

INITIATION AUX TESTS STATISTIQUES ET AUX MODÈLES LINÉAIRES

Corentin M. Barbu

INRAE et Ecole Doctorale ABIES (réseau ADUM)

10 février 2025

Section 1

INTRODUCTION

Subsection 1

PRÉSENTATION

Introduction

Présentation

Modélisation
statistique

CORENTIN M. BARBU

- ▶ 2006-2014, épidémiologie
- ▶ 2014- à l'INRAE, contrôle des maladies et ravageurs (GC)

OBJECTIF DU COURS

Donner les clés conceptuelles et méthodologiques pour réaliser des tests statistiques et modèles linéaires de manière rigoureuse (et lire des articles scientifiques de manière critique).

Bases de stats

Corentin Barbu

Introduction

**Modélisation
statistique**

Bases de régressions
linéaires

Statistiques dans
l'espace et
auto-corrélation

Cas d'étude statistique:
Impact de la distance aux
bois sur les méligethes et
sur leurs parasoïdes

Cas d'étude statistique:
approche paysage
multi-bio-agresseurs

Application

Section 2

MODÉLISATION STATISTIQUE

Introduction

Modélisation statistique

Bases de régressions linéaires

Statistiques dans l'espace et auto-corrélation

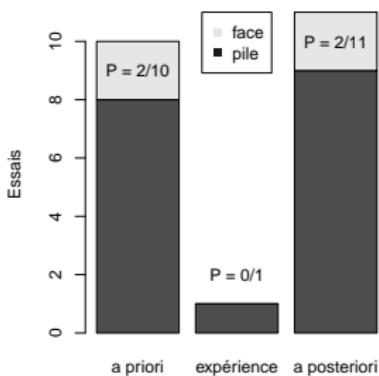
Cas d'étude statistique: Impact de la distance aux bois sur les méligrèthes et sur leurs parasoïdes

Cas d'étude statistique: approche paysage multi-bio-agresseurs

Application

FORMALISMES STATISTIQUES

- fréquentiste (Fisher) : “pas” d'*a priori* sur les données (illusoire)
- bayésien : on spécifie l'*a priori* sur les données



MODÉLISATION STATISTIQUE

Bases de stats

Corentin Barbu

Introduction

Modélisation statistique

Bases de régressions linéaires

Statistiques dans l'espace et auto-corrélation

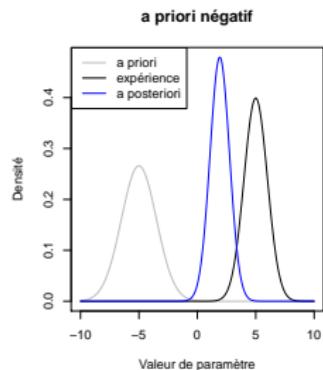
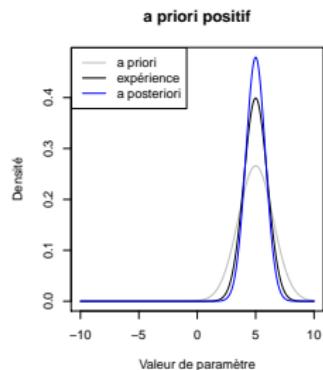
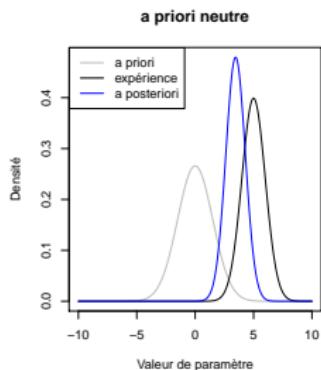
Cas d'étude statistique:
Impact de la distance aux bois sur les méligrèthes et sur leurs parasytoides

Cas d'étude statistique:
approche paysage multi-bio-agresseurs

Application

FORMALISMES STATISTIQUES

- fréquentiste (Fisher) : “pas” d'*a priori* sur les données (illusoire)
- bayésien : on spécifie l'*a priori* sur les données



MODÉLISATION STATISTIQUE

Bases de stats

Corentin Barbu

Introduction

Modélisation statistique

Bases de régressions linéaires

Statistiques dans l'espace et auto-corrélation

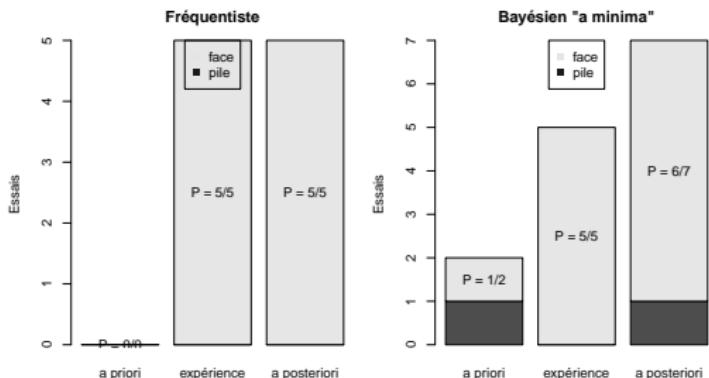
Cas d'étude statistique:
Impact de la distance aux bois sur les méligrèthes et sur leurs parasites

Cas d'étude statistique:
approche paysage multi-bio-agresseurs

Application

FORMALISMES STATISTIQUES

- ▶ fréquentiste (Fisher) : “pas” d'*a priori* sur les données (illusoire)
- ▶ bayésien : on spécifie l'*a priori* sur les données



Voir aussi “Hygiène Mentale”  (Ep 26 Bayesian Thinking) sur YouTube et séries de “science4all”.



Hygiène Mentale
340K subscribers

Introduction

Modélisation
statistique

Bases de régressions
linéaires

GLM et risque de la
sélection de variable

Le LASSO

Fonction de lien

Statistiques dans
l'espace et
auto-corrélation

Cas d'étude statistique:
Impact de la distance aux
bois sur les méligethes et
sur leurs parasitoïdes

Cas d'étude statistique:
approche paysage
multi-bio-agresseurs

Application

Subsection 1

BASES DE RÉGRESSIONS LINÉAIRES

Introduction

Modélisation
statistique

Bases de régressions
linéaires

GLM et risque de la
sélection de variable

Le LASSO

Fonction de lien

Statistiques dans
l'espace et
auto-corrélation

Cas d'étude statistique:
Impact de la distance aux
bois sur les méligethes et
sur leurs parasitoïdes

Cas d'étude statistique:
approche paysage
multi-bio-agresseurs

Application

Subsubsection 1

GLM et risque de la sélection de variable

LES BASES DE RÉGRESSION GLM (GENERALIZED LINEAR MODEL)

Bases de stats

Corentin Barbu

Introduction

Modélisation
statistique

Bases de régressions
linéaires

GLM et risque de la
sélection de variable

Le LASSO

Fonction de lien

Statistiques dans
l'espace et
auto-corrélation

Cas d'étude statistique:
Impact de la distance au
bois sur les méligethes
et sur leurs parasoïdes

Cas d'étude statistique:
approche paysage
multi-bio-agresseurs

Application

LA RÉGRESSION LINÉAIRE DE BASE

On considère que la valeur moyenne d'une variable de sortie y est fonction d'une combinaison linéaire de paramètres x :

$$\hat{y} = \alpha + \beta \cdot x_1 + \gamma \cdot x_2 + \dots$$

Souvent simplifié en utilisant les notations d'algèbres linéaires :

$$\hat{y} = A + BX$$

Autour de cette valeur moyenne on considère dans le cas le plus simple que les valeurs observées ont une distribution Gaussienne (loi normale):

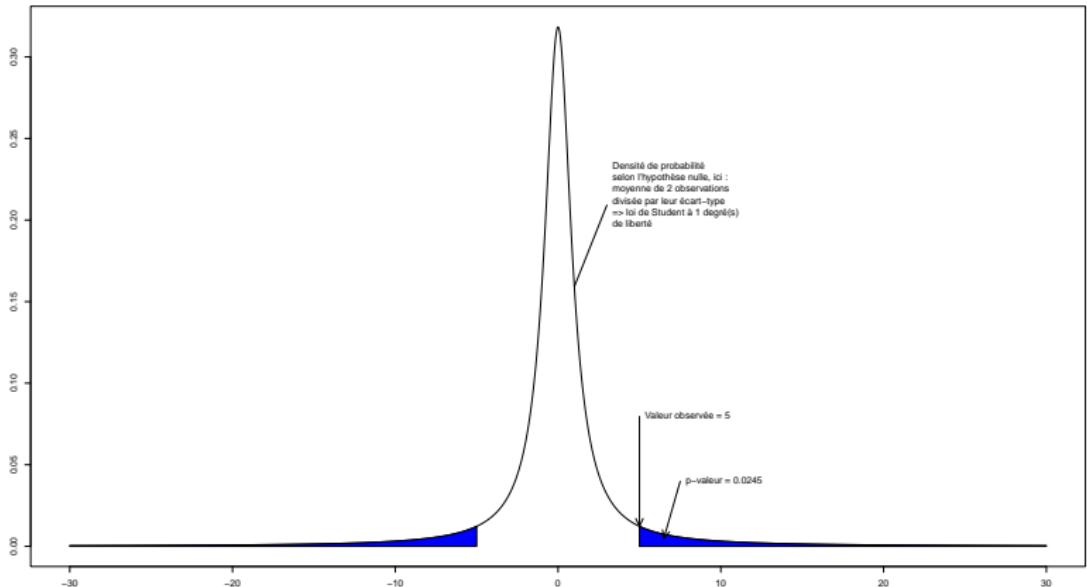
$$y_{obs} \sim N(\hat{y}, \sigma)$$

Le paramètre σ correspondant à l'écart-type des observations autour de la valeur prédite.

LA P-VALEUR (P-VALUE EN ANGLAIS)

DÉFINITION

Probabilité d'observer, par chance quelque chose d'aussi ou plus différent de 0 que ce qui est observé (voir aussi science4all "La plus grosse confusion des sciences, la p-value, Bayes 9").



Introduction

Modélisation statistique

Bases de régressions linéaires

GLM et risque de la sélection de variable

Le LASSO

Fonction de lien

Statistiques dans l'espace et auto-corrélation

Cas d'étude statistique: Impact de la distance au bois sur les méligethes et sur leurs parasitoïdes

Cas d'étude statistique: approche paysage multi-bio-agresseurs

Application

APPLICATIONS DU GLM

Bases de stats

Corentin Barbu

Introduction

Modélisation
statistique

Bases de régressions
linéaires

GLM et risque de la
sélection de variable

Le LASSO

Fonction de lien

Statistiques dans
l'espace et
auto-corrélation

Cas d'étude statistique:
Impact de la distance aux
bois sur les méligethes et
sur leurs parasoïdes

Cas d'étude statistique:
approche paysage
multi-bio-agresseurs

Application

- ▶ exemple dans R: effetSelection.R
- ▶ p-value
- ▶ AIC
- ▶ taille d'effets
- ⇒ Attention, la stratégie “mini-monde” (on met toutes les variables et on vera bien) risque de faire choisir des choses qui n'ont aucun sens, surtout si on se concentre sur les p-values

Bases de stats

Corentin Barbu

Introduction

Modélisation
statistique

Bases de régressions
linéaires

GLM et risque de la
sélection de variable

Le LASSO

Fonction de lien

Statistiques dans
l'espace et
auto-corrélation

Cas d'étude statistique:
Impact de la distance aux
bois sur les méligethes et
sur leurs parasitoïdes

Cas d'étude statistique:
approche paysage
multi-bio-agresseurs

Application

Subsubsection 2

Le LASSO

Introduction

Modélisation
statistique

Bases de régressions
linéaires

GLM et risque de la
sélection de variable

Le LASSO

Fonction de lien

Statistiques dans
l'espace et
auto-corrélation

Cas d'étude statistique:
Impact de la distance aux
bois sur les mélégithes et
sur leurs parasoïdes

Cas d'étude statistique:
approche paysage
multi-bio-agresseurs

Application

suite du fichier effetSelection.R

- ▶ on abandone la p-value
- ▶ on ne laisse s'exprimer que les variables très fortes: comme de la colle sur le zéro (lambda)
- ▶ validation croisée, choix du lambda 1se

Bases de stats

Corentin Barbu

Introduction

Modélisation
statistique

Bases de régressions
linéaires

GLM et risque de la
sélection de variable

Le LASSO

Fonction de lien

Statistiques dans
l'espace et
auto-corrélation

Cas d'étude statistique:
Impact de la distance aux
bois sur les méligethes et
sur leurs parasoïdes

Cas d'étude statistique:
approche paysage
multi-bio-agresseurs

Application

Subsubsection 3

Fonction de lien

GLM : ATTENTION À LA FORME DE LA VARIABLE EXPLIQUÉE

Bases de stats

Corentin Barbu

Introduction

Modélisation statistique

Bases de régressions linéaires

GLM et risque de la sélection de variable

Le LASSO

Fonction de lien

Statistiques dans l'espace et auto-corrélation

Cas d'étude statistique:
Impact de la distance aux bois sur les mélégithées et sur leurs parasites

Cas d'étude statistique:
approche paysage multi-bio-agresseurs

Application

LA FONCTION DE LIEN

Le paramètre “family” dans `glm()` de R :

- ▶ Si mesure continue, éventuellement négative : distribution normale ok (choix par défaut)
 - ▶ $y \sim N(\alpha + \beta x, \sigma)$
 - ▶ `glm(y ~ x, family = "normale")`
- ▶ Si mesure continue mais strictement positive : `glm` sur **log**
 - ▶ $\log(y) \sim N(\alpha + \beta x, \sigma)$
 - ▶ `logy <- log(y); glm(logy ~ x, family = "normale")`
- ▶ Si comptage, théoriquement loi de Poisson mais problématique
⇒ soit **binomiale négative** soit quasi-poisson dans R.
 - ▶ `glm.nb(y ~ x)`
 - ▶ Si plutôt “proportionalité” : `glm.nb(y ~ log(x))`
- ▶ Si binaire (présence/absence): loi de distribution **binomial**

Bases de stats

Corentin Barbu

Introduction

Modélisation
statistique

Bases de régressions
linéaires

Statistiques dans
l'espace et
auto-corrélation

Cas d'étude statistique:
Impact de la distance aux
bois sur les méligethes et
sur leurs parasoïdes

Cas d'étude statistique:
approche paysage
multi-bio-agresseurs

Application

Subsection 2

STATISTIQUES DANS L'ESPACE ET AUTO-CORRÉLATION

AUTO-CORRÉLATION : LE PROBLÈME

Bases de stats

Corentin Barbu

Introduction

Modélisation
statistique

Bases de régressions
linéaires

Statistiques dans
l'espace et
auto-corrélation

Cas d'étude statistique:
Impact de la distance aux
bois sur les méligrèthes et
sur leurs parasoïdes

Cas d'étude statistique:
approche paysage
multi-bio-agresseurs

Application

- ▶ Le problème: une source commune dont on se fiche est corrélée avec deux paramètres
- ▶ ex: position par rapport au flanc de coteaux, terre et exposition au vent
- ▶ si plein de points sur pleins de coteaux différents, il y aura des terres différentes et des expositions au vent différentes donc pas un problème
- ▶ Par contre, si plusieurs points sur le même coteaux, problème d'auto-corrélation spatiale : situation difficile à gérer statistiquement, à éviter dans le plan d'expérimentation.

Subsection 3

**CAS D'ÉTUDE STATISTIQUE: IMPACT DE LA
DISTANCE AUX BOIS SUR LES MÉLIGÈTHES ET SUR
LEURS PARASITOÏDES**

Introduction

Modélisation
statistique

Bases de régressions
linéaires

Statistiques dans
l'espace et
auto-corrélation

Cas d'étude statistique:
Impact de la distance aux
bois sur les méligèthes et
sur leurs parasitoïdes

Cas d'étude statistique:
approche paysage
multi-bio-agresseurs

Application

INTRODUCTION

Comment évolue l'abondance de méligèthe en fonction de la distance aux bois ?

- ▶ Exploration d'un système précis avec une forte connaissance à priori
- ▶ puis de voir où et comment taper
- ▶ description du dispositif

La simplification ici passe par un design de l'expérience pour simplifier les interactions

Subsection 4

CAS D'ÉTUDE STATISTIQUE: APPROCHE PAYSAGE MULTI-BIO-AGRESSEURS

Introduction

Modélisation
statistique

Bases de régressions
linéaires

Statistiques dans
l'espace et
auto-corrélation

Cas d'étude statistique:
Impact de la distance aux
bois sur les méligethes et
sur leurs parasoïdes

Cas d'étude statistique:
approche paysage
multi-bio-agresseurs

Introduction

Méthodes

Résultats

Discussion

Application

Subsubsection 1

Introduction

INTRODUCTION

Bases de stats

Corentin Barbu



- ▶ Limiter l'utilisation de phyto
 - Leviers paysagers ?
- ▶ Beaucoup d'études bio-agresseur par bio-agresseur
 - Multi-bio-agresseurs ?
- ▶ Parapluie phyto-sanitaire
 - C'est la réalité d'aujourd'hui, peut-on voir quelque chose malgré tout ?
- ⇒ Quantification des (dis)services des éléments paysagers

Introduction

Modélisation statistique

Bases de régressions linéaires

Statistiques dans l'espace et auto-corrélation

Cas d'étude statistique:
Impact de la distance aux bois sur les méligrèthes et sur leurs parasitoïdes

Cas d'étude statistique:
approche paysage multi-bio-agresseurs

Introduction

Méthodes

Résultats

Discussion

Application

Introduction

Modélisation
statistique

Bases de régressions
linéaires

Statistiques dans
l'espace et
auto-corrélation

Cas d'étude statistique:
Impact de la distance aux
bois sur les méligethes et
sur leurs parasoïdes

Cas d'étude statistique:
approche paysage
multi-bio-agresseurs

Introduction

Méthodes

Résultats

Discussion

Application

Subsubsection 2

Méthodes

Introduction

Modélisation
statistiqueBases de régressions
linéairesStatistiques dans
l'espace et
auto-corrélationCas d'étude statistique:
Impact de la distance aux
bois sur les méligrèthes et
sur leurs parasoïdesCas d'étude statistique:
approche paysage
multi-bio-agresseurs

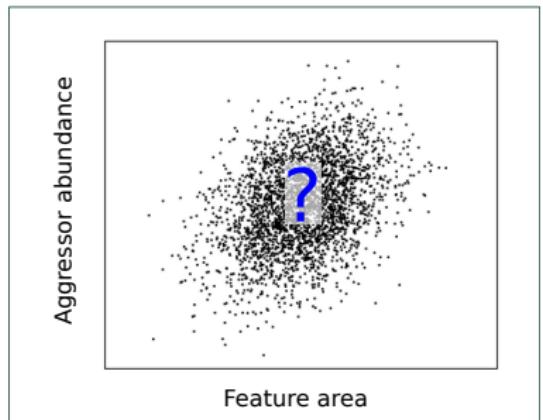
Introduction

Méthodes

Résultats

Discussion

Application



OBJECT

- ▶ France continentale
- ▶ Corrélations entre éléments paysagers and abondances des agresseurs
- ▶ quatre échelles: 200m, 1000m, 5km, 10km

APPROCHES

- ▶ Statistiques sur des données nationales
- ▶ Formalisation de la connaissance existente (élicitation)
 - ▶ Experts
 - ▶ Bibliographie

Introduction

Modélisation
statistique

Bases de régressions
linéaires

Statistiques dans
l'espace et
auto-corrélation

Cas d'étude statistique:
Impact de la distance aux
bois sur les méligrâthes et
sur leurs parasytoides

Cas d'étude statistique:
approche paysage
multi-bio-agresseurs

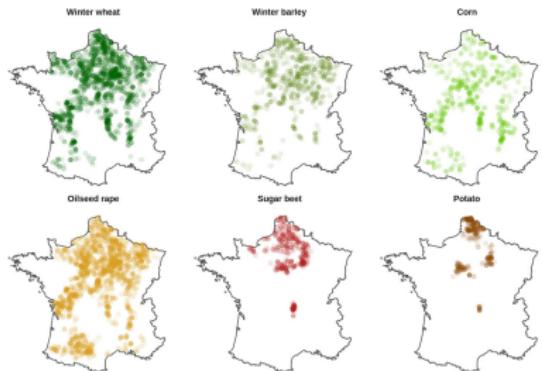
Introduction

Méthodes

Résultats

Discussion

Application



Points de relevés (2009-2017)

PAYSAGE

- ▶ RPG (PAC) → grandes cultures & prairies
- ▶ couche végétation BDTOPO (IGN) → bois, vergers, haies, landes

BIO-AGRESSEURS: VIGICULTURE®

- ▶ engagement annuel,
surveillance hebdomadaire
- ▶ vigilance sur tous les
bio-agresseurs → BSV
- ▶ souvent plusieurs mesures
par bio-agresseur
- ▶ géolocalisés à la parcelle

Introduction

Modélisation statistique

Bases de régressions linéaires

Statistiques dans l'espace et auto-corrélation

Cas d'étude statistique: Impact de la distance aux bois sur les méligrèthes et sur leurs parasytoides

Cas d'étude statistique: approche paysage multi-bio-agresseurs

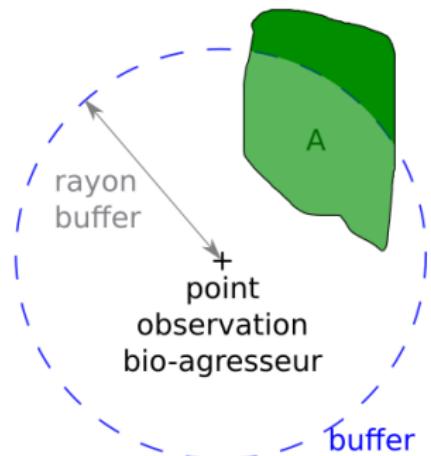
Introduction

Méthodes

Résultats

Discussion

Application



MÉTRIQUES PAR POINT ANNÉE

- ▶ abondance bio-agresseurs : $n_{obs} > seuil$
- ▶ métriques paysagères : aires dans 4 buffer

MODÈLES DE POISSON OU BINOMIAL NÉGATIF

$$\log(Y) = \alpha \cdot \text{région agro-climatique} + \sum_e \beta_e \cdot \log(Area_e)$$

- ▶ Univarié + régions
- ▶ Lasso ($\lambda = \lambda_{1se}$)

Introduction

Modélisation
statistique

Bases de régressions
linéaires

Statistiques dans
l'espace et
auto-corrélation

Cas d'étude statistique:
Impact de la distance aux
bois sur les méligethes et
sur leurs parasoïdes

Cas d'étude statistique:
approche paysage
multi-bio-agresseurs

Introduction

Méthodes

Résultats

Discussion

Application

Subsubsection 3

Résultats

Introduction

Modélisation
statistiqueBases de régressions
linéairesStatistiques dans
l'espace et
auto-corrélationCas d'étude statistique:
Impact de la distance aux
bois sur les méligrèthes et
sur leurs parasytoidesCas d'étude statistique:
approche paysage
multi-bio-agresseurs

Introduction

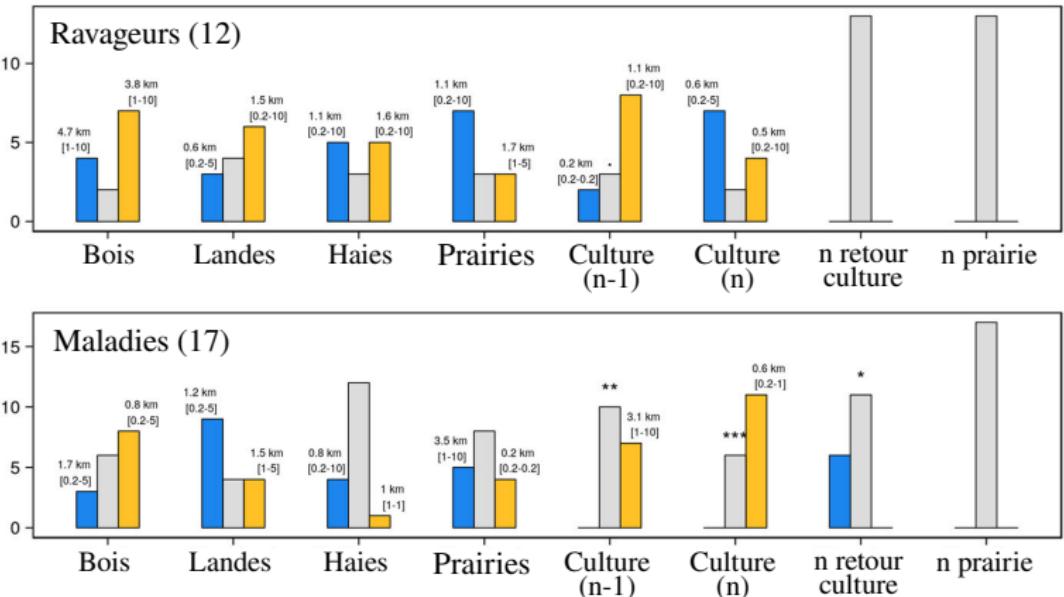
Méthodes

Résultats

Discussion

Application

STATISTIQUES PAYSAGÈRES



N bioagresseurs ayant une corrélation négative (bleu), absente (gris) ou positive (orange) avec chaque facteur paysager ou temporel

- ▶ statistiques: il y a de la puissance prédictive, accord pas fantastique mais incertitude expert énorme
- ▶ l'effet du paysage varie fortement entre bio-agresseurs
- ▶ En tout cas des cas avec service d'autre avec disservice

Bases de stats

Corentin Barbu

Introduction

Modélisation
statistique

Bases de régressions
linéaires

Statistiques dans
l'espace et
auto-corrélation

Cas d'étude statistique:
Impact de la distance aux
bois sur les méligethes et
sur leurs parasoïdes

Cas d'étude statistique:
approche paysage
multi-bio-agresseurs

Introduction

Méthodes

Résultats

Discussion

Application

Subsubsection 4

Discussion

CONCLUSIONS PRINCIPALES

Bases de stats

Corentin Barbu

Introduction

Modélisation
statistique

Bases de régressions
linéaires

Statistiques dans
l'espace et
auto-corrélation

Cas d'étude statistique:
Impact de la distance aux
bois sur les méligrèthes et
sur leurs parasoïdes

Cas d'étude statistique:
approche paysage
multi-bio-agresseurs

Introduction

Méthodes

Résultats

Discussion

Application

ECOSYSTEM (DIS)SERVICE DES ÉLÉMENTS PAYSAGERS

- ▶ Service/disservice dépend du bio-agresseurs
- ⇒ supposer par défaut un service de régulation écosystémique paraît déraisonnable

Bases de stats

Corentin Barbu

Introduction

Modélisation
statistique

Bases de régressions
linéaires

Statistiques dans
l'espace et
auto-corrélation

Cas d'étude statistique:
Impact de la distance aux
bois sur les méligethes et
sur leurs parasoïdes

Cas d'étude statistique:
approche paysage
multi-bio-agresseurs

Application

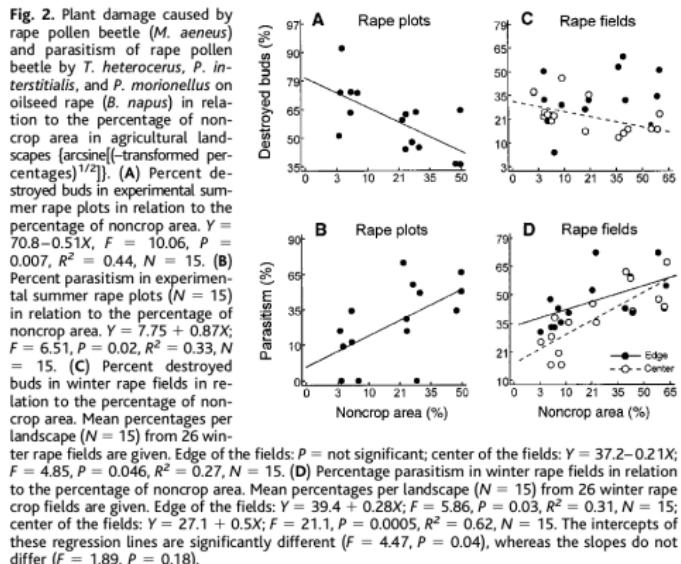
Subsection 5

APPLICATION

THIES ET TSCHARNTKE 1999, SCIENCE

Landscape structure and biological control in agroecosystems

RÉSULTATS



CONCLUSION:

"These results provide evidence that complex landscapes with a high density and connectivity of uncultivated, perennial habitats may enhance populations of natural enemies, which immigrate into neighboring annual crop fields, attack pest insects, and contribute significantly to the reduction of pest populations below an economic threshold." ?!?

Introduction

Modélisation statistique

Bases de régressions linéaires

Statistiques dans l'espace et auto-corrélation

Cas d'étude statistique: Impact de la distance aux bois sur les méligrâthes et sur leurs parasytoides

Cas d'étude statistique: approche paysage multi-bio-agresseurs

Application

[...] natural habitat provides important ecosystem services including pest control (...) Vraiment?

- ▶ Landis et al. 2000: Habitat management to conserve natural enemies of arthropod pests in agriculture
 - ▶ dizaines d'articles ↗ ennemis naturels.
 - ▶ ~ rien sur succès de contrôle
 - ▶ quelques exemples de favorisation des ravageurs...
- ▶ Bianchi et al. 2006: Sustainable pest regulation in agricultural landscapes: a review on landscape composition, biodiversity and natural pest control
 - ▶ 24 articles sur ennemis naturels. (↗ 18/ = 5/ ↘ 1)
 - ▶ 10 articles sur contrôle (↘ 4.5/=4/↗ 1.5)
- ▶ Karp et al. 2013: Forest bolsters bird abundance, pest control and coffee yield
 - ▶ oiseaux ↘ ravageurs (exclusion)
 - ▶ oiseaux mangent ravageurs et forêts ↗ ravageurs
 - ▶ forêts ↘ ravageurs : significatif à 125m (quelques buissons)
- ▶ Shackelford et al. 2013: Comparison of pollinators and natural enemies: A meta-analysis of landscape and local effects on abundance and richness in crops.
 - ▶ Rien sur les paysage/bioagresseurs, tout sur ennemis naturels
 - ▶ Effet espaces semi-naturels sur ennemis naturels < effet sur polliniseurs
- ▶ Milligan et al., 2016: Quantifying pest control services by birds and ants in Kenyan coffee farms
 - ▶ Prédation sur cartes décroît avec distance à la forêts (oiseaux) (12 points AIC)
 - ▶ prédit: décroît de 16% à 5% de prédation sur 100 m

Introduction

Modélisation statistique

Bases de régressions linéaires

Statistiques dans l'espace et auto-corrélation

Cas d'étude statistique:
Impact de la distance aux bois sur les méligrèthes et sur leurs parasytoides

Cas d'étude statistique:
approche paysage multi-bio-agresseurs

Application

When natural habitat fails to enhance biological pest control – Five hypotheses

ABSTRACT

Ecologists and farmers often have contrasting perceptions about the value of natural habitat in agricultural production landscapes, which so far has been little acknowledged in ecology and conservation. Ecologists and conservationists often appreciate the contribution of natural habitat to biodiversity and potential ecosystem services such as biological pest control, whereas many farmers see habitat remnants as a waste of cropland or source of pests. While natural habitat has been shown to increase pest control in many systems, we here identify five hypotheses for when and why natural habitat can fail to support biological pest control, and illustrate each with case studies from the literature: (1) pest populations have no effective natural enemies in the region, (2) natural habitat is a greater source of pests than natural enemies, (3) crops provide more resources for natural enemies than does natural habitat, (4) natural habitat is insufficient in amount, proximity, composition, or configuration to provide large enough enemy populations needed for pest control, and (5) agricultural practices counteract enemy establishment and biocontrol provided by natural habitat. In conclusion, we show that the relative importance of natural habitat for biocontrol can vary dramatically depending on type of crop, pest, predator, land management, and landscape structure. This variation needs to be considered when designing measures aimed at enhancing biocontrol services through restoring or maintaining natural habitat.

- ▶ (6) et si les ennemis naturels étaient tout simplement proportionnels aux ravageurs ?
- ▶ (7) et si le service écosystémique était l'exception plutôt que la règle ?

Attention à l'objectivité de ceux qui écrivent...

Introduction

Modélisation
statistique

Bases de régressions
linéaires

Statistiques dans
l'espace et
auto-corrélation

Cas d'étude statistique:
Impact de la distance aux
bois sur les méligrèthes et
sur leurs parasytoides

Cas d'étude statistique:
approche paysage
multi-bio-agresseurs

Application

Introduction

Modélisation
statistique

Bases de régressions
linéaires

Statistiques dans
l'espace et
auto-corrélation

Cas d'étude statistique:
Impact de la distance aux
bois sur les méligethes et
sur leurs parasoïdes

Cas d'étude statistique:
approche paysage
multi-bio-agresseurs

Application

► envoi des codes?