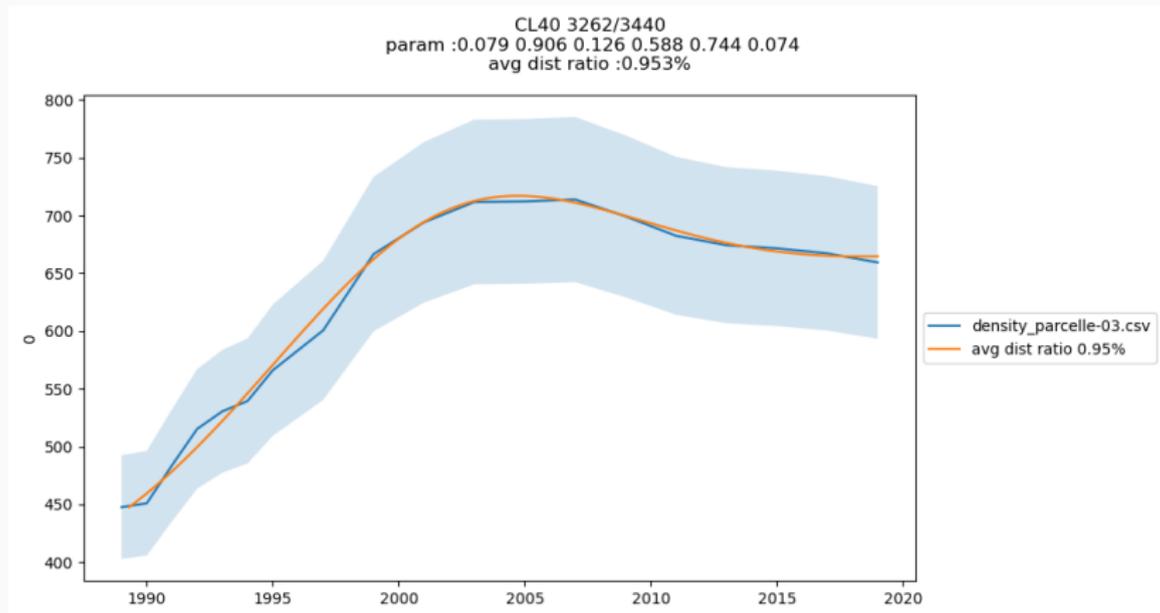
Petit prog C++ utilisant le modèle, les résultats et le tunnel pour générer du json

Demo visu résultats

- Parcelles individuelles
- Groupe de parcelles
- *Exploiter la biodiversité*

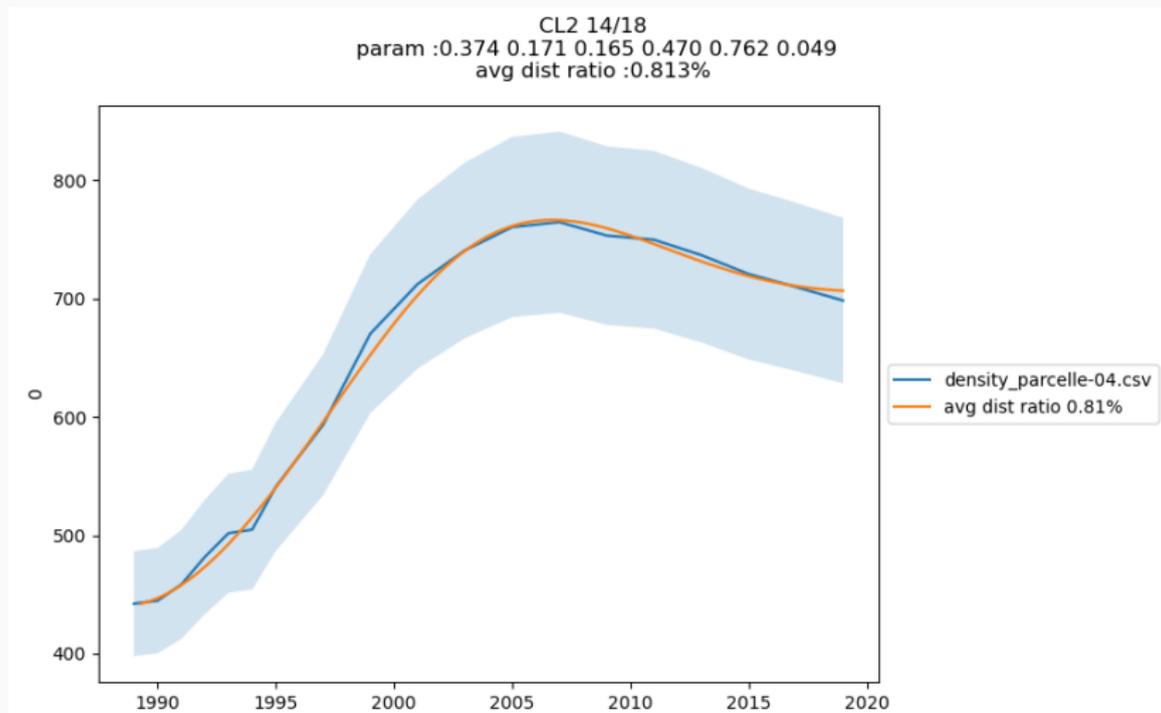
Visu résultats mono parcelle - Captures de la démo

Parcelles individuelles - Capture



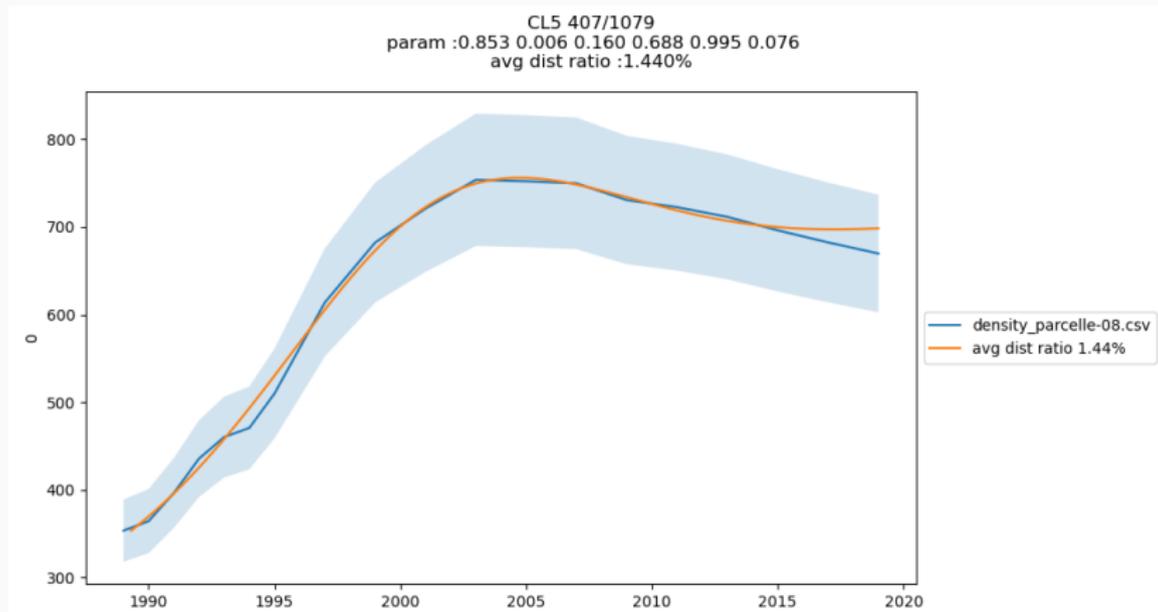
Visu résultats mono parcelle - Captures de la démo

Parcelles individuelles - Capture



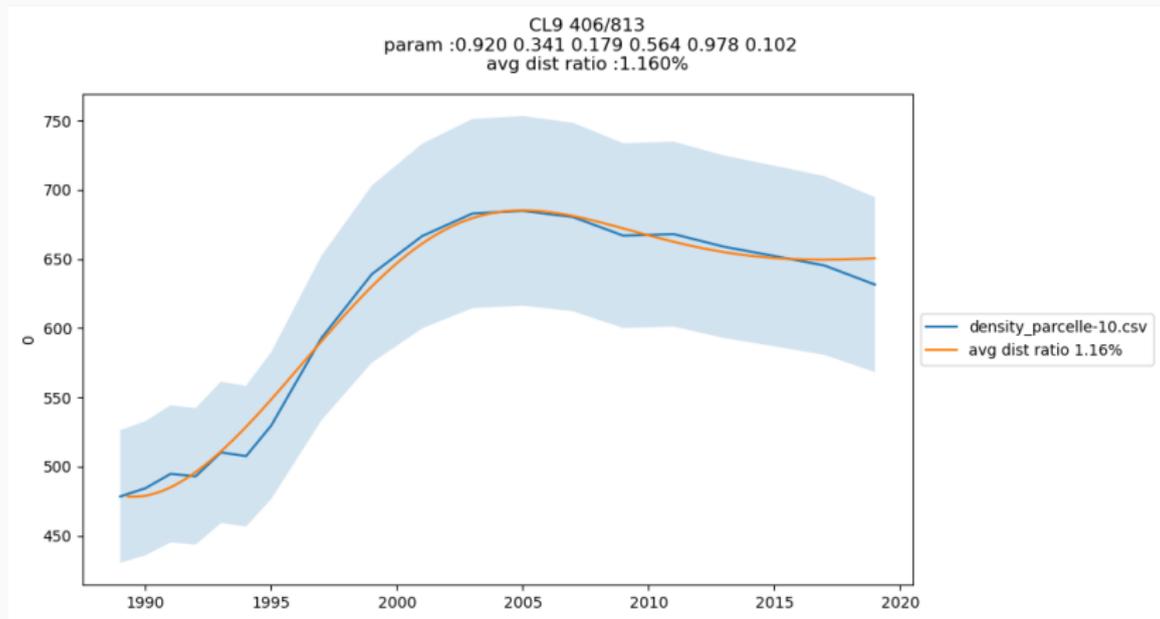
Visu résultats mono parcelle - Captures de la démo

Parcelles individuelles - Capture

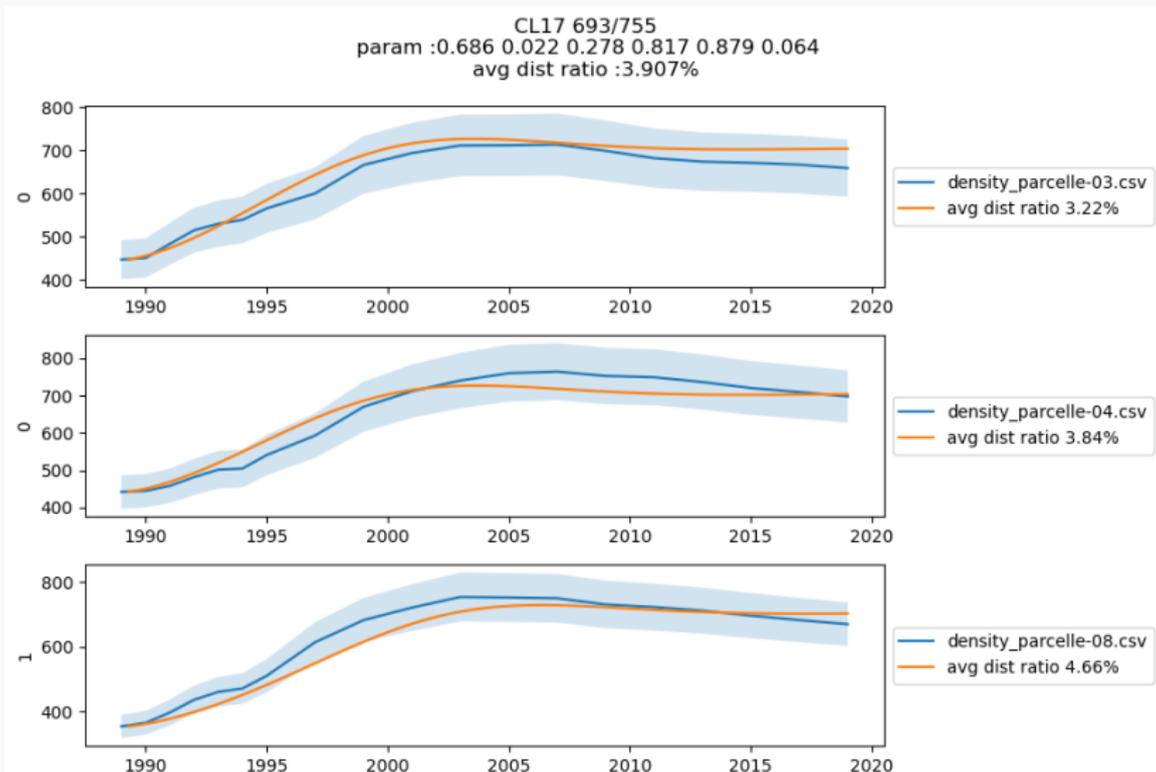


Visu résultats mono parcelle - Captures de la démo

Parcelles individuelles - Capture

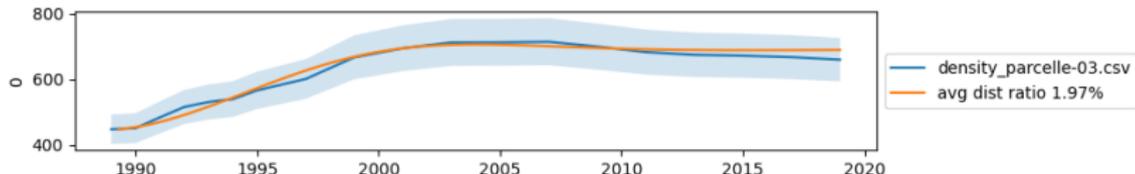
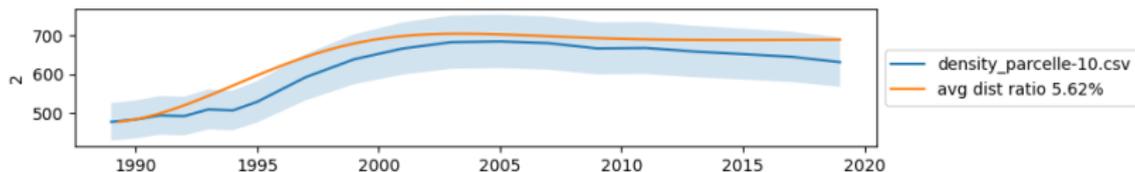
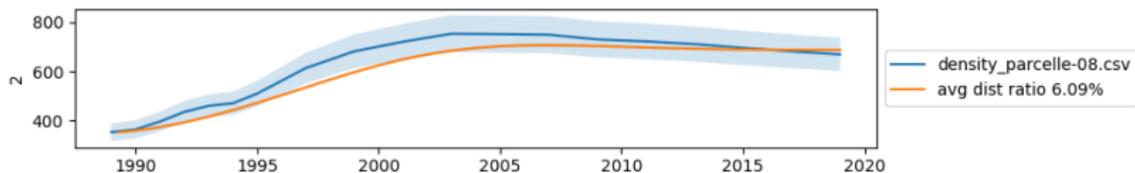
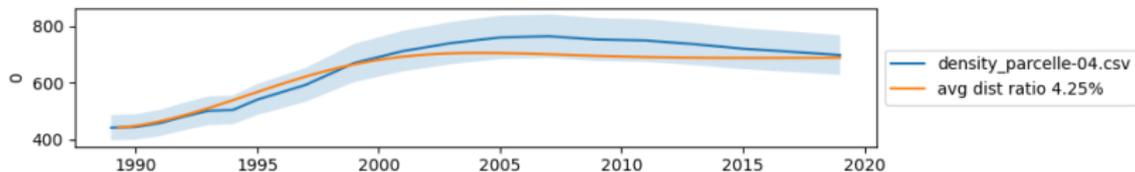


Visu résultats 3 parcelles - Captures de la démo

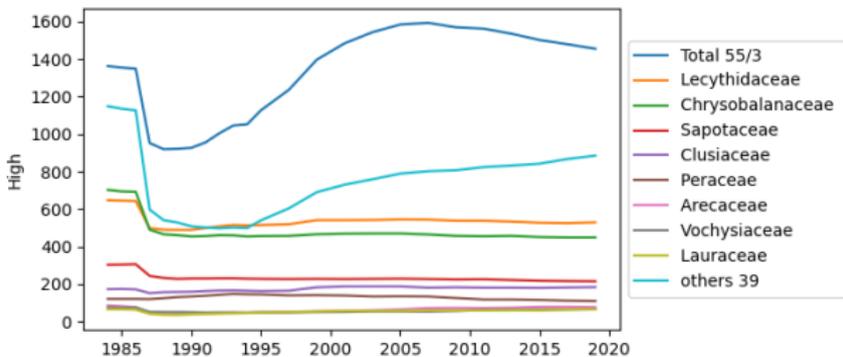
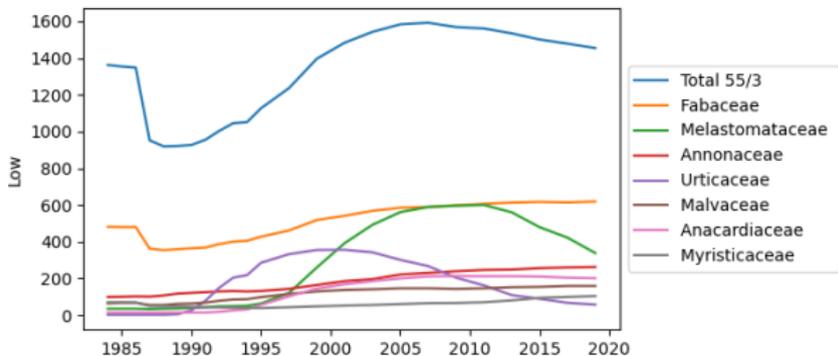


Visu résultats 4 parcelles - Captures de la démo

CL1 29/40
param :0.999 0.052 0.298 0.851 0.954 0.078
avg dist ratio :4.482%



Exploration de la variabilité par famille : parcelle 4



TODO

- Refactor, plus de généricité, intégration d'outils (AMICI...)
- Exploiter la variabilité :

TODO

- Refactor, plus de généricité, intégration d'outils (AMICI...)
- Exploiter la variabilité :
 - analyse de sensibilité pour déterminer un découpage
 - a posteriori sur cellules intéressantes ?
 - durant le run : sur quel critère, à quel coût ? quasi aléatoire plus pertinent ?

TODO

- Refactor, plus de généricité, intégration d'outils (AMICI...)
- Exploiter la variabilité :
 - analyse de sensibilité pour déterminer un découpage
 - a posteriori sur cellules intéressantes ?
 - durant le run : sur quel critère, à quel coût ? quasi aléatoire plus pertinent ?
- Estimation du temps (run sans cut, allocation de budget temps)

TODO

- Refactor, plus de généricité, intégration d'outils (AMICI...)
- Exploiter la variabilité :
 - analyse de sensibilité pour déterminer un découpage
 - a posteriori sur cellules intéressantes ?
 - durant le run : sur quel critère, à quel coût ? quasi aléatoire plus pertinent ?
- Estimation du temps (run sans cut, allocation de budget temps)
- ajustements dynamiques (SPRT, stringence propriété ? : risqué)