

Variations géographiques et sociodémographiques de la liaison en français

Une modélisation statistique

Basilio Calderone

basilio.calderone@univ-tlse2.fr

CNRS

&

CLLE-ERSS Université de Toulouse-Jean-Jaurès

- Objectif du travail
- La liaison - une petite présentation
- Un corpus phonologique : PFC
- Les outils statistiques (RL et GAM)
- Le model final et le résultats
- Limites du travail et perspectives futures

- Objectif du travail
- La liaison - une petite présentation
- Un corpus phonologique : PFC
- Les outils statistiques (RL et GAM)
- Le model final et le résultats
- Limites du travail et perspectives futures

1- Le present travail ne fait pas appelle aux théories/approches phonologiques sur la liaison en français ni sur des approches basées sur des contraintes syntaxiques (cadre HPSG et similaires) ou sur les aspects prosodiques

1- Le present travail ne fait pas appelle aux théories/approches phonologiques sur la liaison en français ni sur des approches basées sur des contraintes syntaxiques (cadre HPSG et similaires) ou sur les aspects prosodiques

2- Il est basé sur une analyse statistique des probabilités de réalisation (ou pas) des liaison dans un corpus d'observations qui espère être représentatif

- 1- Le present travail ne fait pas appelle aux théories/approches phonologiques sur la liaison en français ni sur des approches basées sur des contraintes syntaxiques (cadre HPSG et similaires) ou sur les aspects prosodiques
- 2- Il est basé sur une analyse statistique des probabilités de réalisation (ou pas) des liaison dans un corpus d'observations qui espère être représentatif
- 3-Il ne s'agit pas d'un modèle classifieur où les variables explicatives rentrent dans une boîte noire sans pouvoir les controller

1. Objectif

- Le present travail vise à réaliser un modèle prédictif de la liaison en français à partir d'un ensemble de variables sociolinguistiques, grammaticales et de surface

1. Objectif

- Le present travail vise à réaliser un modèle prédictif de la liaison en français à partir d'un ensemble de variables sociolinguistiques, grammaticales et de surface
- Dans notre approche on essaie de se baser sur les données observées pour synthétiser dans un modèle unique les variables explicatives qui sont significatives et participent à la réalisation (ou à la non réalisation) de la liaison

1. Objectif

- Le present travail vise à réaliser un modèle prédictif de la liaison en français à partir d'un ensemble de variables sociolinguistiques, grammaticales et de surface
- Dans notre approche on essaie de se baser sur les données observées pour synthétiser dans un modèle unique les variables explicatives qui sont significatives et participent à la réalisation (ou à la non réalisation) de la liaison
- Nous faisons des inférences statistiques : induire les caractéristiques inconnues d'une population (toutes les réalisations possibles de liaison) à partir d'un échantillon issu de cette population (notre jeu de données)

1. Objectif

- Le present travail vise à réaliser un modèle prédictif de la liaison en français à partir d'un ensemble de variables sociolinguistiques, grammaticales et de surface
- Dans notre approche on essaie de se baser sur les données observées pour synthétiser dans un modèle unique les variables explicatives qui sont significatives et participent à la réalisation (ou à la non réalisation) de la liaison
- Nous faisons des inférences statistiques : induire les caractéristiques inconnues d'une population (toutes les réalisations possibles de liaison) à partir d'un échantillon issu de cette population (notre jeu de données)
- Ce modèle nous permet d'identifier les variables qui ont un poids plus important parmi les variables prises en examen (choisies par moi même) et de les ordonner selon une échelle de “puissance explicative”

1. Objectif

- Le present travail vise à réaliser un modèle prédictif de la liaison en français à partir d'un ensemble de variables sociolinguistiques, grammaticales et de surface
- Dans notre approche on essaie de se baser sur les données observées pour synthétiser dans un modèle unique les variables explicatives qui sont significatives et participent à la réalisation (ou à la non réalisation) de la liaison
- Nous faisons des inférences statistiques : induire les caractéristiques inconnues d'une population (toutes les réalisations possibles de liaison) à partir d'un échantillon issu de cette population (notre jeu de données)
- Ce modèle nous permet d'identifier les variables qui ont un poids plus important parmi les variables prises en examen (choisies par moi même) et de les ordonner selon une échelle de “puissance explicative”
- Il s'agit d'un modèle prédictif, permettant de tirer des généralisations sur des données inconnues éventuelles

- Objectif du travail
- La liaison - une petite présentation
- Un corpus phonologique : PFC
- Les outils statistiques (RL et GAM)
- Le model final et le résultats
- Limites du travail et perspectives futures

2. La liaison en français

La liaison est un type de phénomène de *sandhi* qui consiste à prononcer la consonne finale d'un mot lorsque celui-ci précède un mot comportant une voyelle initiale.

Par exemple, entre le déterminant <les> /le/ et le substantif <amis> /ami/, tout locuteur natif insérera un phonème /z/ dit de liaison \longrightarrow \lɛz_a.mi\

2. La liaison en français

La liaison est un type de phénomène de *sandhi* qui consiste à prononcer la consonne finale d'un mot lorsque celui-ci précède un mot comportant une voyelle initiale.

Par exemple, entre le déterminant <les> /le/ et le substantif <amis> /ami/, tout locuteur natif insérera un phonème /z/ dit de liaison \longrightarrow \lɛz_a.mi\

il est

\i.l_ɛ\

2. La liaison en français

La liaison est un type de phénomène de *sandhi* qui consiste à prononcer la consonne finale d'un mot lorsque celui-ci précède un mot comportant une voyelle initiale.

Par exemple, entre le déterminant <les> /le/ et le substantif <amis> /ami/, tout locuteur natif insérera un phonème /z/ dit de liaison \longrightarrow \lɛz_a.mi\

il est	\i.l_ɛ\
nous avons	\nu.z_a.vɔ̃\

2. La liaison en français

La liaison est un type de phénomène de *sandhi* qui consiste à prononcer la consonne finale d'un mot lorsque celui-ci précède un mot comportant une voyelle initiale.

Par exemple, entre le déterminant <les> /le/ et le substantif <amis> /ami/, tout locuteur natif insérera un phonème /z/ dit de liaison \longrightarrow \lɛz_a.mi\

il est	\i.l_ɛ\
nous avons	\nu.z_a.vɑ̃\
les oiseaux	\le.z_wa.zo\

2. La liaison en français

La liaison est un type de phénomène de *sandhi* qui consiste à prononcer la consonne finale d'un mot lorsque celui-ci précède un mot comportant une voyelle initiale.

Par exemple, entre le déterminant <les> /le/ et le substantif <amis> /ami/, tout locuteur natif insérera un phonème /z/ dit de liaison \longrightarrow \lɛz_a.mi\

il est	\i.l_ɛ\
nous avons	\nu.z_a.vɑ̃\
les oiseaux	\le.z_wa.zo\
petit avion	\pe.ti.t_a.vjɑ̃\

2. La liaison en français

La liaison est un type de phénomène de *sandhi* qui consiste à prononcer la consonne finale d'un mot lorsque celui-ci précède un mot comportant une voyelle initiale.

Par exemple, entre le déterminant <les> /le/ et le substantif <amis> /ami/, tout locuteur natif insérera un phonème /z/ dit de liaison \longrightarrow \lɛz_a.mi\

il est	\i.l_ɛ\
nous avons	\nu.z_a.vɑ̃\
les oiseaux	\le.z_wa.zo\
petit avion	\pe.ti.t_a.vjɑ̃\
trop avancé	\tʁo.p_a.vɑ̃.se\

2. La liaison en français

La liaison est un type de phénomène de *sandhi* qui consiste à prononcer la consonne finale d'un mot lorsque celui-ci précède un mot comportant une voyelle initiale.

Par exemple, entre le déterminant <les> /le/ et le substantif <amis> /ami/, tout locuteur natif insérera un phonème /z/ dit de liaison \longrightarrow \lɛz_a.mi\

il est	\i.l_ɛ\
nous avons	\nu.z_a.vɑ̃\
les oiseaux	\le.z_wa.zo\
petit avion	\pe.ti.t_a.vjɑ̃\
trop avancé	\tʁo.p_a.vɑ̃.se\
pays imaginaires	\pe.i.z_i.ma.ʒi.nɛʁ\

2. La liaison en français

La liaison est un type de phénomène de *sandhi* qui consiste à prononcer la consonne finale d'un mot lorsque celui-ci précède un mot comportant une voyelle initiale.

Par exemple, entre le déterminant <les> /le/ et le substantif <amis> /ami/, tout locuteur natif insérera un phonème /z/ dit de liaison \longrightarrow \lɛz_a.mi\

il est	\i.l_ɛ\
nous avons	\nu.z_a.vɑ̃\
les oiseaux	\le.z_wa.zo\
petit avion	\pe.ti.t_a.vjɑ̃\
trop avancé	\tʁo.p_a.vɑ̃.se\
pays imaginaires	\pe.i.z_i.ma.ʒi.nɛʁ\
beaucoup à dire	\bo.ku.p_a diʁ\

2. La liaison en français

La liaison est un type de phénomène de *sandhi* qui consiste à prononcer la consonne finale d'un mot lorsque celui-ci précède un mot comportant une voyelle initiale.

Par exemple, entre le déterminant <les> /le/ et le substantif <amis> /ami/, tout locuteur natif insérera un phonème /z/ dit de liaison \longrightarrow \lɛz_a.mi\

il est	\i.l_ɛ\
nous avons	\nu.z_a.vɑ̃\
les oiseaux	\le.z_wa.zo\
petit avion	\pe.ti.t_a.vjɑ̃\
trop avancé	\tʁo.p_a.vɑ̃.se\
pays imaginaires	\pe.i.z_i.ma.ʒi.nɛʁ\
beaucoup à dire	\bo.ku.p_a diʁ\
pendant un	\pɑ̃.dɑ̃.t_ɑ̃\

2.1 La liaison en français

En littérature, on distingue a) la liaison obligatoire, b) la liaison facultative et c) la liaison interdite.

Liaison interdite

Liaison facultative

Liaison obligatoire



et encore

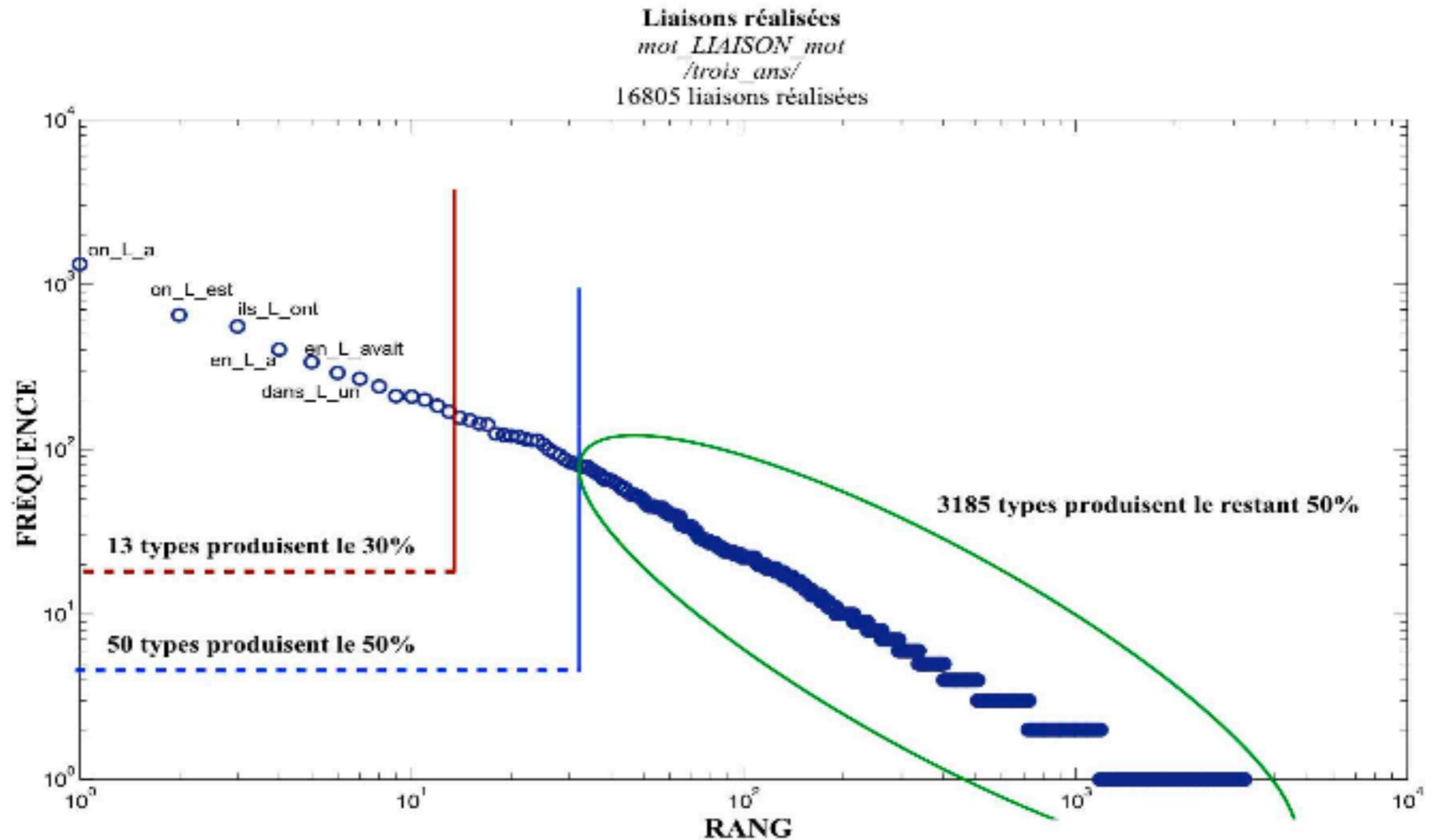
pas encore

les enfants

2.2 La liaison en français - analyses quantitatives

Analyse fréquentielle de la liaison en termes de *type* et *token* entre les deux mots (G et D)

De Jong 1996; Laks et al. 2015 et 2016



Analyse des correspondances (Benzécri, 1982) entre les POS des deux mots (G et D)
Laks et al. 2014

POS DROITE

POS GAUCHE

		ADJ	ADV	DET : ART	DET : POS	KON	NOM	INT	NUM	...	VER : ppres
ADJ		29	59	645	0	111	22	114	22	...	1
ADV		444	7	12	98	31	76	32	201	...	99
DET : ART		730	34	21	16	211	2	90	4	...	24
DET : POS		36	87	12	23	1321	44	333	38	...	65
KON		16	34	5	98	10	32	11	201	...	176
NOM		1544	13	12	91	110	78	12	81	...	92
INT		0	2	120	12	8	11	13	11	...	88
NUM		24	80	70	5	3	55	32	20	...	67
...	
VER : ppres		0	1	12	76	21	98	9	14	...	3

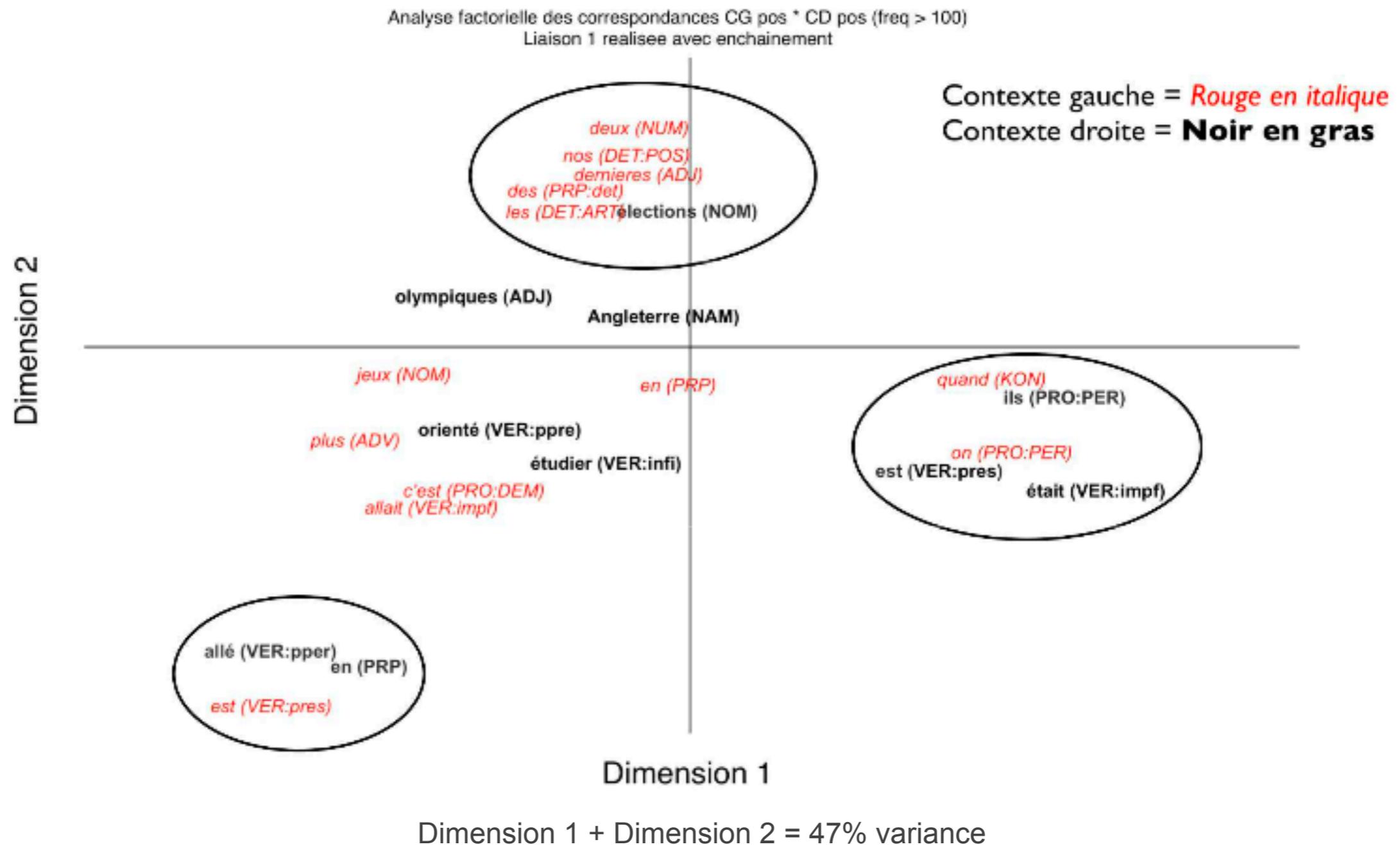
Analyse des correspondances (Benzécri, 1982) de la liaison en termes de *type* et *token* entre les deux mots (G et D)

Laks et al. 2014

2.2 La liaison en français - analyses quantitatives

Analyse des correspondances (Benzécri, 1982) de la liaison en termes de *type* et *token* entre les deux mots (G et D)

Laks et al. 2014



- Objectif du travail
- La liaison - une petite présentation
- Un corpus phonologique : PFC
- Les outils statistiques (RL et GAM)
- Le model final et le résultats
- Limites du travail et perspectives futures

3. La base de données PFC

<https://www.projet-pfc.net/>

PFC (Phonologie du Français Contemporain) est un programme de recherche offrant une base de données de français oral contemporain dans l'espace francophone



The screenshot shows the homepage of the PFC website. At the top, there is a navigation menu with the following items: ACCUEIL, PROTOCOLE, ENQUÊTES, CORPUS, REQUÊTES, OUTILS, IPFC, PFC-EF, COLLOQUES, PUBLICATIONS, ÉQUIPE, and HISTORIQUE. Below the menu is a search bar with the text "Search...". The main heading reads "PHONOLOGIE DU FRANÇAIS CONTEMPORAIN (PFC) : USAGES, VARIÉTÉS, STRUCTURE". Below this, a paragraph describes the project: "PFC (Phonologie du Français Contemporain) est un programme de recherche offrant une base de données de français oral contemporain dans l'espace francophone. Il s'adresse à un triple public :". A bulleted list follows, detailing the target audience: "chercheurs", "enseignants/apprenants de français", and "grand public". At the bottom, a paragraph states: "Le corpus PFC contient actuellement 16 enquêtes anonymisées (soit 164 locuteurs) publiées en ligne sur le site PFC et Ortolang. Plus de 40 autres enquêtes sont en cours de traitement. En 2008, PFC s'est élargi aux apprenants de français langue étrangère dans le sous-projet". On the right side, there are two sections: "PFC RECHERCHE" with links to "ACCÈS À LA BASE DE DONNÉES", "CORPUS PUBLIC", and "CORPUS RECHERCHE"; and "PFC ENSEIGNEMENT" with links to "LE PROJET PFC-EF", "PARTICIPANTS", "LE FRANÇAIS ILLUSTRÉ", "LE FRANÇAIS EXPLIQUÉ", and "RESSOURCES LINGUISTIQUES".

FLORAL

ACCUEIL PROTOCOLE ENQUÊTES CORPUS REQUÊTES OUTILS IPFC PFC-EF COLLOQUES
PUBLICATIONS ÉQUIPE HISTORIQUE

PHONOLOGIE DU FRANÇAIS CONTEMPORAIN (PFC) : USAGES, VARIÉTÉS, STRUCTURE

PFC (Phonologie du Français Contemporain) est un programme de recherche offrant une base de données de français oral contemporain dans l'espace francophone. Il s'adresse à un triple public :

- **chercheurs** : phonétique, phonologie, syntaxe, pragmatique, sociolinguistique, analyse conversationnelle, etc.
- **enseignants/apprenants de français** : langue étrangère, première ou seconde
- **grand public** : intéressé par les accents du français ou par le patrimoine linguistique francophone en général

Le corpus PFC contient actuellement 16 enquêtes anonymisées (soit 164 locuteurs) publiées en ligne sur le site PFC et Ortolang. Plus de 40 autres enquêtes sont en cours de traitement. En 2008, PFC s'est élargi aux apprenants de français langue étrangère dans le sous-projet

Search...

PFC RECHERCHE

- > ACCÈS À LA BASE DE DONNÉES
- > CORPUS PUBLIC
- > CORPUS RECHERCHE

PFC ENSEIGNEMENT

- > LE PROJET PFC-EF
- > PARTICIPANTS
- > LE FRANÇAIS ILLUSTRÉ
- > LE FRANÇAIS EXPLIQUÉ
- > RESSOURCES LINGUISTIQUES

3. La base de données PFC

La base de données

Durand, Laks & Lyche (2002)

The screenshot shows a web browser window with the URL 'projet-pfc.net'. The website has a dark navigation bar with the 'FLORAL' logo and menu items: ACCUEIL, PROTOCOLE, ENQUÊTES, CORPUS, REQUÊTES, OUTILS, IPFC, PFC-EF, COLLOQUES, PUBLICATIONS, ÉQUIPE, and HISTORIQUE. A search icon is on the right. The main content area is titled 'ACCÈS À LA BASE DE DONNÉES PFC' and contains the following text:

Vous pouvez consulter la base de données via le moteur de recherche ou le menu
Transcription gratuitement et sans enregistrement préalable

Pour avoir accès aux fonctions avancées de la base PFC, vous devez avoir un compte sur le site et vous identifier. Si vous n'avez pas de compte vous pouvez en créer un immédiatement.

Veillez remplir et envoyer la convention d'accès aux données :

Université Paris X Nanterre
Laboratoire MoDyCo
200 av de la République
92001 Nanterre
bât L R12C

CONVENTION D'ACCES AUX DONNEES SOURCES DU PROJET

PHONOLOGIE DU FRANÇAIS CONTEMPORAIN

Entre Jacques Durand, Bernard Laks et Chantal Lyche, responsables scientifiques du projet
Phonologie du Français Contemporain, ci après désignés comme la direction du projet PFC
Et
Nom, Prénom, Titre

On the right side, there is a search bar and two sections: 'PFC RECHERCHE' with links to 'ACCÈS À LA BASE DE DONNÉES', 'CORPUS PUBLIC', and 'CORPUS RECHERCHE'; and 'PFC ENSEIGNEMENT' with links to 'LE PROJET PFC-EF', 'PARTICIPANTS', 'LE FRANÇAIS ILLUSTRÉ', 'LE FRANÇAIS EXPLIQUÉ', 'RESSOURCES LINGUISTIQUES', 'RESSOURCES DIDACTIQUES', and 'DVD OPHRYS'. At the bottom right, there is a section for 'COLLOQUES PFC'.

L'interface

Recherche dans la base PFC - Liaisons PFC database liaison search

[\[Enquêtes\]](#) [\[Transcriptions\]](#) [\[Liaisons\]](#) [\[Schwas\]](#) [\[Logout\]](#)

[◀ Site PFC](#) • [PFC Site](#) [◀ Début](#) • [First](#) [? Aide](#) • [Help](#)

@liaison_consonne z

Recherche - Search

Recherche dans la transcription :

Tous les mots/All of these words | Un de ces mots/Any of these words | Cette phrase/This phrase)

Recherche dans tous les champs (transcription, enquêtes, régions, etc.)

Tous les mots/All of these words :

Un de ces mots/Any of these words :

(or at least of the words)

Cette phrase/This phrase :

Match whole field only

Aucun de ces mots/None of these words :

(must have at least one of the above)

Allow stemming (bridge matches bridging, bridges etc too)

Exact word matching

A gauche de la liaison :

All / Any / Phrase / Whole Field / Field ending with)

Mot de la liaison :

3.1 Resultats de la recherche: @consonne z @ Dijon @ Sexe M

VILLE

LIAISON RÉALISÉE

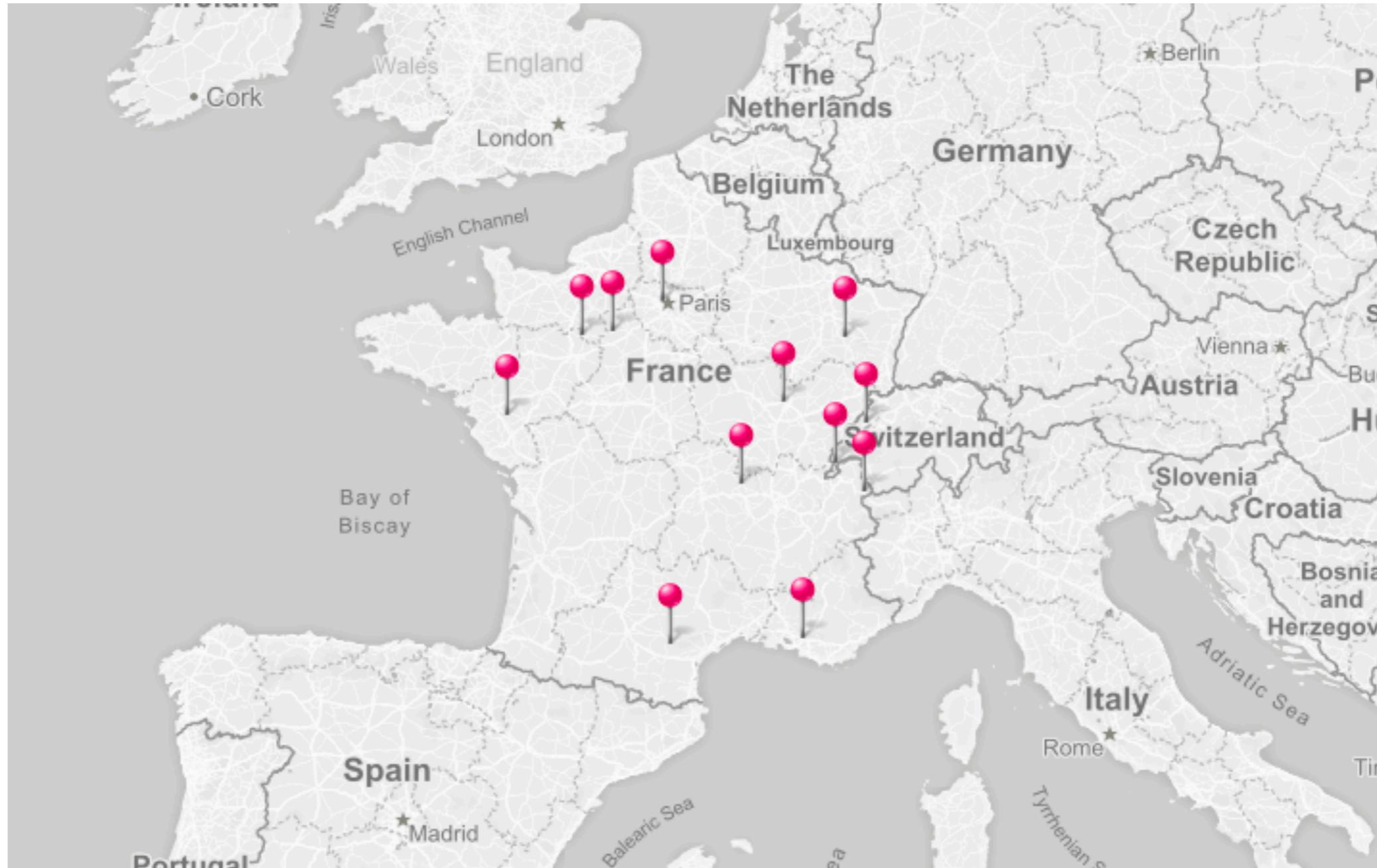
ID LOCUTEUR

La recherche pour @liaison_consonne z @loc_sexe m @enquete \"dijon\" a trouvé 85 réponses

Exporter toutes les réponses au format [CSV](#)

Info	Locuteur	Enquête	Contexte Gauche	Liaison	Ecoute	Contexte droit	L1	L2	L3	L4	L5
4035	21abl1	Dijon	quatrième aux jeux	jeux [z] olympiques		olympiques de Berlin en 1936 et plus récemment son usine de pâtes20 italiennes	1	z			
4042	21abl1	Dijon	Le maire de Beaulieu - Marc Blanc - est11t en revanche très	très [z] inquiet		inquiet	1	z			
4043	21abl1	Dijon	La cote du Premier Ministre ne cesse de baisser depuis les	les [z] élections		élections	1	z			
4046	21abl1	Dijon	La côte escarpée du Mont Saint-Pierre qui mène au village connaît des barrages chaque fois que les	les [z] opposants		opposants de tous les bords manifestent leur colère	1	z			
4053	21abl1	Dijon	ous	Nous [z] avons		avons le soutien du village entier	1	z			
4054	21abl1	Dijon	De plus quelques	quelques [z] articles		articles parus dans La Dépêche du Centre L'Express uest France et Le ouvel bserveateur	1	z			
4055	21abl1	Dijon	indiqueraient que des	des [z] activistes		activistes des communes voisines préparent20 une journée chaude au Premier Ministre	1	z			
9900	21abl1	Dijon	BL: Ben c'était euh c'était hard quoi En plus le système quand tu tu prends10 une nounou tu as des	des [z] aides		aides mais qui tombent10 au bout du de trois mois	1	z			
9904	21abl1	Dijon	BL: les	les [z] aides		aides10 au logement euh le RMI on touche à peu près cinq mille et quelques francs cinq mille trois cent francs	1	z			
9911	21abl1	Dijon	BL: Bon après c'est vrai que on a quand c'était vraiment très galère on j'ai pu on j'ai pu avoir des	des [z] aides		aides quoi de ma mère quoi	1	z			

3.2 Points d'enquêtes (limités à la France) de PFC



3.3 Exemples de données et dataset

3410	21aml1	Dijon	ML : Poser des plans tout le temps quoi tu vois une fois de temps	temps [z] en	en temps euh accord moi ça m'arrive aussi 1 z mais euh
22100	44ajs2	Nantes	JS: Euh ça aggrave pas mais elle a pas récupéré mais par contre elle est	est [] en	en forme 0

Détail dataset de PFC (après nettoyage semi-automatique)

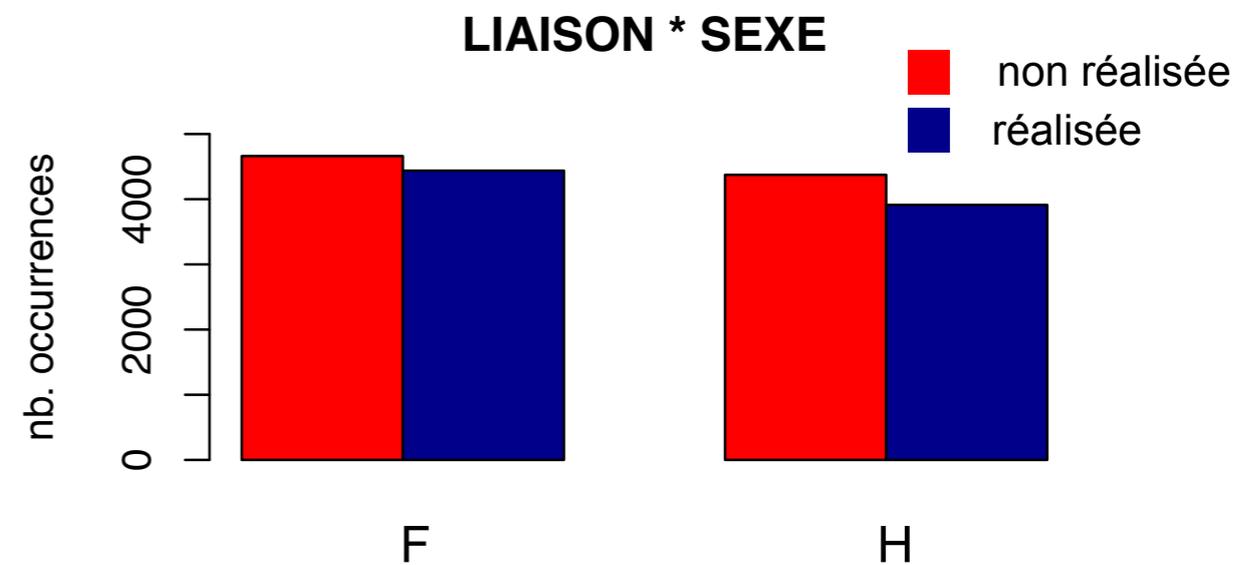
LIAISON	LOCUTEUR	ÉDUCATION	SEXE	AGE	VILLE	MOT_GAUCHE	MOT_DROITE	POS_GAUCHE	POS_DROITE
OUI	21aml1	20	F	65	Dijon	temps	en	NOM	PRP
NON	44ajs2	14	F	59	Nantes	est	en	VER	PRP
...
oui = 8350 no = 9036 TOT = 17386	niveaux = 192	ÉDUCATION = 15,41	F = 9100 H = 8286	AGE = 48,11	niveaux = 20	niveaux = 487	niveaux = 487	niveaux = 11	niveaux = 7

- Objectif du travail
- La liaison - une petite présentation
- Un corpus phonologique : PFC
- Les outils statistiques (RL et GAM)
- Le model final et le résultats
- Limites du travail et perspectives futures

4. Premières analyses statistiques

- Croiser deux variables catégorielles
- Tableaux de contingence LIAISON * SEXE

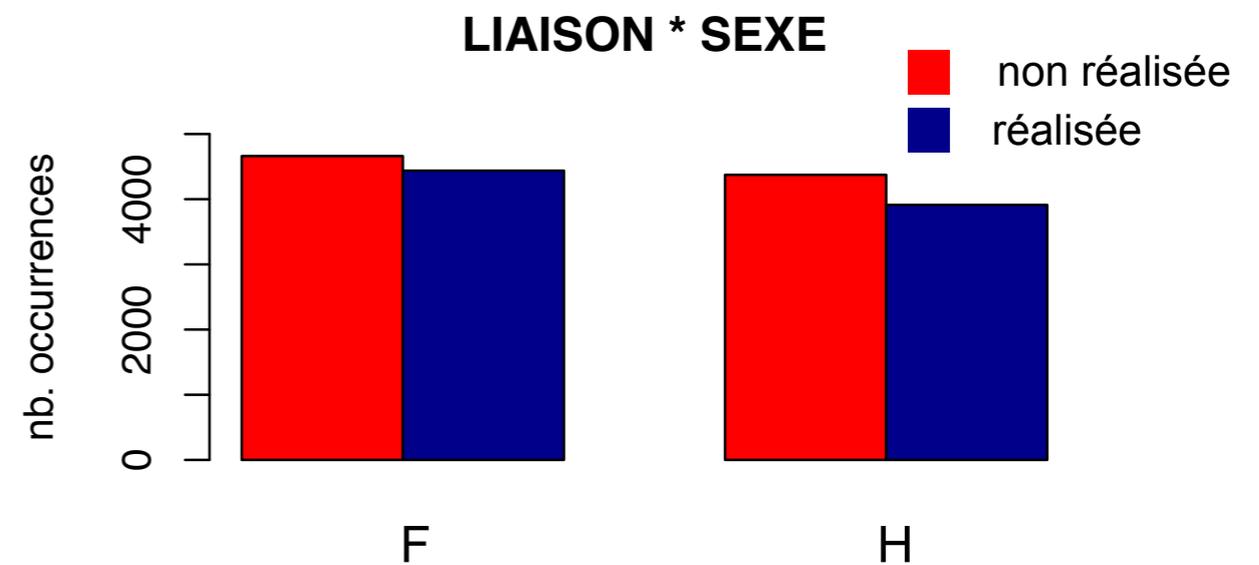
	F	H	Tot
Réalisée	4438	3912	8350
NON réalisée	4662	4374	9036
Tot	9100	8286	17386



4. Premières analyses statistiques

- Croiser deux variables catégorielles
- Tableaux de contingence LIAISON * SEXE

	F	H	Tot
Réalisée	4438	3912	8350
NON réalisée	4662	4374	9036
Tot	9100	8286	17386

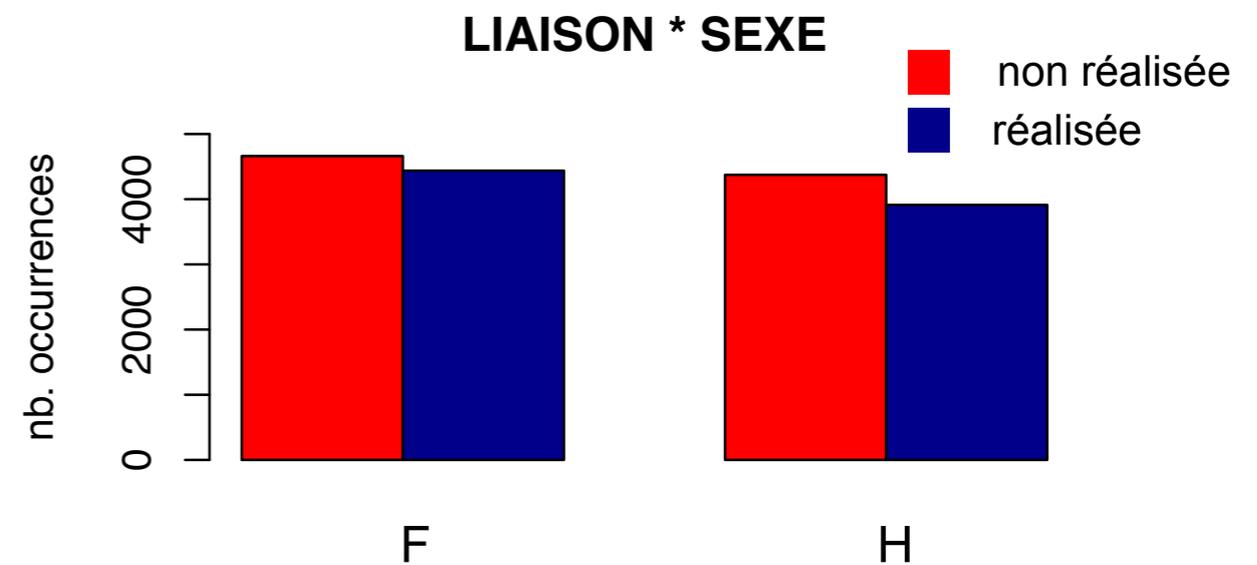


Pour savoir si le sexe a une influence significative sur la réalisation/non-réalisation de la liaison, nous pouvons effectuer un test d'indépendance du χ^2

4. Premières analyses statistiques

- Croiser deux variables catégorielles
- Tableaux de contingence LIAISON * SEXE

	F	H	Tot
Réalisée	4438	3912	8350
NON réalisée	4662	4374	9036
Tot	9100	8286	17386



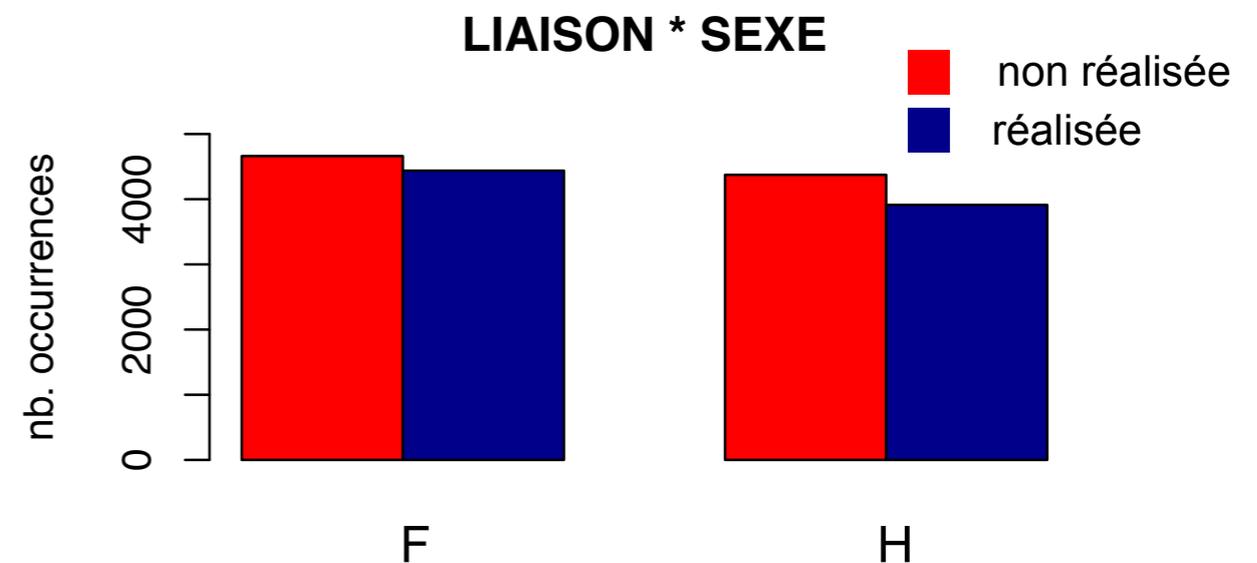
Pour savoir si le sexe a une influence significative sur la réalisation/non-réalisation de la liaison, nous pouvons effectuer un test d'indépendance du χ^2

Le test est significatif : $\chi^2 (1) = 4.1503, p\text{-value} < 0.05^*$

4. Premières analyses statistiques

- Croiser deux variables catégorielles
- Tableaux de contingence LIAISON * SEXE

	F	H	Tot
Réalisée	4438	3912	8350
NON réalisée	4662	4374	9036
Tot	9100	8286	17386



Pour savoir si le sexe a une influence significative sur la réalisation/non-réalisation de la liaison, nous pouvons effectuer un test d'indépendance du χ^2

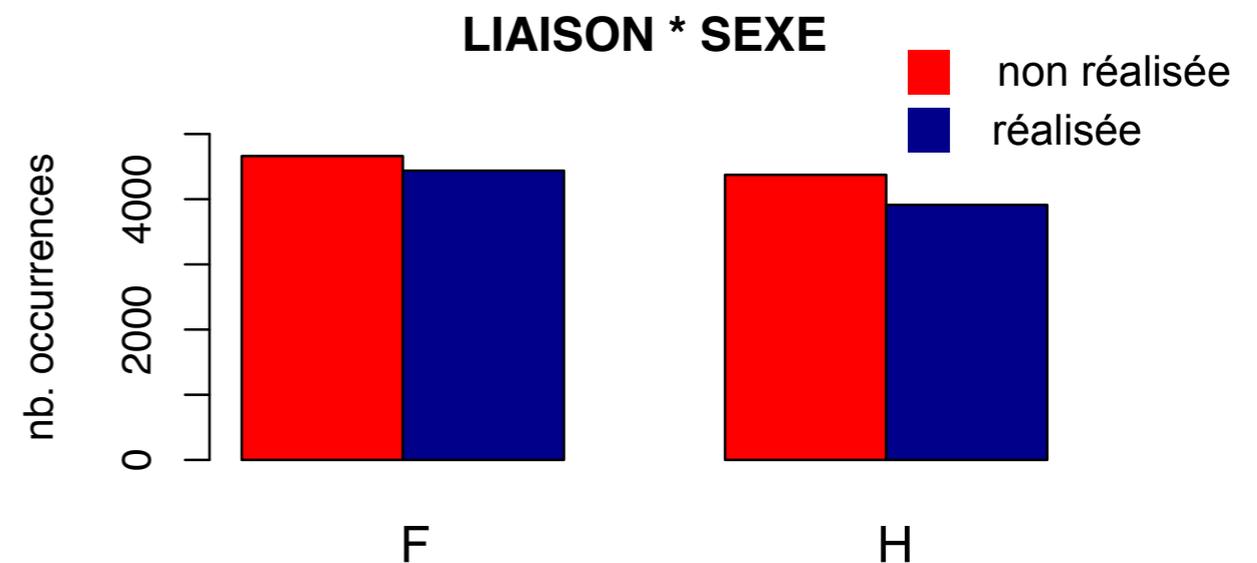
Le test est significatif : $\chi^2 (1) = 4.1503, p\text{-value} < 0.05^*$

Il existe une dépendance significative entre le nombre de liaisons et le sexe

4. Premières analyses statistiques

- Croiser deux variables catégorielles
- Tableaux de contingence LIAISON * SEXE

	F	H	Tot
Réalisée	4438	3912	8350
NON réalisée	4662	4374	9036
Tot	9100	8286	17386



Pour savoir si le sexe a une influence significative sur la réalisation/non-réalisation de la liaison, nous pouvons effectuer un test d'indépendance du χ^2

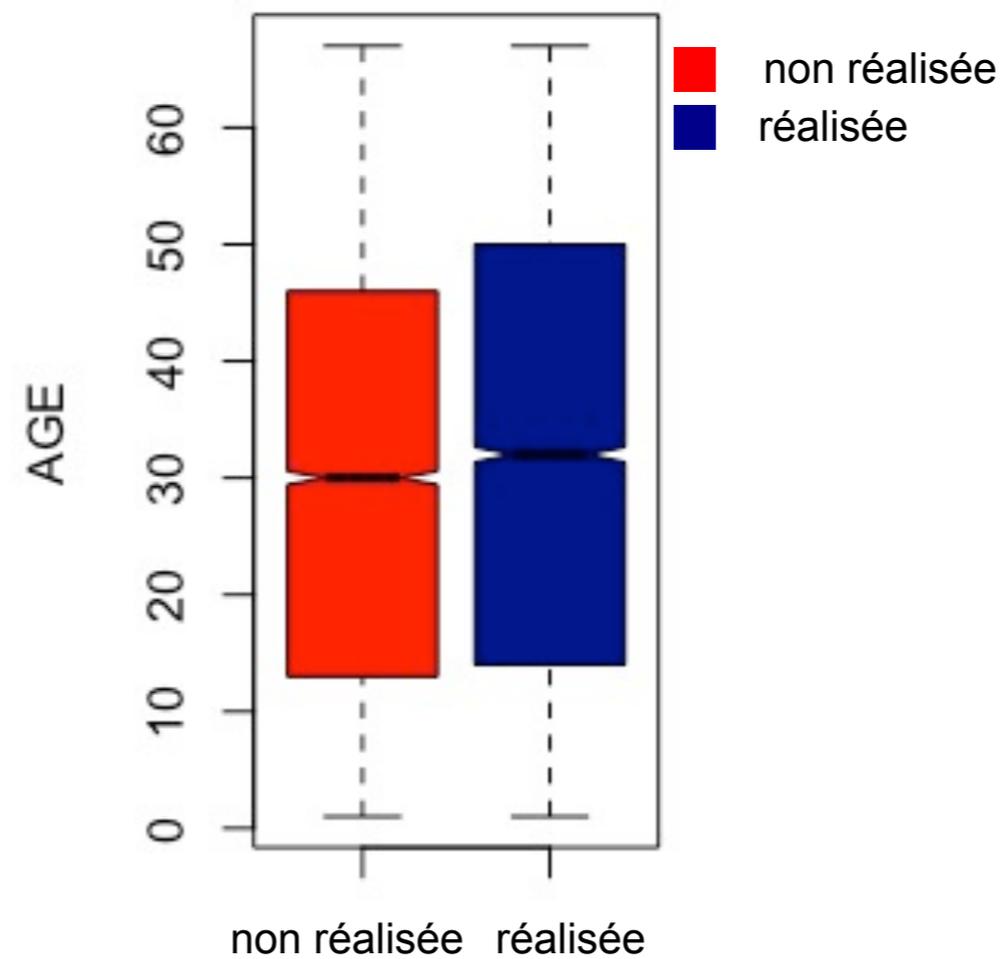
Le test est significatif : $\chi^2 (1) = 4.1503, p\text{-value} < 0.05^*$

Il existe une dépendance significative entre le nombre de liaisons et le sexe

C'est un test plutôt limitant et fortement lié aux données, qui ne permet pas une généralisation prédictive

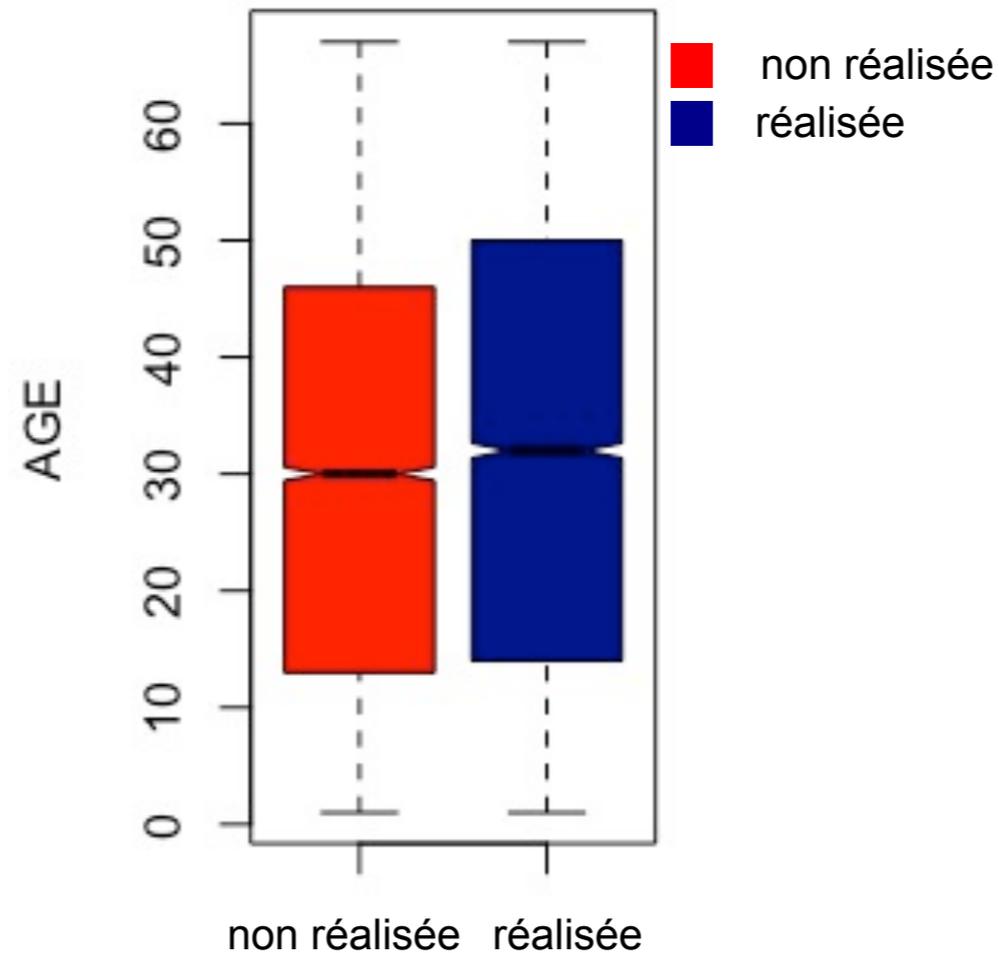
4. Premières analyses statistiques

- Croiser une variables catégorielle (LIAISON) avec une variable quantitative (AGE)



4. Premières analyses statistiques

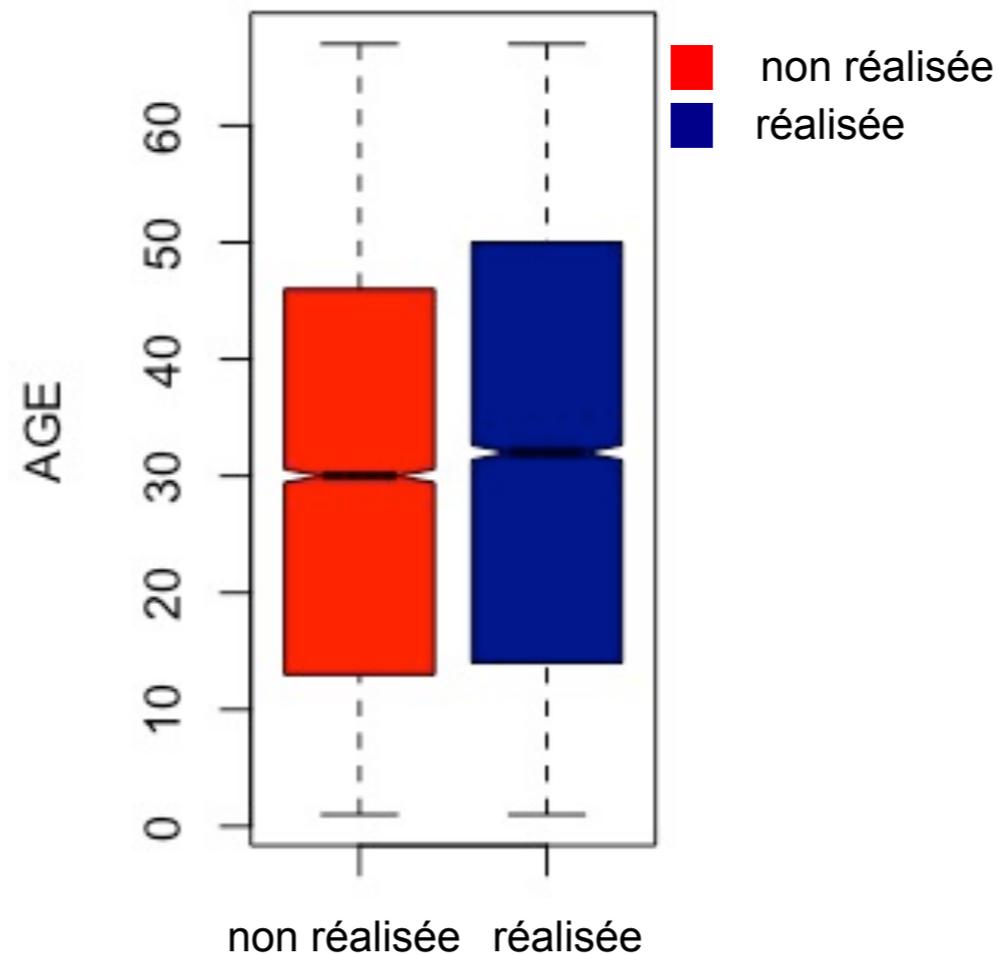
- Croiser une variables catégorielle (LIAISON) avec une variable quantitative (AGE)



Pour savoir si les deux groupes (réalisée vs. non réalisée) sont différents par rapport à l'âge des locuteurs, nous pouvons effectuer un test de variance (ANOVA)

4. Premières analyses statistiques

- Croiser une variables catégorielle (LIAISON) avec une variable quantitative (AGE)

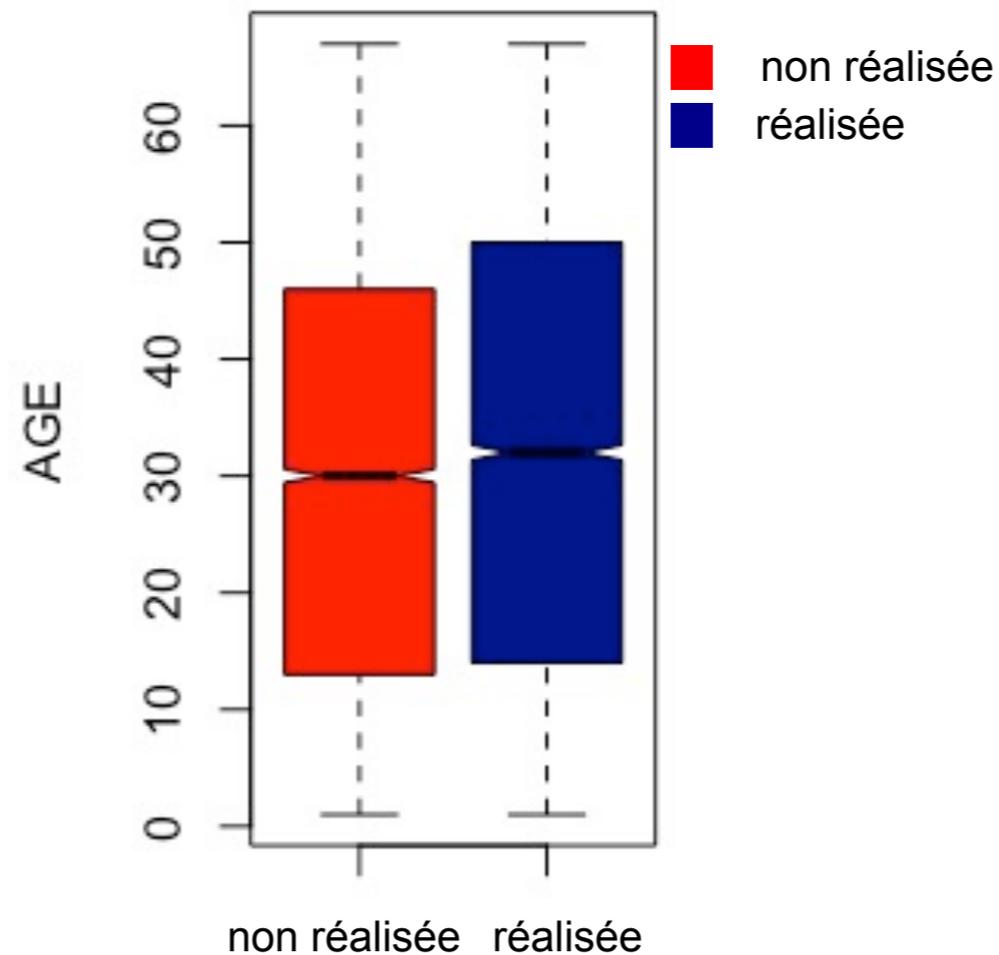


Pour savoir si les deux groupes (réalisée vs. non réalisée) sont différents par rapport à l'âge des locuteurs, nous pouvons effectuer un test de variance (ANOVA)

Le test est significatif : ANOVA (1,15863) = 37.084, $p\text{-value} < 0.001^{***}$

4. Premières analyses statistiques

- Croiser une variables catégorielle (LIAISON) avec une variable quantitative (AGE)



Pour savoir si les deux groupes (réalisée vs. non réalisée) sont différents par rapport à l'âge des locuteurs, nous pouvons effectuer un test de variance (ANOVA)

Le test est significatif : ANOVA (1,15863) = 37.084, $p\text{-value} < 0.001^{***}$

Les deux groupes sont statistiquement différents mais le test est limité à une statistique descriptive

Il nous faut une approche prédictive globale, capable de synthétiser dans un modèle unique toutes les variables à tester (qualitatives et quantitatives) et fournir des valeurs de significativité

Il nous faut une approche prédictive globale, capable de synthétiser dans un modèle unique toutes les variables à tester (qualitatives et quantitatives) et fournir des valeurs de significativité

Nécessité d'explicitier la nature de la relation entre variables (et pas seulement tester si deux variables sont différentes ou pas)

Il nous faut une approche prédictive globale, capable de synthétiser dans un modèle unique toutes les variables à tester (qualitatives et quantitatives) et fournir des valeurs de significativité

Nécessité d'explicitier la nature de la relation entre variables (et pas seulement tester si deux variables sont différentes ou pas)

Possibilité de faire de prédiction à partir des données analysées

4.1 La régression logistique (LR)

- La régression logistique (RL) permet de mesurer l'association entre la survenue (probabilité) d'un évènement (à expliquer) et les facteurs susceptibles de l'influencer (variables explicatives)

4.1 La régression logistique (LR)

- La régression logistique (RL) permet de mesurer l'association entre la survenue (probabilité) d'un évènement (à expliquer) et les facteurs susceptibles de l'influencer (variables explicatives)
- Le but de la RL est de caractériser les relations entre la variable dépendante \mathbf{Y} (ou variable à expliquer) et plusieurs variables \mathbf{X}_n prises en compte simultanément (variables explicatives)

4.1 La régression logistique (LR)

- La régression logistique (RL) permet de mesurer l'association entre la survenue (probabilité) d'un évènement (à expliquer) et les facteurs susceptibles de l'influencer (variables explicatives)
- Le but de la RL est de caractériser les relations entre la variable dépendante \mathbf{Y} (ou variable à expliquer) et plusieurs variables \mathbf{X}_n prises en compte simultanément (variables explicatives)
- La RL s'applique lorsque la variable à expliquer (\mathbf{Y}) est qualitative (variable catégorielle avec nb. factors ≥ 2)

4.1 La régression logistique (LR)

- La régression logistique (RL) permet de mesurer l'association entre la survenue (probabilité) d'un évènement (à expliquer) et les facteurs susceptibles de l'influencer (variables explicatives)
- Le but de la RL est de caractériser les relations entre la variable dépendante \mathbf{Y} (ou variable à expliquer) et plusieurs variables \mathbf{X}_n prises en compte simultanément (variables explicatives)
- La RL s'applique lorsque la variable à expliquer (\mathbf{Y}) est qualitative (variable catégorielle avec nb. factors ≥ 2)

$$PLUIE_{OUI/NON} = NEBULOSITE$$

4.1 La régression logistique (LR)

- La régression logistique (RL) permet de mesurer l'association entre la survenue (probabilité) d'un évènement (à expliquer) et les facteurs susceptibles de l'influencer (variables explicatives)
- Le but de la RL est de caractériser les relations entre la variable dépendante \mathbf{Y} (ou variable à expliquer) et plusieurs variables \mathbf{X}_n prises en compte simultanément (variables explicatives)
- La RL s'applique lorsque la variable à expliquer (Y) est qualitative (variable catégorielle avec nb. factors ≥ 2)

$$PLUIE_{OUI/NON} = NEBULOSITE$$

$$REUSSIR_{OUI/NON} = TEMPS - TRAVAIL$$

4.1 La régression logistique (LR)

- La régression logistique (RL) permet de mesurer l'association entre la survenue (probabilité) d'un évènement (à expliquer) et les facteurs susceptibles de l'influencer (variables explicatives)
- Le but de la RL est de caractériser les relations entre la variable dépendante Y (ou variable à expliquer) et plusieurs variables X_n prises en compte simultanément (variables explicatives)
- La RL s'applique lorsque la variable à expliquer (Y) est qualitative (variable catégorielle avec nb. factors ≥ 2)

$$PLUIE_{OUI/NON} = NEBULOSITE$$

$$REUSSIR_{OUI/NON} = TEMPS - TRAVAIL$$

$$LIAISON_{OUI/NON} = variables X_n$$

$$Y = \beta_0 + \beta_1 \alpha_1 + \epsilon$$

$$Y = \beta_0 + \beta_1 \alpha_1 + \epsilon$$

La variable à expliquer (Y) est exprimée en fonction d'un intercept (ou ordonnée à l'origine) β_0 , d'une variable explicative (α_1) rattachée à son coefficient β_1 et d'un terme de bruit ϵ

$$Y = \beta_0 + \beta_1 \alpha_1 + \epsilon$$

La variable à expliquer (Y) est exprimée en fonction d'un intercept (ou ordonnée à l'origine) β_0 , d'une variable explicative (α_1) rattachée à son coefficient β_1 et d'un terme de bruit ϵ

Y est la variable catégorielle à expliquer : $PLUIE_{OUI/NON}$, $REUSSIR_{OUI/NON}$

4.2 La régression logistique - définition mathématique

$$Y = \beta_0 + \beta_1 \alpha_1 + \epsilon$$

La variable à expliquer (Y) est exprimée en fonction d'un intercept (ou ordonnée à l'origine) β_0 , d'une variable explicative (α_1) rattachée à son coefficient β_1 et d'un terme de bruit ϵ

Y est la variable catégorielle à expliquer : $PLUIE_{OUI/NON}$, $REUSSIR_{OUI/NON}$

Le paramètre α_1 représente la variable explicative du modèle:
NÉBULOSITÉ, TEMPS_DE_TRAVAIL

4.2 La régression logistique - définition mathématique

$$Y = \beta_0 + \beta_1 \alpha_1 + \epsilon$$

La variable à expliquer (Y) est exprimée en fonction d'un intercept (ou ordonnée à l'origine) β_0 , d'une variable explicative (α_1) rattachée à son coefficient β_1 et d'un terme de bruit ϵ

Y est la variable catégorielle à expliquer : $PLUIE_{OUI/NON}$, $REUSSIR_{OUI/NON}$

Le paramètre α_1 représente la variable explicative du modèle:
NÉBULOSITÉ, TEMPS_DE_TRAVAIL

Le coefficient β_1 mesure l'association entre la variable explicative (α_1) et la variable expliquée (Y). Il correspond au logarithme de l'odds ratio
 $\beta = \log(\text{OR})$

4.2 La régression logistique - définition mathématique

$$Y = \beta_0 + \beta_1 \alpha_1 + \epsilon$$

La variable à expliquer (Y) est exprimée en fonction d'un intercept (ou ordonnée à l'origine) β_0 , d'une variable explicative (α_1) rattachée à son coefficient β_1 et d'un terme de bruit ϵ

Y est la variable catégorielle à expliquer : $PLUIE_{OUI/NON}$, $REUSSIR_{OUI/NON}$

Le paramètre α_1 représente la variable explicative du modèle:
NÉBULOSITÉ, TEMPS_DE_TRAVAIL

Le coefficient β_1 mesure l'association entre la variable explicative (α_1) et la variable expliquée (Y). Il correspond au logarithme de l'odds ratio
 $\beta = \log(\text{OR})$

L'intercept β_0 correspond aux valeurs baseline du modèle lorsque les variables explicatives sont égales à 0 ($\beta_n=0$).

4.2 La régression logistique - définition mathématique

$$Y = \beta_0 + \beta_1 \alpha_1 + \epsilon$$

La variable à expliquer (Y) est exprimée en fonction d'un intercept (ou ordonnée à l'origine) β_0 , d'une variable explicative (α_1) rattachée à son coefficient β_1 et d'un terme de bruit ϵ

Y est la variable catégorielle à expliquer : $PLUIE_{OUI/NON}$, $REUSSIR_{OUI/NON}$

Le paramètre α_1 représente la variable explicative du modèle:
NÉBULOSITÉ, TEMPS_DE_TRAVAIL

Le coefficient β_1 mesure l'association entre la variable explicative (α_1) et la variable expliquée (Y). Il correspond au logarithme de l'odds ratio
 $\beta = \log(\text{OR})$

L'intercept β_0 correspond aux valeurs baseline du modèle lorsque les variables explicatives sont égales à 0 ($\beta_n=0$).

ϵ est un terme d'erreur non contrôlé qui doit impérativement suivre une distribution normale et de même variance. C'est un paramètre négligeable

4.2 La régression logistique - définition mathématique

Dans le cas de nos données sur la liaison, on essaie de trouver une relation entre l'AGE du locuteur et la réalisation de la liaison

$$Liaison = \beta_0 + \beta_1 AGE + \epsilon$$

4.2 La régression logistique - définition mathématique

Dans le cas de nos données sur la liaison, on essaie de trouver une relation entre l'AGE du locuteur et la réalisation de la liaison

$$Liaison = \beta_0 + \beta_1 AGE + \epsilon$$

$$LOGIT[P(Liaison = oui) | AGE] = \beta_0 + \beta_1 AGE + \epsilon$$

4.2 La régression logistique - définition mathématique

Dans le cas de nos données sur la liaison, on essaie de trouver une relation entre l'AGE du locuteur et la réalisation de la liaison

$$Liaison = \beta_0 + \beta_1 AGE + \epsilon$$

$$LOGIT[P(Liaison = oui) | AGE] = \beta_0 + \beta_1 AGE + \epsilon$$

Trouver les coefficients:

β_0 (Intercept)	-0.2713395
β_1	0.0043476

4.2 La régression logistique - définition mathématique

Dans le cas de nos données sur la liaison, on essaie de trouver une relation entre l'AGE du locuteur et la réalisation de la liaison

$$Liaison = \beta_0 + \beta_1 AGE + \epsilon$$

$$LOGIT[P(Liaison = oui) | AGE] = \beta_0 + \beta_1 AGE + \epsilon$$

Trouver les coefficients:

β_0 (Intercept)	-0.2713395
β_1	0.0043476

β_1 = odds ratio de la probabilité d'avoir *Liaison=oui* quand AGE augmente d'une unité

4.2 La régression logistique - définition mathématique

Dans le cas de nos données sur la liaison, on essaie de trouver une relation entre l'AGE du locuteur et la réalisation de la liaison

$$Liaison = \beta_0 + \beta_1 AGE + \epsilon$$

$$LOGIT[P(Liaison = oui) | AGE] = \beta_0 + \beta_1 AGE + \epsilon$$

Trouver les coefficients:

β_0 (Intercept)	-0.2713395
β_1	0.0043476

β_1 = odds ratio de la probabilité d'avoir *Liaison=oui* quand AGE augmente d'une unité

Regarder le signe : POSITIF

4.2 La régression logistique - définition mathématique

Dans le cas de nos données sur la liaison, on essaie de trouver une relation entre l'AGE du locuteur et la réalisation de la liaison

$$Liaison = \beta_0 + \beta_1 AGE + \epsilon$$

$$LOGIT[P(Liaison = oui) | AGE] = \beta_0 + \beta_1 AGE + \epsilon$$

Trouver les coefficients:

$$\begin{array}{ll} \beta_0(\text{Intercept}) & -0.2713395 \\ \beta_1 & 0.0043476 \end{array}$$

β_1 = odds ratio de la probabilité d'avoir *Liaison=oui* quand AGE augmente d'une unité

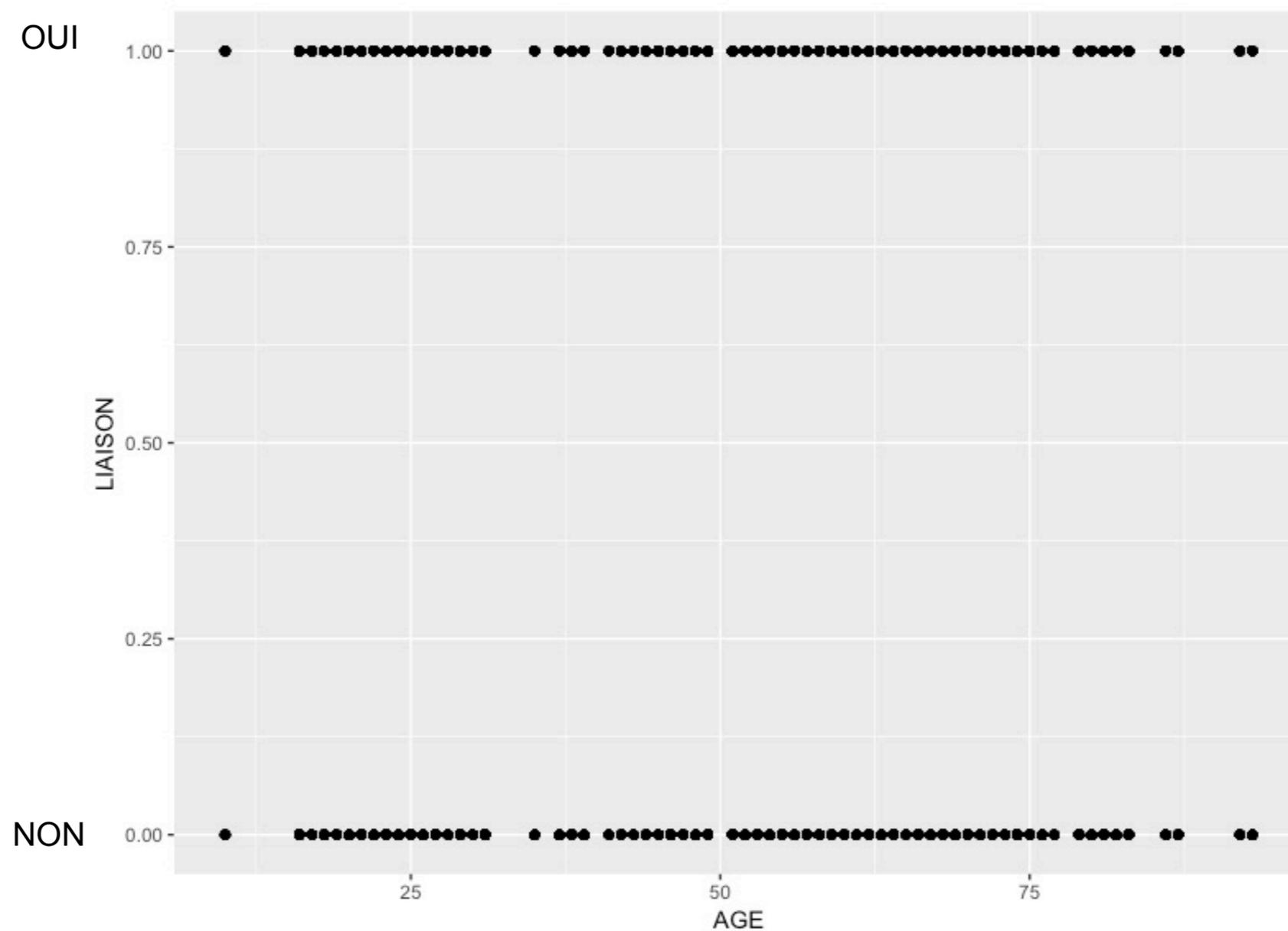
Regarder le signe : POSITIF

$$LOGIT[P(Liaison = oui) | AGE] = -0.2713395 + 0.0043476 AGE + \epsilon$$

A partir de cette équation on peut tracer la fonction de régression $Liaison \sim AGE$

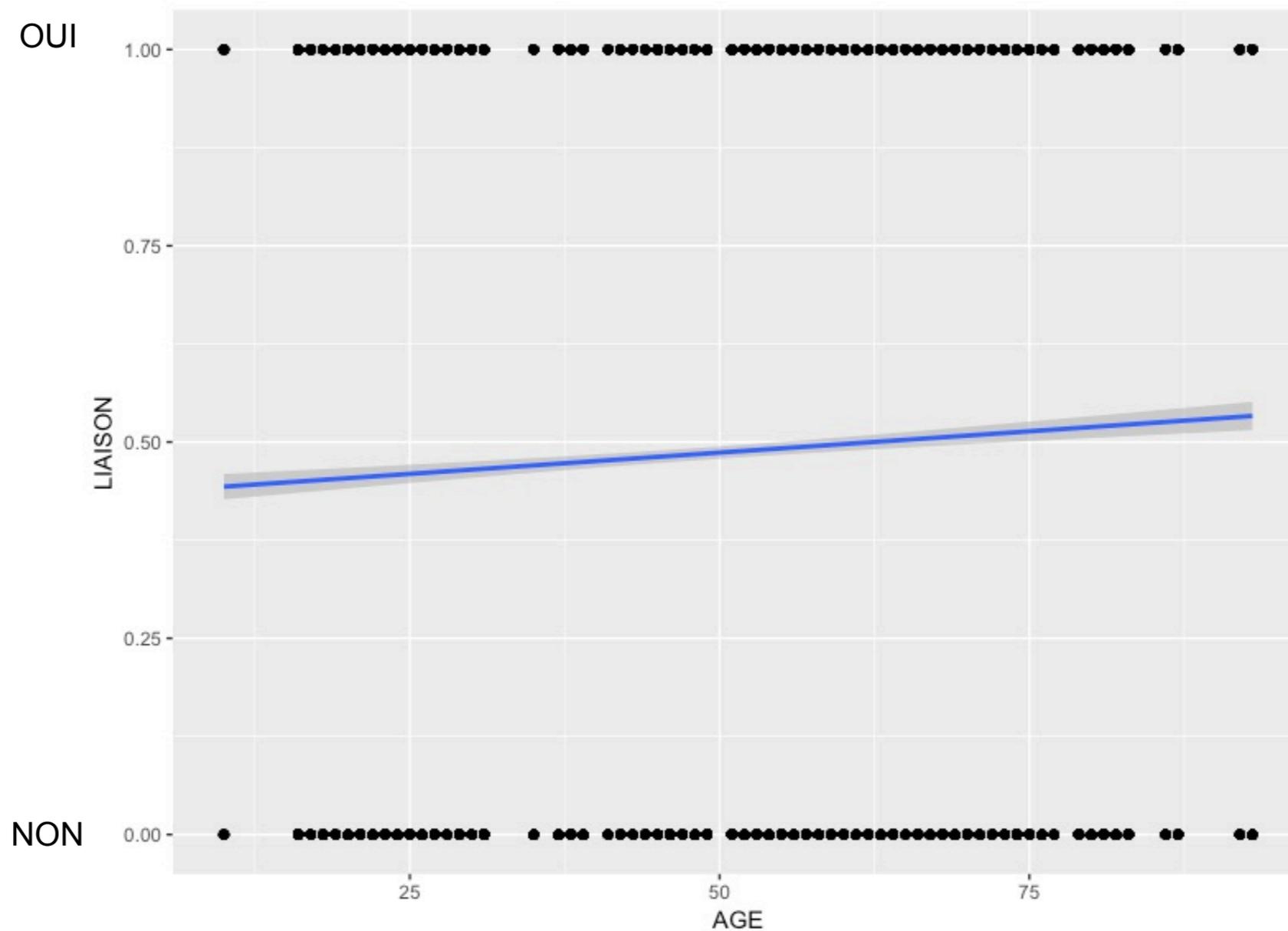
4.3 Un premier modèle: Liaison ~ AGE

$$\text{LOGIT}[P(\text{Liaison} = \text{oui} \mid \text{AGE})] = -0.2713395 + 0.0043476 \text{ AGE} + \epsilon$$

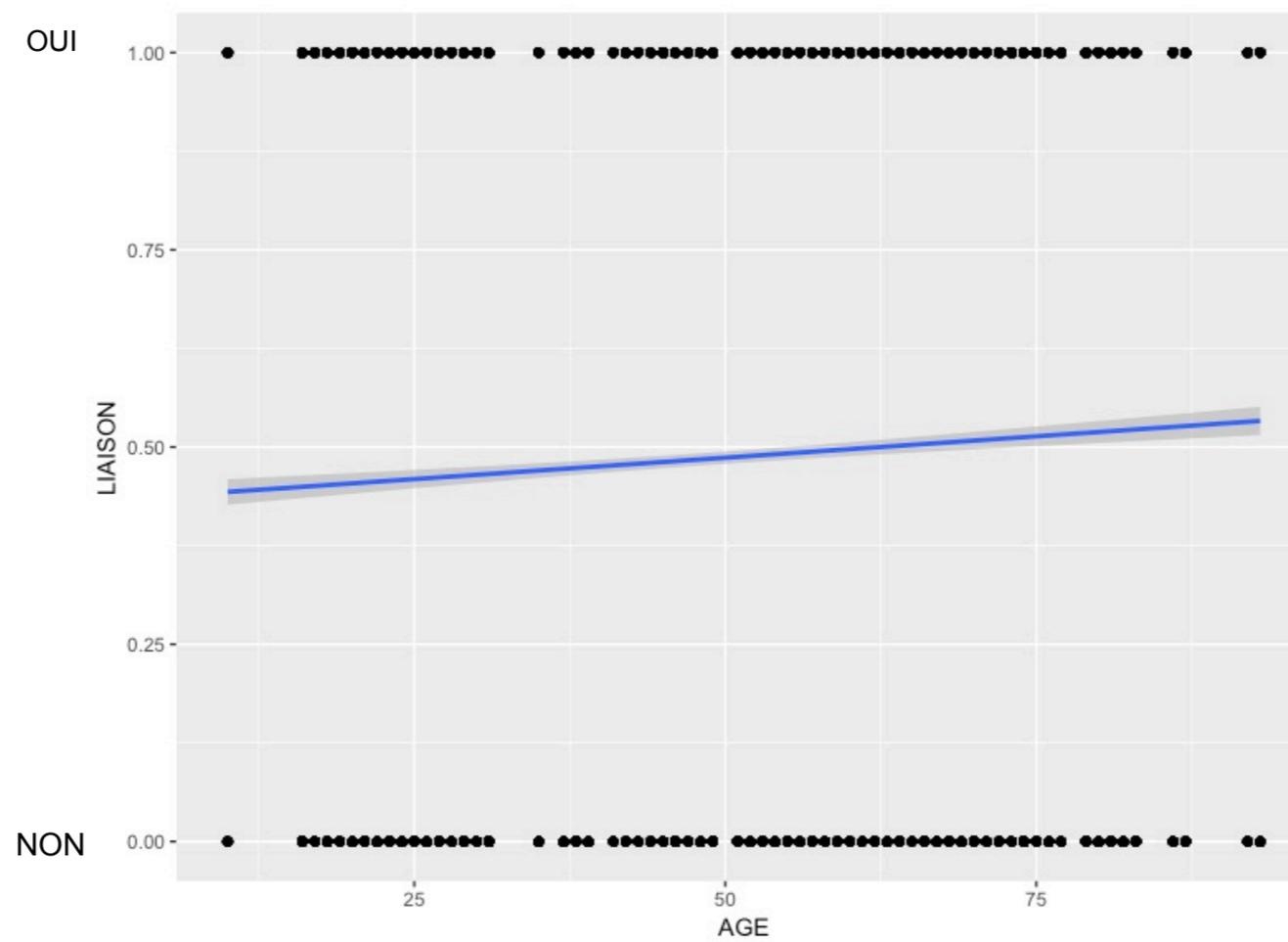


4.3 Un premier modèle: Liaison ~ AGE

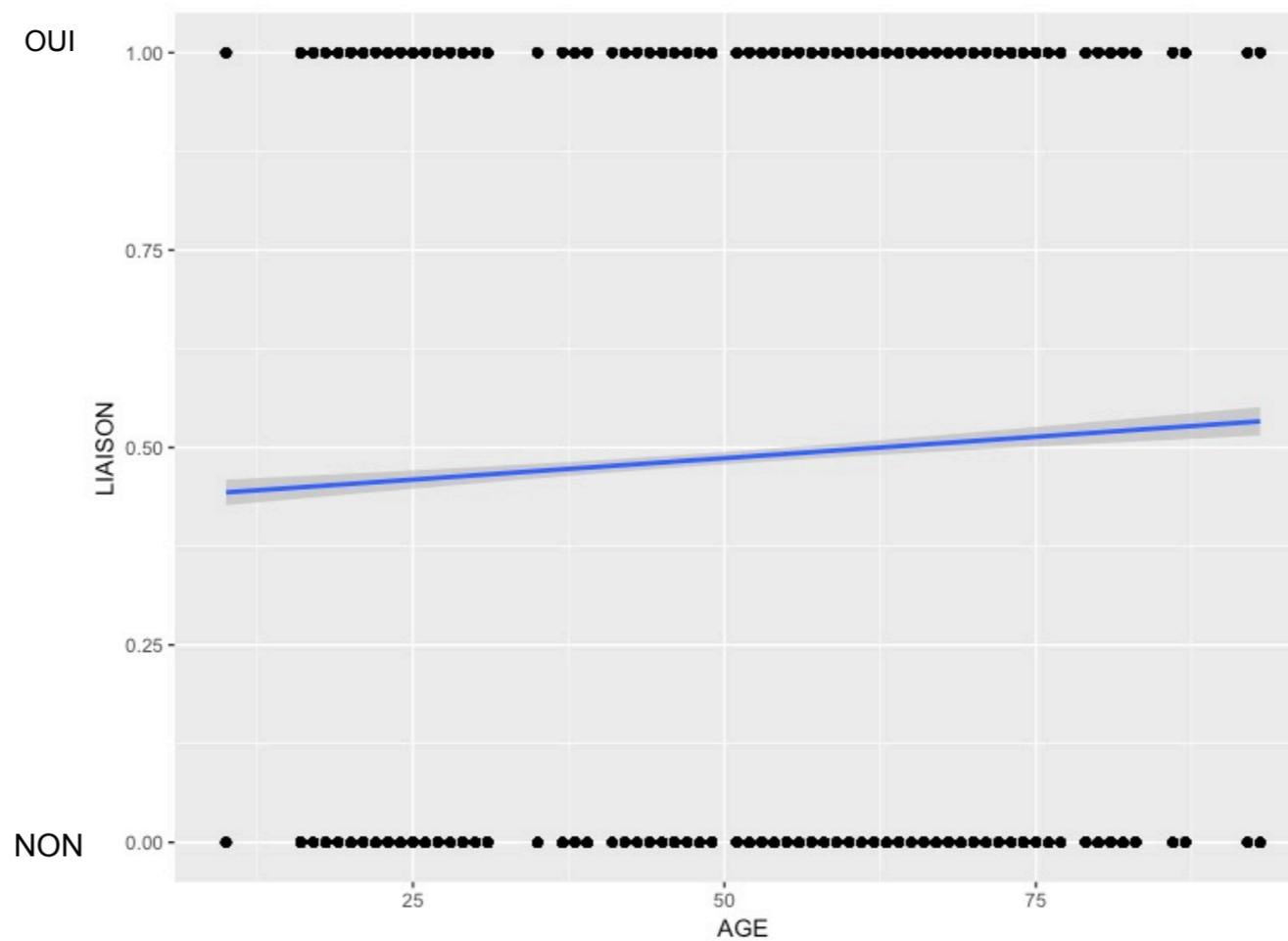
$$\text{LOGIT}[P(\text{Liaison} = \text{oui} \mid \text{AGE})] = -0.2713395 + 0.0043476 \text{ AGE} + \epsilon$$



L'augmentation de l'AGE détermine une plus forte probabilité d'avoir une liaison réalisée

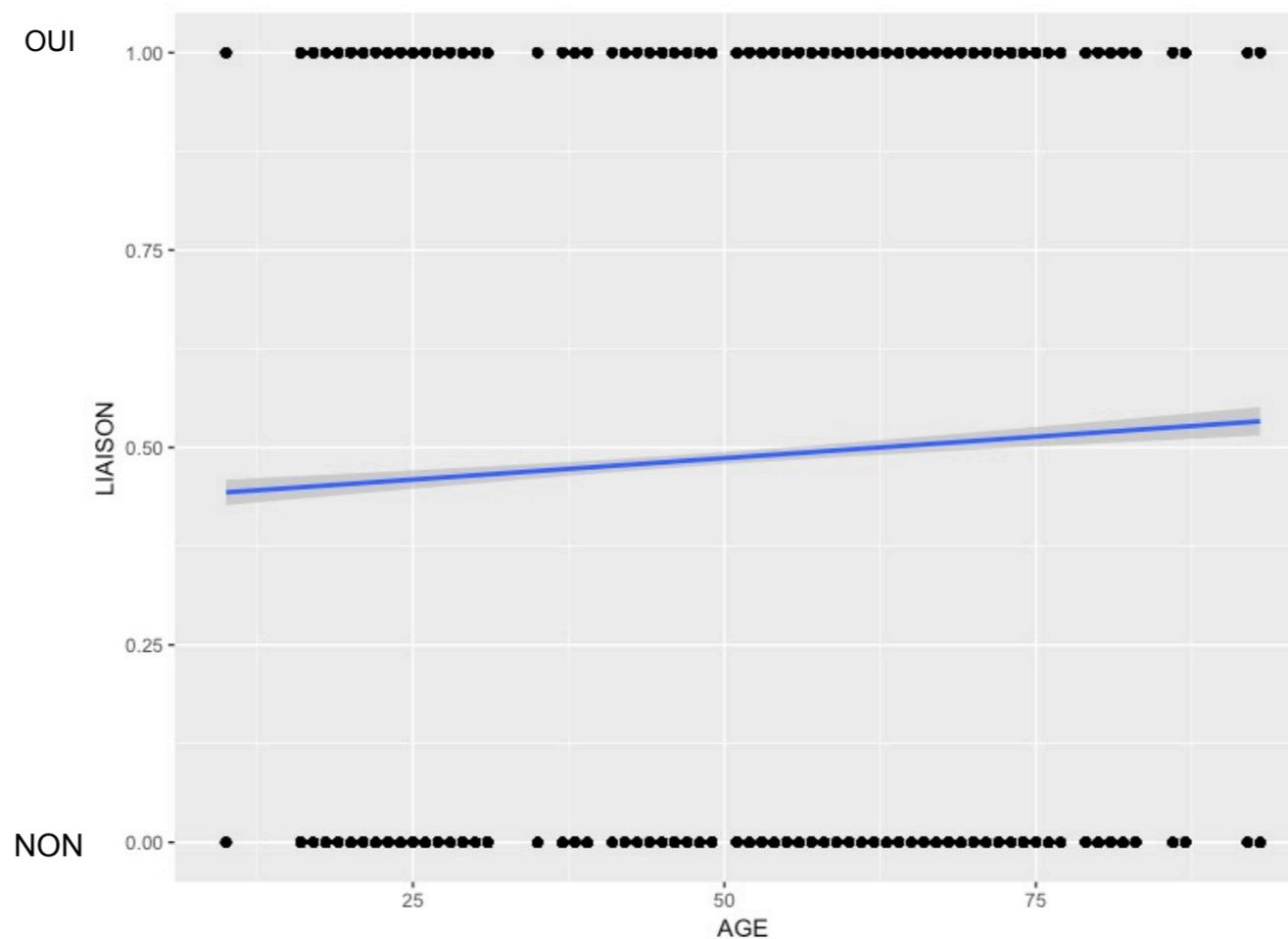


L'augmentation de l'AGE détermine une plus forte probabilité d'avoir une liaison réalisée



Deux mesures d'évaluation d'un model logistique:

L'augmentation de l'AGE détermine une plus forte probabilité d'avoir une liaison réalisée

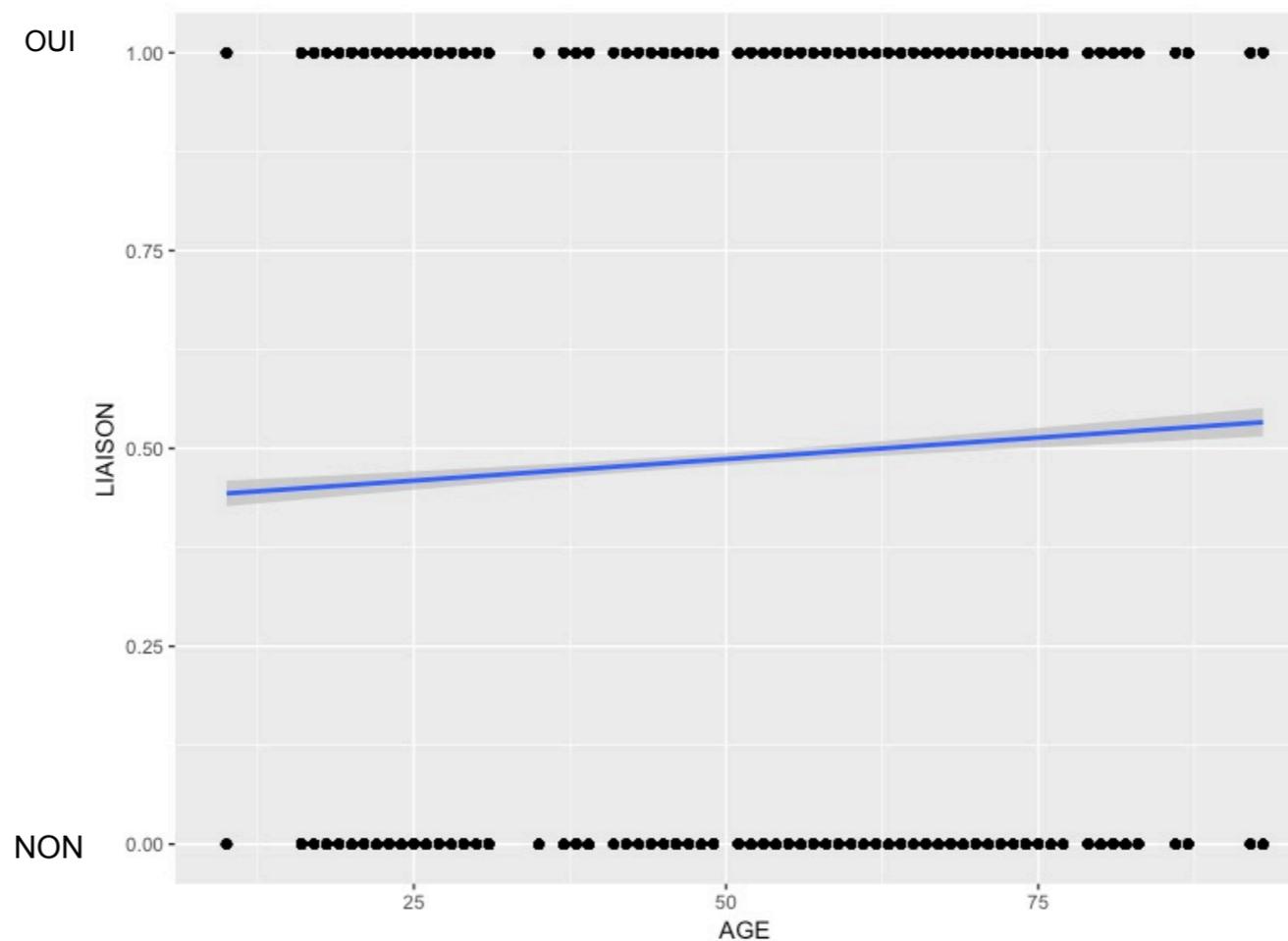


Deux mesures d'évaluation d'un modèle logistique:

1- **AIC**, critère d'information d'Akaike est une mesure de la qualité d'un modèle statistique

$AIC = 2k - 2\log(L)$ où k est le nombre de paramètres à estimer du modèle et L est le maximum de la fonction de vraisemblance du modèle. On choisit le modèle avec le AIC le plus faible

L'augmentation de l'AGE détermine une plus forte probabilité d'avoir une liaison réalisée



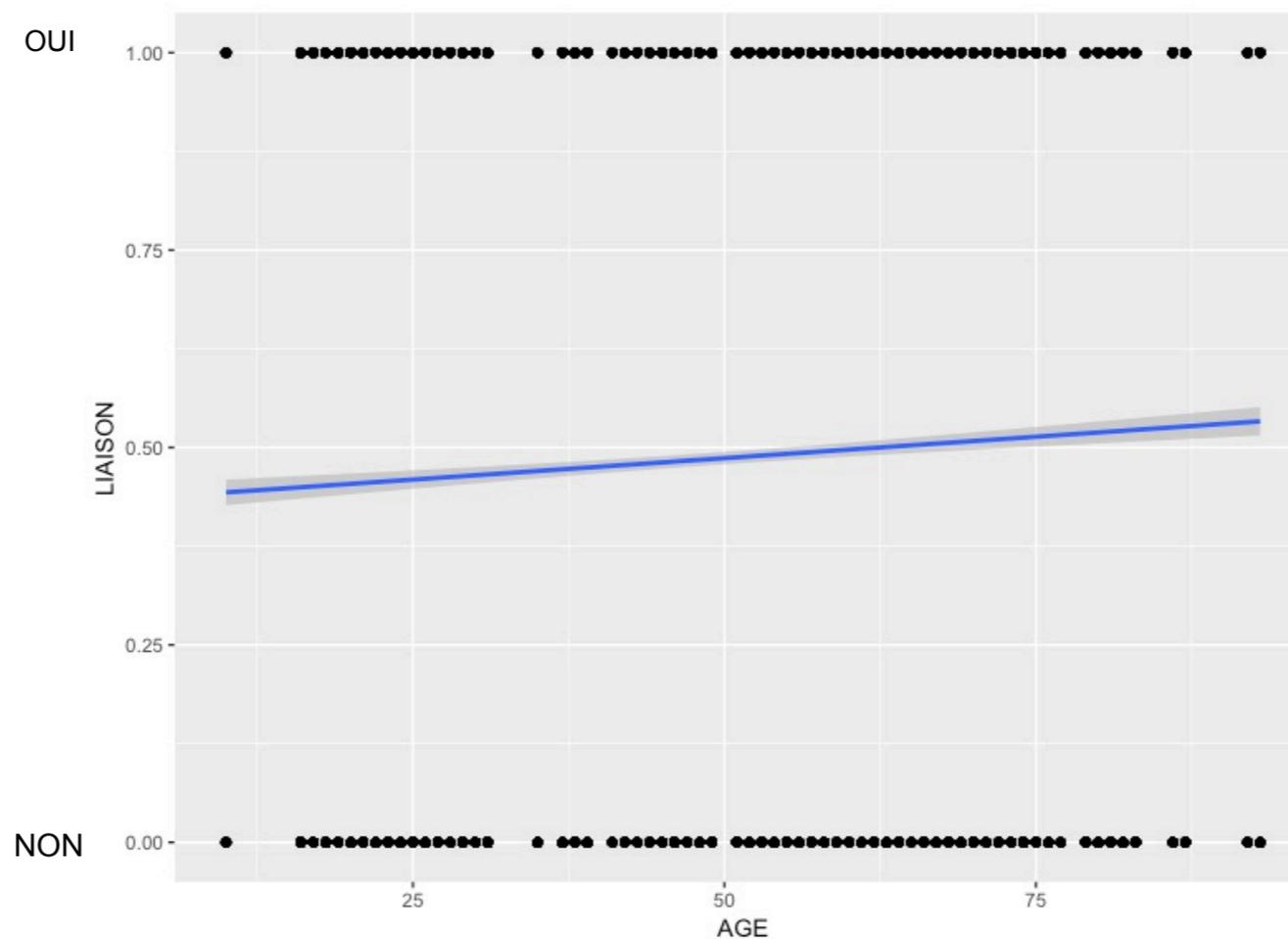
Deux mesures d'évaluation d'un modèle logistique:

1- **AIC**, critère d'information d'Akaike est une mesure de la qualité d'un modèle statistique

$AIC = 2k - 2\log(L)$ où k est le nombre de paramètres à estimer du modèle et L est le maximum de la fonction de vraisemblance du modèle. On choisit le modèle avec le AIC le plus faible

2- **La déviance expliquée du modèle (%)**, est une mesure de la qualité de la prédiction de la régression et représente la variabilité de la variable dépendante Y décrite par le modèle

L'augmentation de l'AGE détermine une plus forte probabilité d'avoir une liaison réalisée



AIC = 21947.7

Déviante expliquée = 0.15%

Deux mesures d'évaluation d'un modèle logistique:

1- **AIC**, critère d'information d'Akaike est une mesure de la qualité d'un modèle statistique

$AIC = 2k - 2\log(L)$ où k est le nombre de paramètres à estimer du modèle et L est le maximum de la fonction de vraisemblance du modèle. On choisit le modèle avec le AIC le plus faible

2- **La déviance expliquée du modèle (%)**, est une mesure de la qualité de la prédiction de la régression et représente la variabilité de la variable dépendante Y décrite par le modèle

4.3.1 Un premier modèle: Liaison \sim AGE- predictions

4.3.1 Un premier modèle: Liaison \sim AGE- predictions

Probabilité de réaliser une liaison par une personne de 50 ans

$$\text{AGE} = 50$$

4.3.1 Un premier modèle: Liaison ~ AGE- predictions

Probabilité de réaliser une liaison par une personne de 50 ans

$$\text{AGE} = 50$$

$$P(\text{Liaison} = \text{oui} \mid \text{AGE}) = \frac{\exp(\beta_0 + \beta_1 \text{AGE} + \epsilon)}{1 + \exp(\beta_0 + \beta_1 \text{AGE} + \epsilon)}$$

4.3.1 Un premier modèle: Liaison ~ AGE- predictions

Probabilité de réaliser une liaison par une personne de 50 ans

$$AGE = 50$$

$$P(Liaison = oui \mid AGE) = \frac{\exp(\beta_0 + \beta_1 AGE + \epsilon)}{1 + \exp(\beta_0 + \beta_1 AGE + \epsilon)}$$

$$P(Liaison = oui \mid AGE = 50) = \frac{\exp(-0.2713395 + 0.0043476 \cdot 50) + \epsilon}{1 + \exp(-0.2713395 + 0.0043476 \cdot 50) + \epsilon}$$

4.3.1 Un premier modèle: Liaison \sim AGE- predictions

Probabilité de réaliser une liaison par une personne de 50 ans

$$AGE = 50$$

$$P(Liaison = oui \mid AGE) = \frac{\exp(\beta_0 + \beta_1 AGE + \epsilon)}{1 + \exp(\beta_0 + \beta_1 AGE + \epsilon)}$$

$$P(Liaison = oui \mid AGE = 50) = \frac{\exp(-0.2713395 + 0.0043476 \cdot 50) + \epsilon}{1 + \exp(-0.2713395 + 0.0043476 \cdot 50) + \epsilon}$$

$$P(Liaison = oui \mid AGE = 50) = 0.491$$

4.3.1 Un premier modèle: Liaison ~ AGE- predictions

Probabilité de réaliser une liaison par une personne de 50 ans

$$AGE = 50$$

$$P(Liaison = oui \mid AGE) = \frac{\exp(\beta_0 + \beta_1 AGE + \epsilon)}{1 + \exp(\beta_0 + \beta_1 AGE + \epsilon)}$$

$$P(Liaison = oui \mid AGE = 50) = \frac{\exp(-0.2713395 + 0.0043476 \cdot 50) + \epsilon}{1 + \exp(-0.2713395 + 0.0043476 \cdot 50) + \epsilon}$$

$$P(Liaison = oui \mid AGE = 50) = 0.491$$

$$P(Liaison = oui \mid AGE = 20) = 0.462$$

4.3.1 Un premier modèle: Liaison \sim AGE- predictions

Probabilité de réaliser une liaison par une personne de 50 ans

$$AGE = 50$$

$$P(Liaison = oui \mid AGE) = \frac{\exp(\beta_0 + \beta_1 AGE + \epsilon)}{1 + \exp(\beta_0 + \beta_1 AGE + \epsilon)}$$

$$P(Liaison = oui \mid AGE = 50) = \frac{\exp(-0.2713395 + 0.0043476 \cdot 50) + \epsilon}{1 + \exp(-0.2713395 + 0.0043476 \cdot 50) + \epsilon}$$

$$P(Liaison = oui \mid AGE = 50) = 0.491$$

$$P(Liaison = oui \mid AGE = 20) = 0.462$$

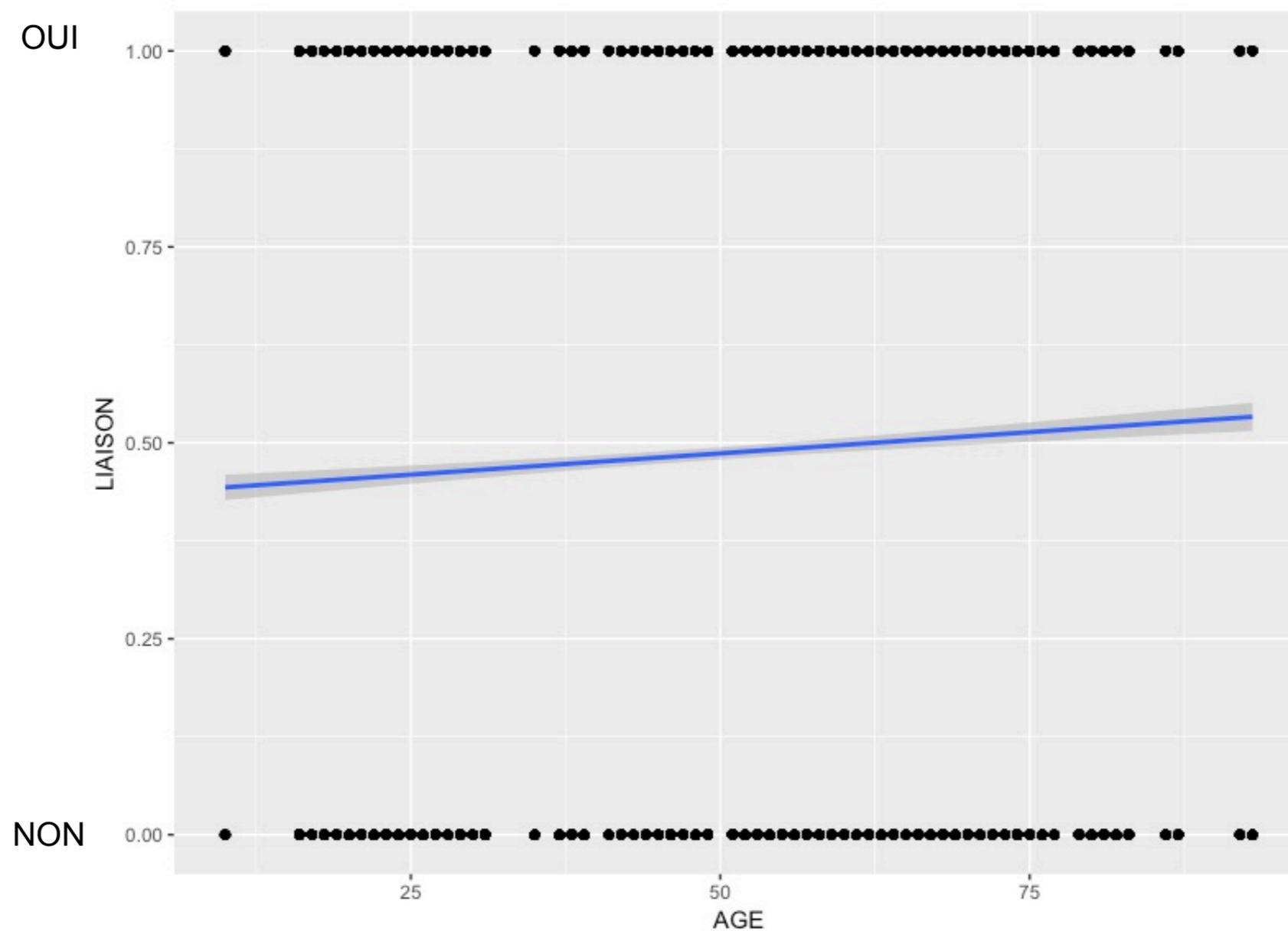
$$P(Liaison = oui \mid AGE = 90) = 0.523$$

4.3.1 Un premier modèle: Liaison ~ AGE- predictions

$$P(\text{Liaison} = \text{oui} \mid \text{AGE} = 20) = 0.462$$

$$P(\text{Liaison} = \text{oui} \mid \text{AGE} = 50) = 0.491$$

$$P(\text{Liaison} = \text{oui} \mid \text{AGE} = 90) = 0.523$$

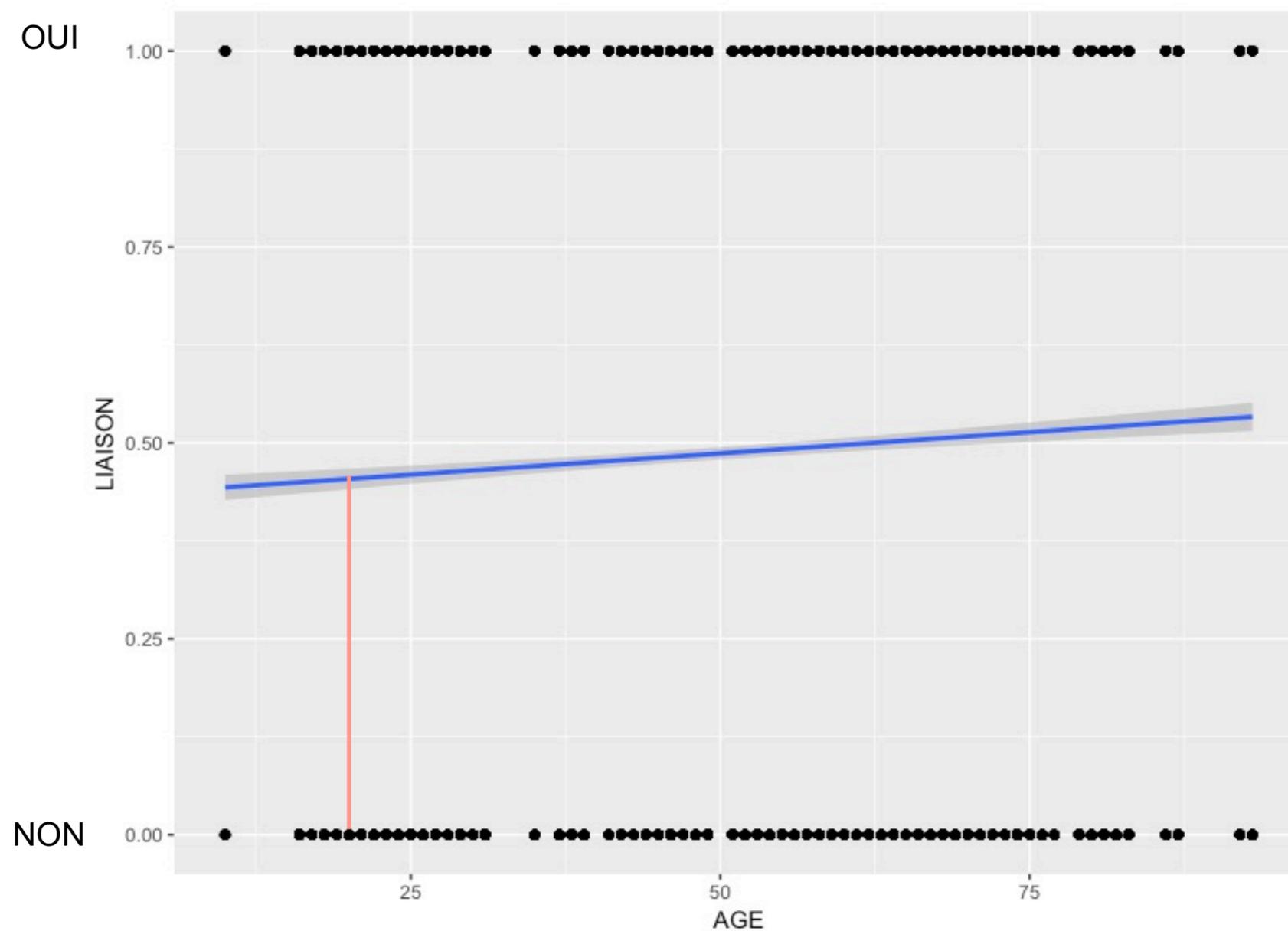


4.3.1 Un premier modèle: $Liaison \sim AGE$ - predictions

$$P(Liaison = oui \mid AGE = 20) = 0.462$$

$$P(Liaison = oui \mid AGE = 50) = 0.491$$

$$P(Liaison = oui \mid AGE = 90) = 0.523$$

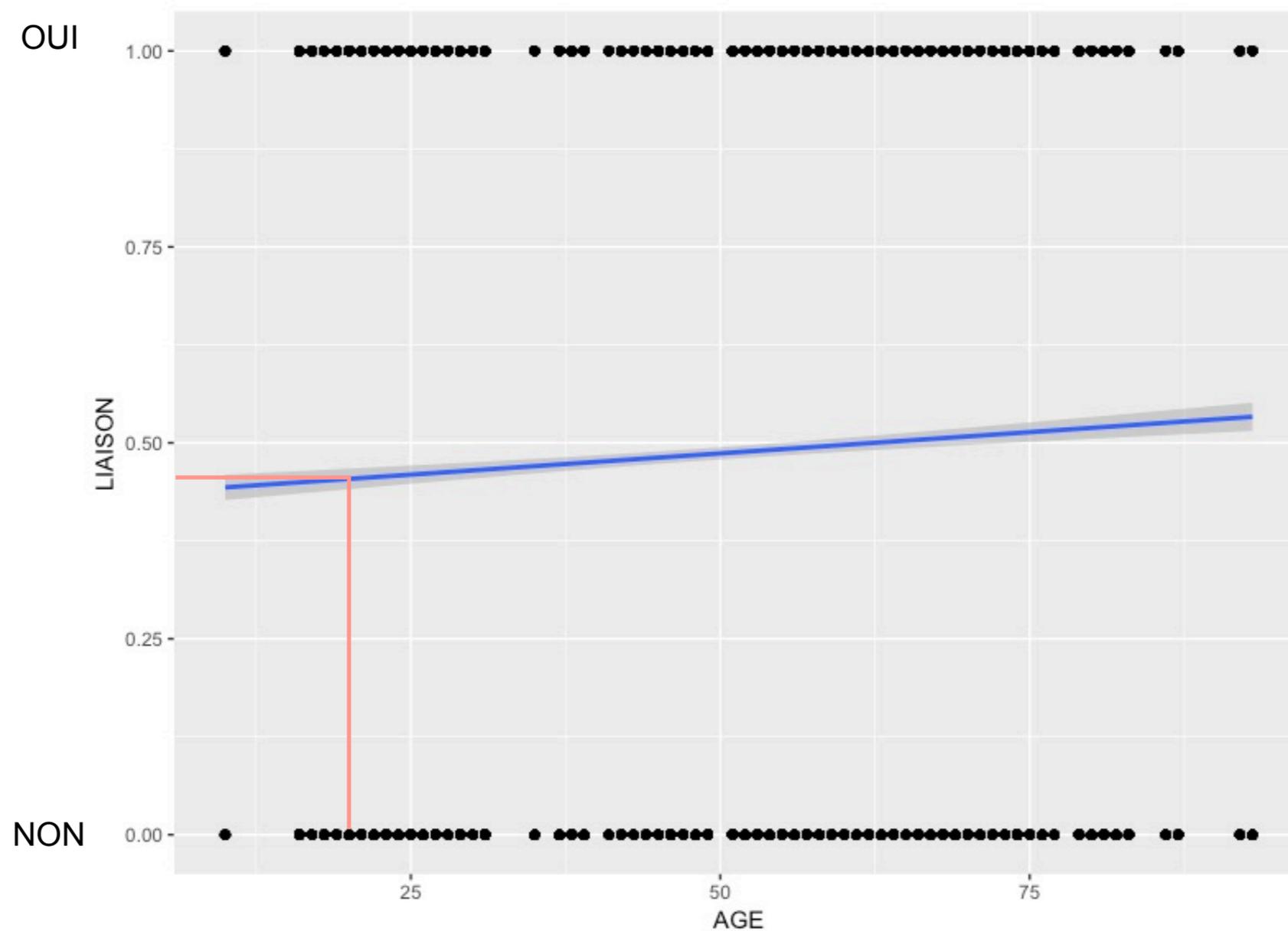


4.3.1 Un premier modèle: Liaison ~ AGE- predictions

$$P(\text{Liaison} = \text{oui} \mid \text{AGE} = 20) = 0.462$$

$$P(\text{Liaison} = \text{oui} \mid \text{AGE} = 50) = 0.491$$

$$P(\text{Liaison} = \text{oui} \mid \text{AGE} = 90) = 0.523$$

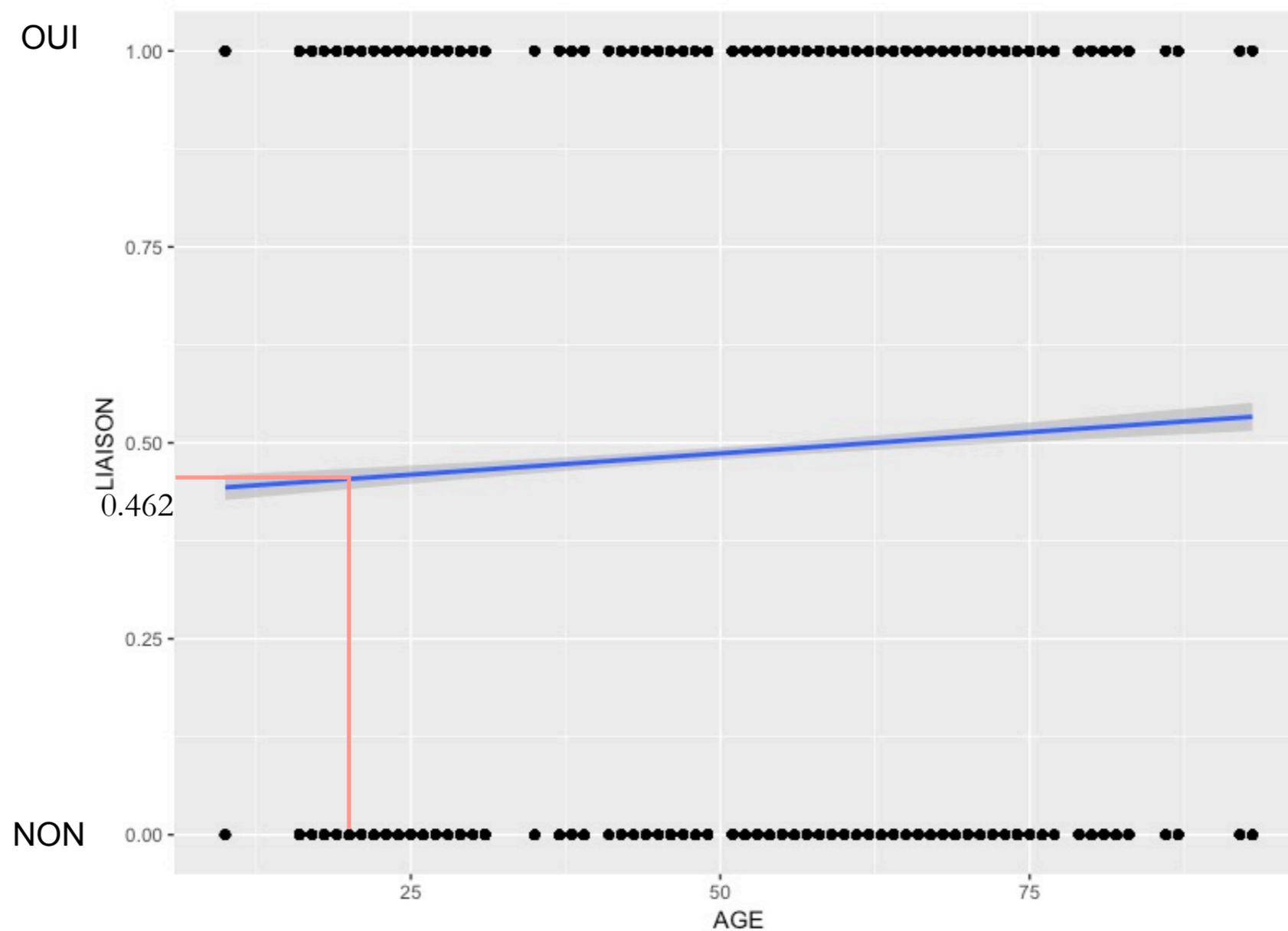


4.3.1 Un premier modèle: Liaison ~ AGE- predictions

$$P(\text{Liaison} = \text{oui} \mid \text{AGE} = 20) = 0.462$$

$$P(\text{Liaison} = \text{oui} \mid \text{AGE} = 50) = 0.491$$

$$P(\text{Liaison} = \text{oui} \mid \text{AGE} = 90) = 0.523$$

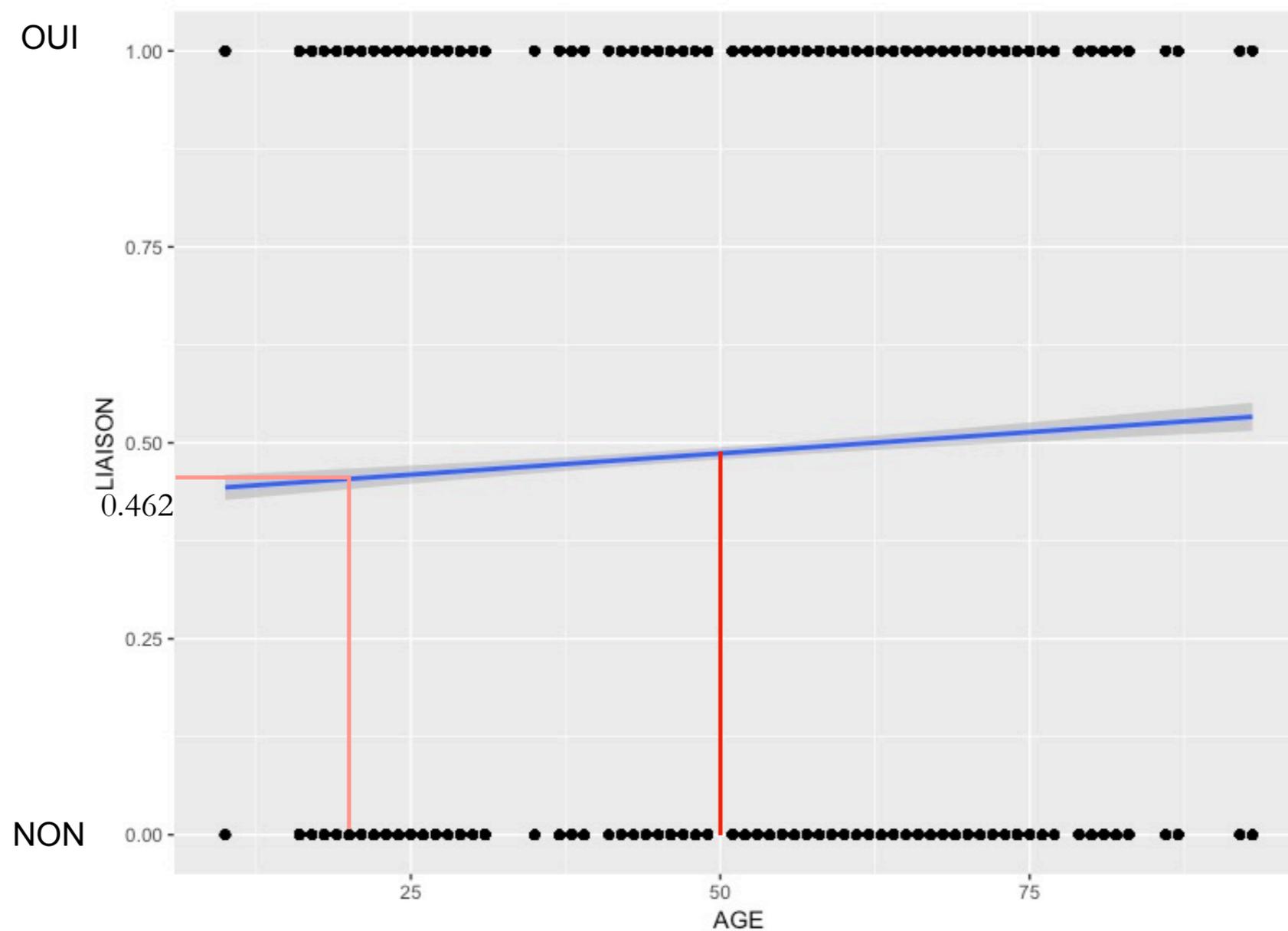


4.3.1 Un premier modèle: Liaison ~ AGE- predictions

$$P(\text{Liaison} = \text{oui} \mid \text{AGE} = 20) = 0.462$$

$$P(\text{Liaison} = \text{oui} \mid \text{AGE} = 50) = 0.491$$

$$P(\text{Liaison} = \text{oui} \mid \text{AGE} = 90) = 0.523$$

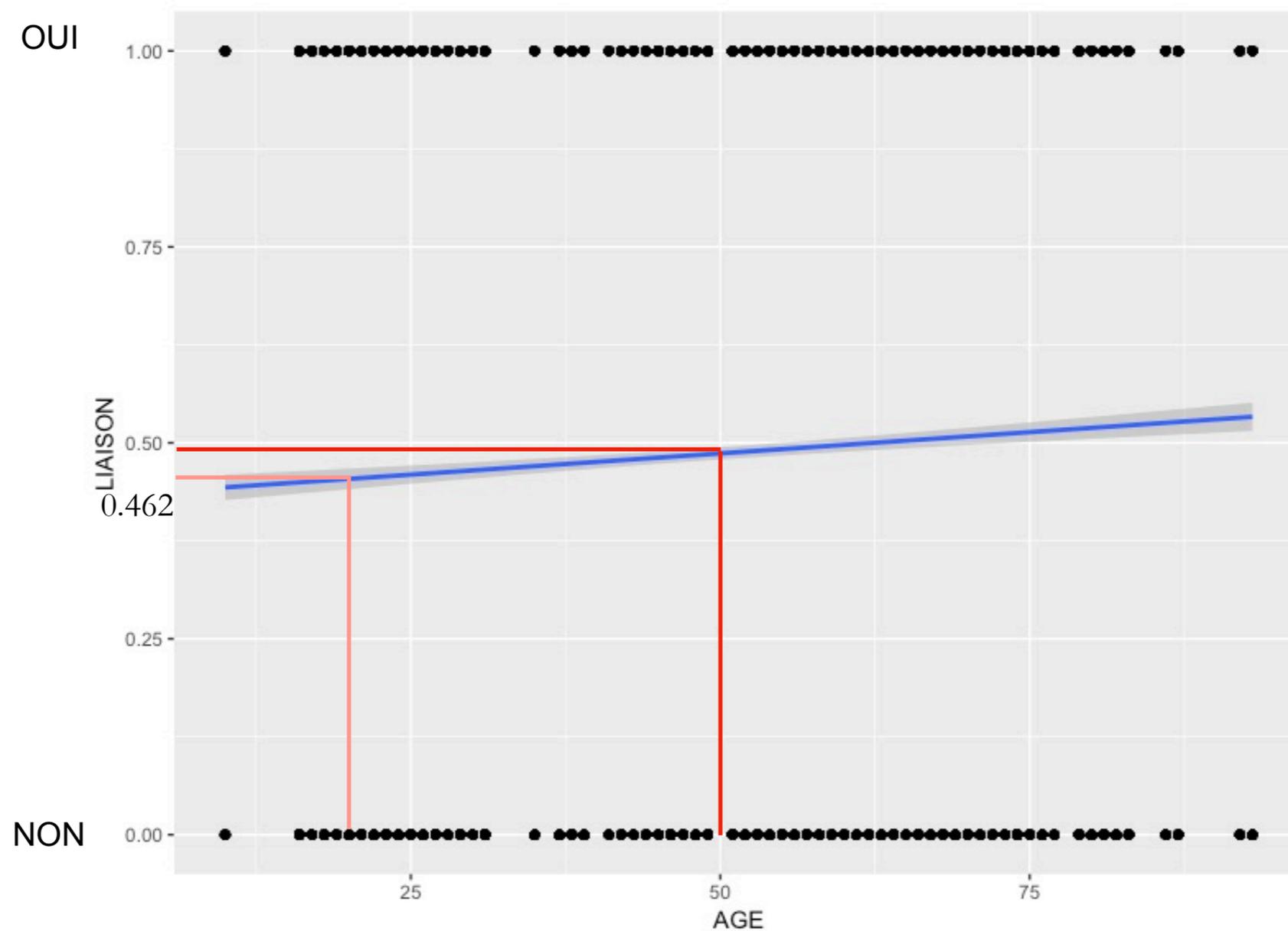


4.3.1 Un premier modèle: $Liaison \sim AGE$ - predictions

$$P(Liaison = oui \mid AGE = 20) = 0.462$$

$$P(Liaison = oui \mid AGE = 50) = 0.491$$

$$P(Liaison = oui \mid AGE = 90) = 0.523$$

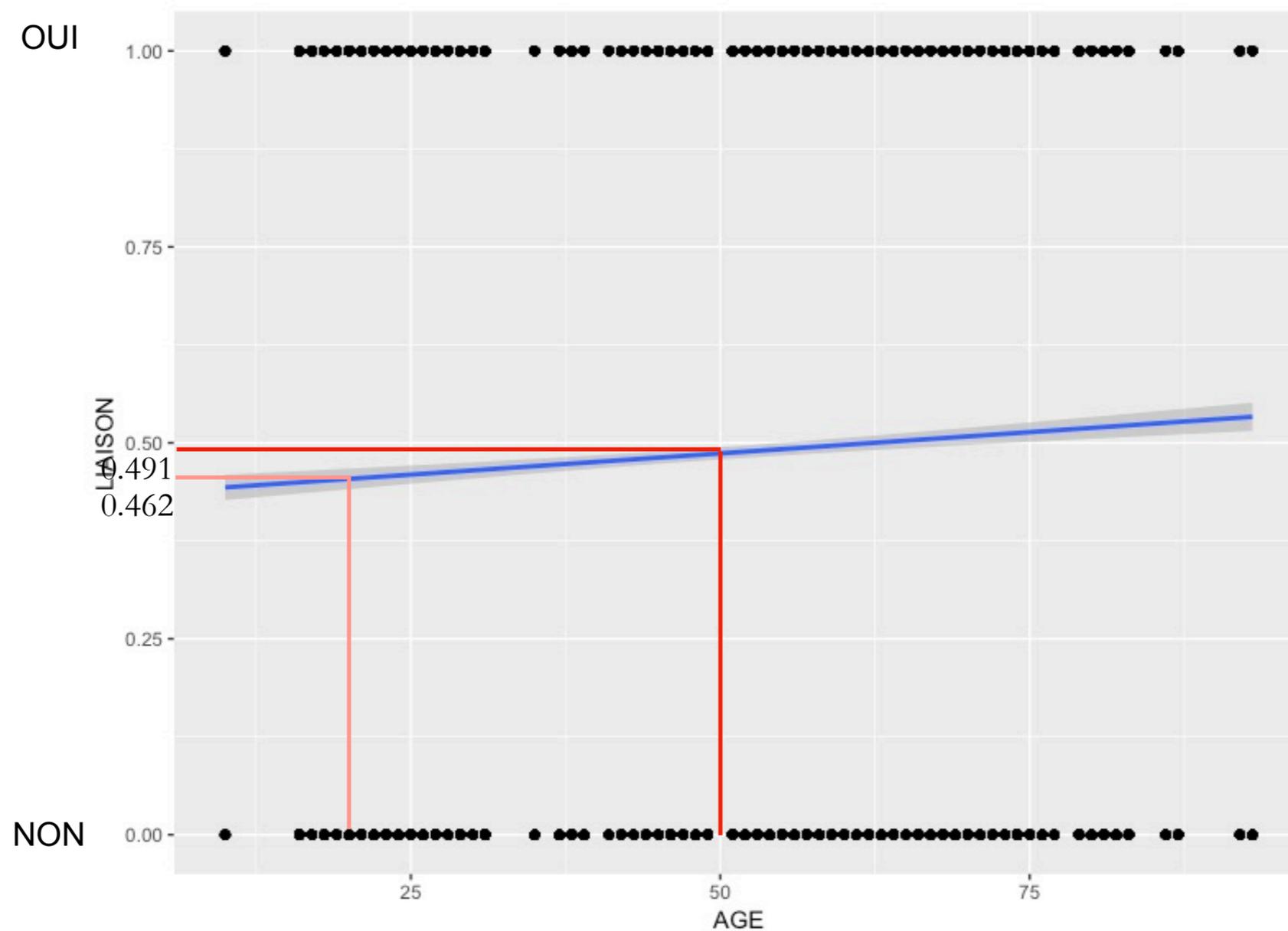


4.3.1 Un premier modèle: Liaison ~ AGE- predictions

$$P(\text{Liaison} = \text{oui} \mid \text{AGE} = 20) = 0.462$$

$$P(\text{Liaison} = \text{oui} \mid \text{AGE} = 50) = 0.491$$

$$P(\text{Liaison} = \text{oui} \mid \text{AGE} = 90) = 0.523$$

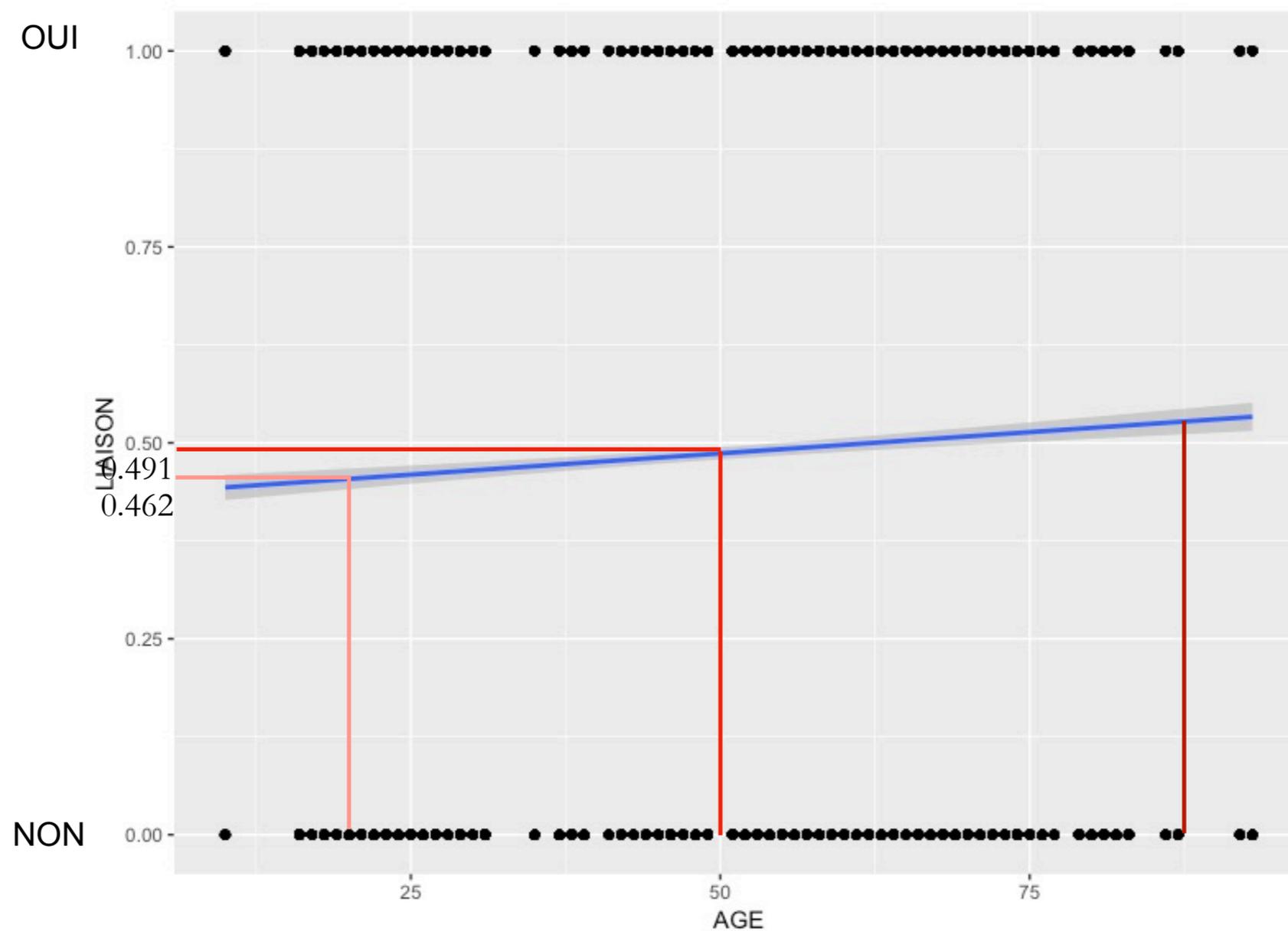


4.3.1 Un premier modèle: Liaison ~ AGE- predictions

$$P(\text{Liaison} = \text{oui} \mid \text{AGE} = 20) = 0.462$$

$$P(\text{Liaison} = \text{oui} \mid \text{AGE} = 50) = 0.491$$

$$P(\text{Liaison} = \text{oui} \mid \text{AGE} = 90) = 0.523$$

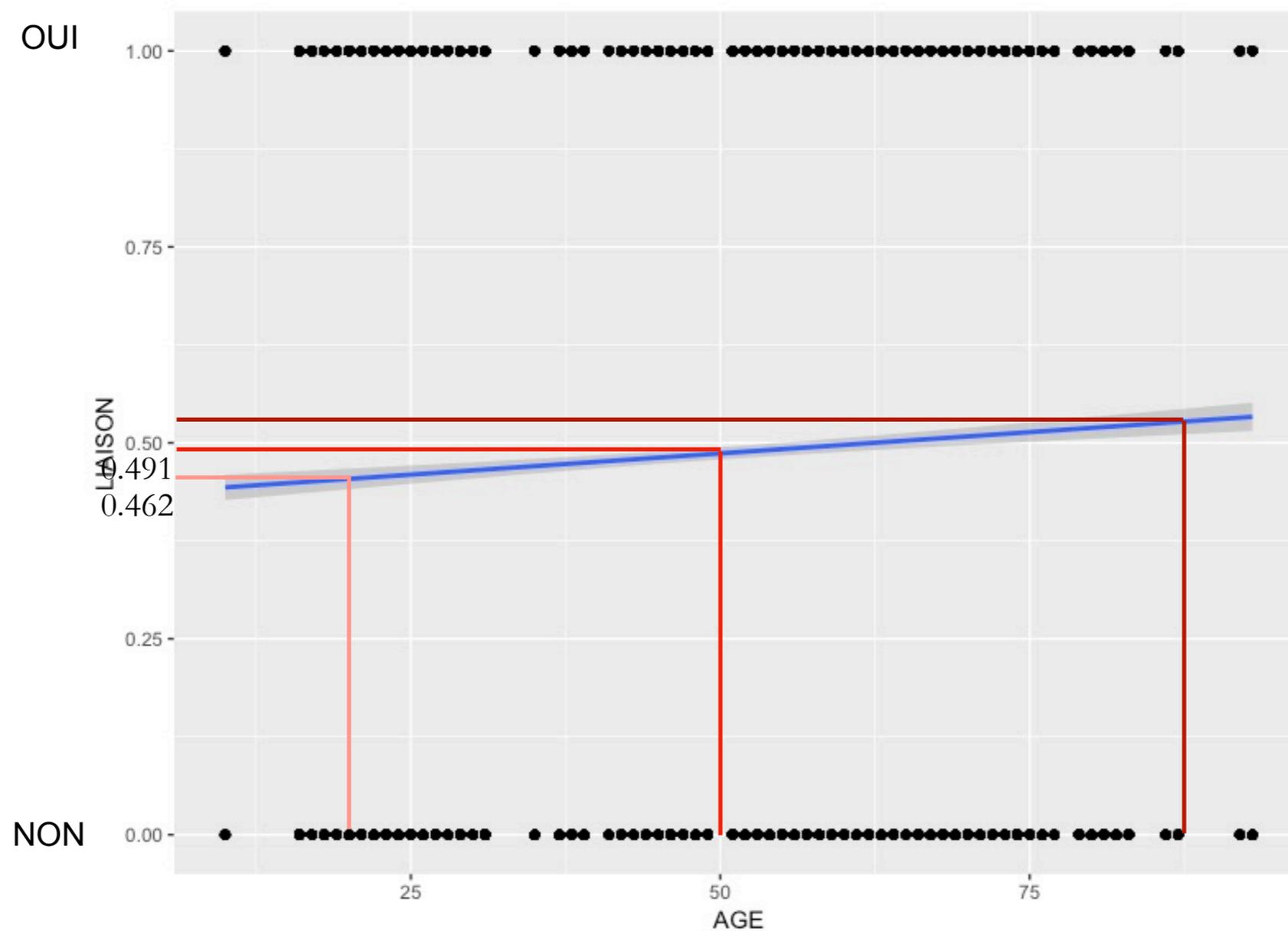


4.3.1 Un premier modèle: Liaison ~ AGE- predictions

$$P(\text{Liaison} = \text{oui} \mid \text{AGE} = 20) = 0.462$$

$$P(\text{Liaison} = \text{oui} \mid \text{AGE} = 50) = 0.491$$

$$P(\text{Liaison} = \text{oui} \mid \text{AGE} = 90) = 0.523$$

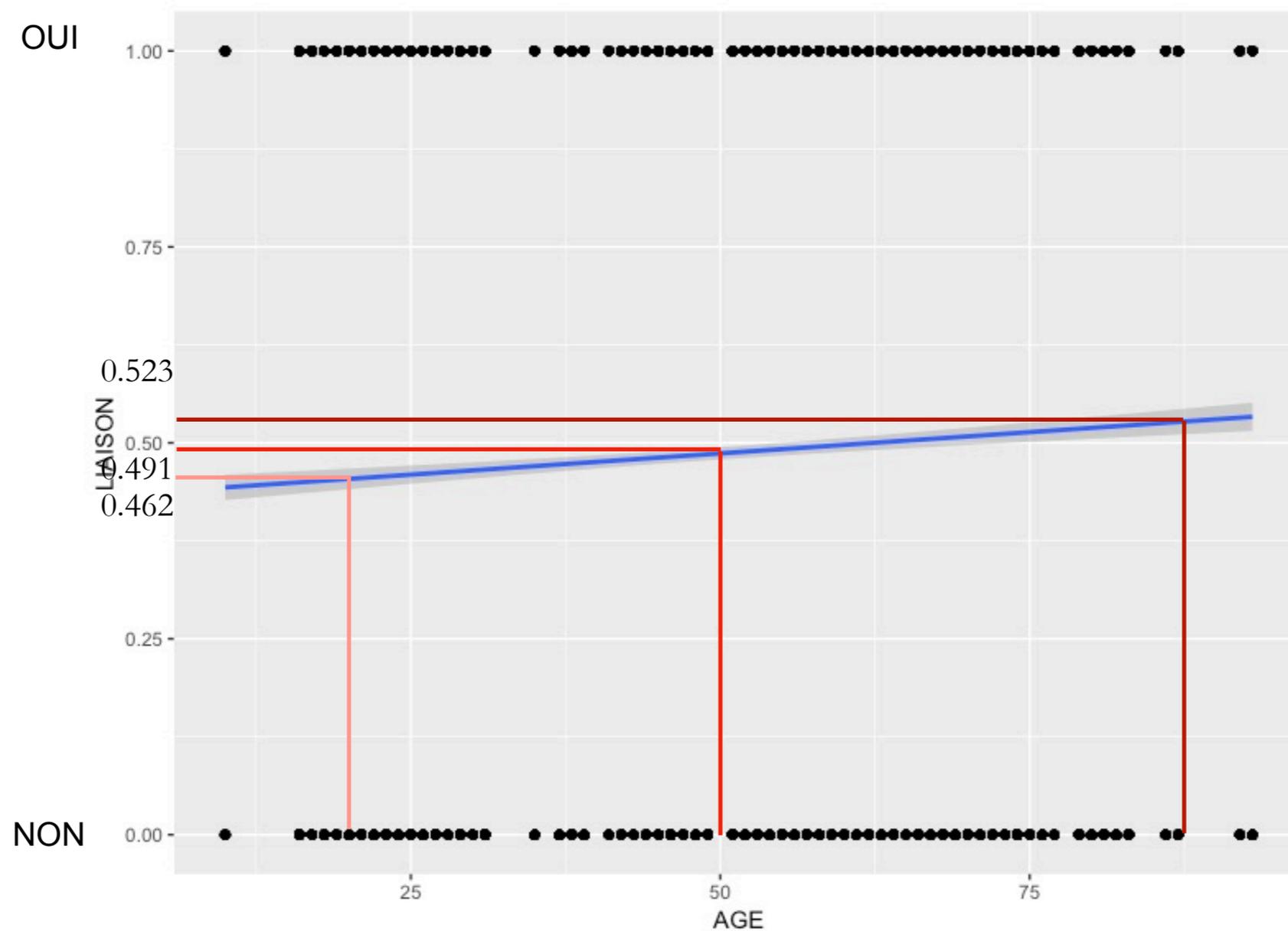


4.3.1 Un premier modèle: Liaison ~ AGE- predictions

$$P(\text{Liaison} = \text{oui} \mid \text{AGE} = 20) = 0.462$$

$$P(\text{Liaison} = \text{oui} \mid \text{AGE} = 50) = 0.491$$

$$P(\text{Liaison} = \text{oui} \mid \text{AGE} = 90) = 0.523$$



4.4 Un modèle plus complexe : Liaison \sim AGE + SEXE

$$Y = \beta_0 + \beta_1\alpha_1 + \beta_2\alpha_2 + \epsilon$$

4.4 Un modèle plus complexe : Liaison \sim AGE + SEXE

$$Y = \beta_0 + \beta_1\alpha_1 + \beta_2\alpha_2 + \epsilon$$

On essaie de trouver une relation entre l'AGE et le SEXE du locuteur et la réalisation de la liaison

$$Liaison = \beta_0 + \beta_1 AGE + \beta_2 SEXE + \epsilon$$

4.4 Un modèle plus complexe : Liaison \sim AGE + SEXE

$$Y = \beta_0 + \beta_1\alpha_1 + \beta_2\alpha_2 + \epsilon$$

On essaie de trouver un relation entre l'AGE et le SEXE du locuteur et la réalisation de la liaison

$$Liaison = \beta_0 + \beta_1 AGE + \beta_2 SEXE + \epsilon$$

Coefficients:

(Intercept)	-0.2515769
AGE	0.0044065
SEXE-H	-0.0476284

4.4 Un modèle plus complexe : Liaison \sim AGE + SEXE

$$Y = \beta_0 + \beta_1\alpha_1 + \beta_2\alpha_2 + \epsilon$$

On essaie de trouver une relation entre l'AGE et le SEXE du locuteur et la réalisation de la liaison

$$Liaison = \beta_0 + \beta_1 AGE + \beta_2 SEXE + \epsilon$$

Coefficients:

(Intercept)	-0.2515769
AGE	0.0044065
SEXE-H	-0.0476284

$$Y = -0.2515769 + 0.0044065 AGE + (-0.0476284) SEXE-H + \epsilon$$

4.4 Un modèle plus complexe : Liaison ~ AGE + SEXE

$$Y = \beta_0 + \beta_1 \alpha_1 + \beta_2 \alpha_2 + \epsilon$$

On essaie de trouver une relation entre l'AGE et le SEXE du locuteur et la réalisation de la liaison

$$Liaison = \beta_0 + \beta_1 AGE + \beta_2 SEXE + \epsilon$$

Coefficients:

(Intercept)	-0.2515769
AGE	0.0044065
SEXE-H	-0.0476284

$$Y = -0.2515769 + 0.0044065 AGE + (-0.0476284) SEXE-H + \epsilon$$

signe positif

+AGE +LIAISON (AGE favorise la réalisation)

|

variable quantitative

4.4 Un modèle plus complexe : Liaison ~ AGE + SEXE

$$Y = \beta_0 + \beta_1 \alpha_1 + \beta_2 \alpha_2 + \epsilon$$

On essaie de trouver une relation entre l'AGE et le SEXE du locuteur et la réalisation de la liaison

$$Liaison = \beta_0 + \beta_1 AGE + \beta_2 SEXE + \epsilon$$

Coefficients:

(Intercept)	-0.2515769
AGE	0.0044065
SEXE-H	-0.0476284

$$Y = -0.2515769 + 0.0044065 AGE + (-0.0476284) SEXE-H + \epsilon$$

signe positif

+AGE +LIAISON (AGE favorise la réalisation)

|

variable quantitative

signe négatif

-SEXE-H -LIAISON (SEXE-H défavorise la réalisation)

|

variable qualitative

4.4 Un modèle plus complexe : $Liaison \sim AGE + SEXE$ - prédictions

Probabilité de réaliser une liaison par une femme de 70 ans

$$P(Liaison = oui \mid AGE=70, SEXE=F) = 0.4385201$$

4.4 Un modèle plus complexe : $Liaison \sim AGE + SEXE$ - prédictions

Probabilité de réaliser une liaison par une femme de 70 ans

$$P(Liaison = oui \mid AGE=70, SEXE=F) = 0.4385201$$

Probabilité de réaliser une liaison par un homme de 70 ans

$$P(Liaison = oui \mid AGE=70, SEXE=H) = 0.4268295$$

4.4 Un modèle plus complexe : $Liaison \sim AGE + SEXE$ - prédictions

Probabilité de réaliser une liaison par une femme de 70 ans

$$P(Liaison = oui \mid AGE=70, SEXE=F) = 0.4385201$$

Probabilité de réaliser une liaison par un homme de 70 ans

$$P(Liaison = oui \mid AGE=70, SEXE=H) = 0.4268295$$

Probabilité de réaliser une liaison par un homme de 18 ans

$$P(Liaison = oui \mid AGE=18, SEXE=H) = 0.3834355$$

Limites de la RL classique

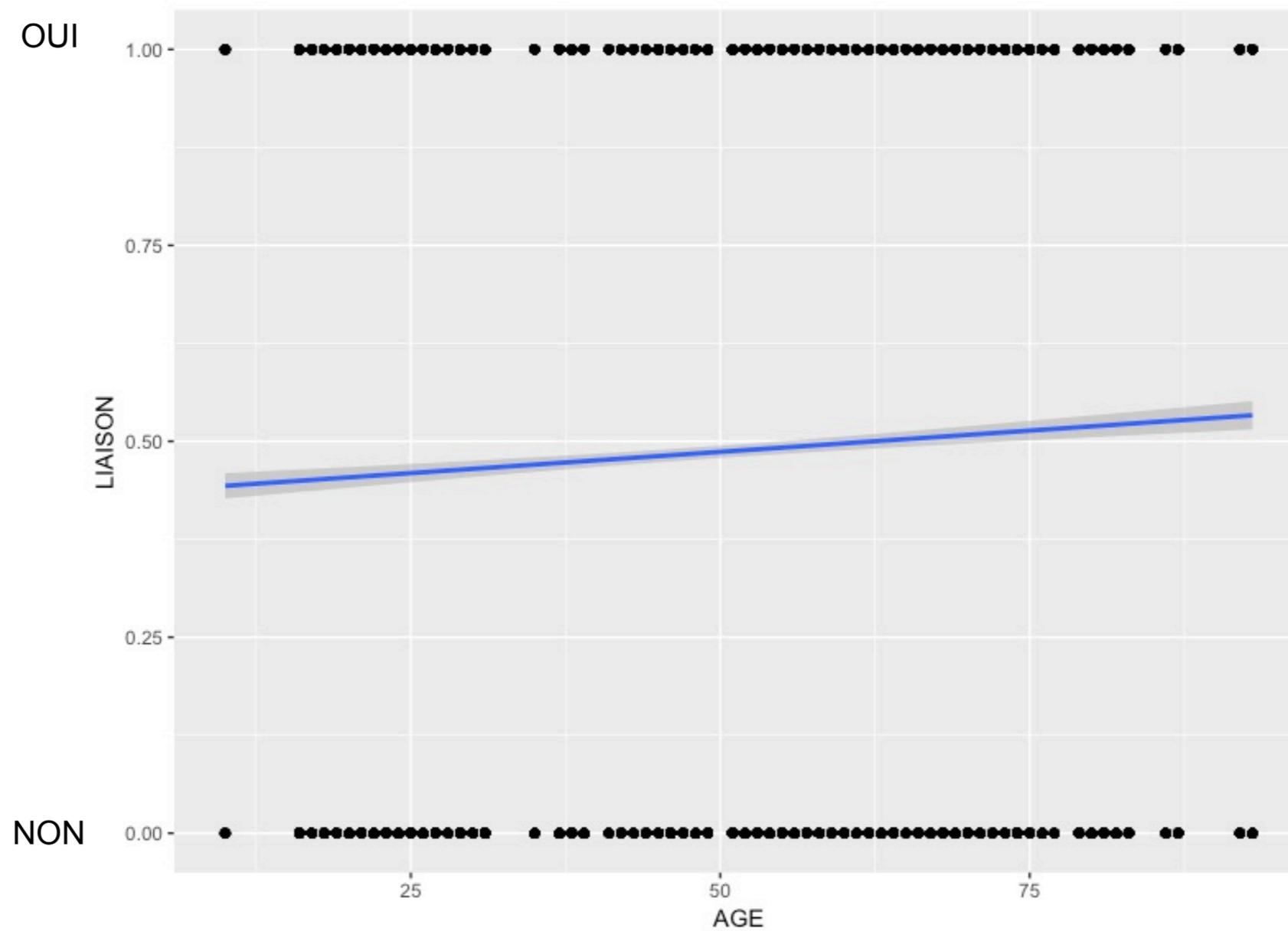
La RL se base sur des fonctions de relation entre Y (variable à expliquer) et X_n (variables explicatives) qui sont essentiellement linéaires

Limites de la RL classique

La RL se base sur des fonctions de relation entre Y (variable à expliquer) et X_n (variables explicatives) qui sont essentiellement linéaires

Possibilité d'utiliser des fonctions de lissage (*smooth functions*) qui ne sont pas forcément linéaires

Limites de la RL classique



$$Y = \beta_0 + \beta_1 \alpha_1 + \epsilon$$

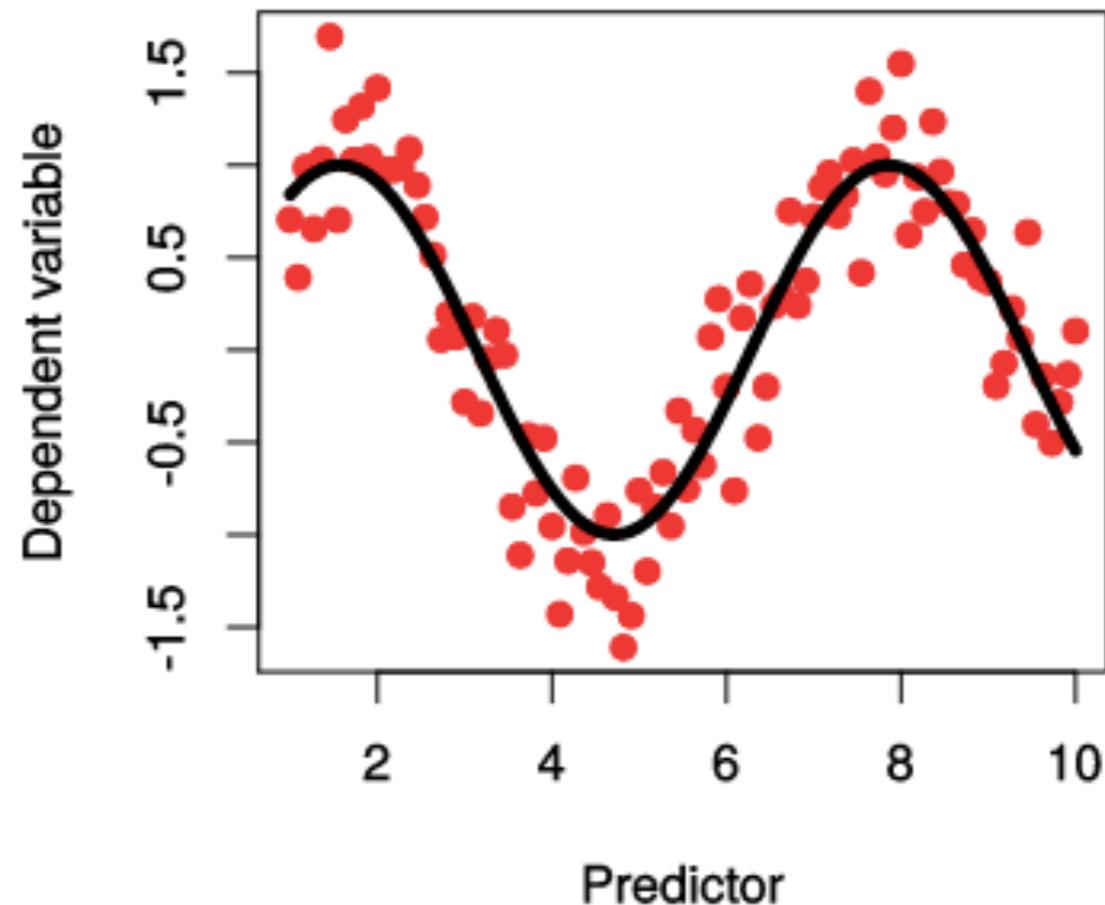
Les méthodes GAM s'appuient sur des transformations empiriques des variables explicatives par des techniques de *lissage*

4.5 Modèle additif généralisé (GAM, Wood 2016)

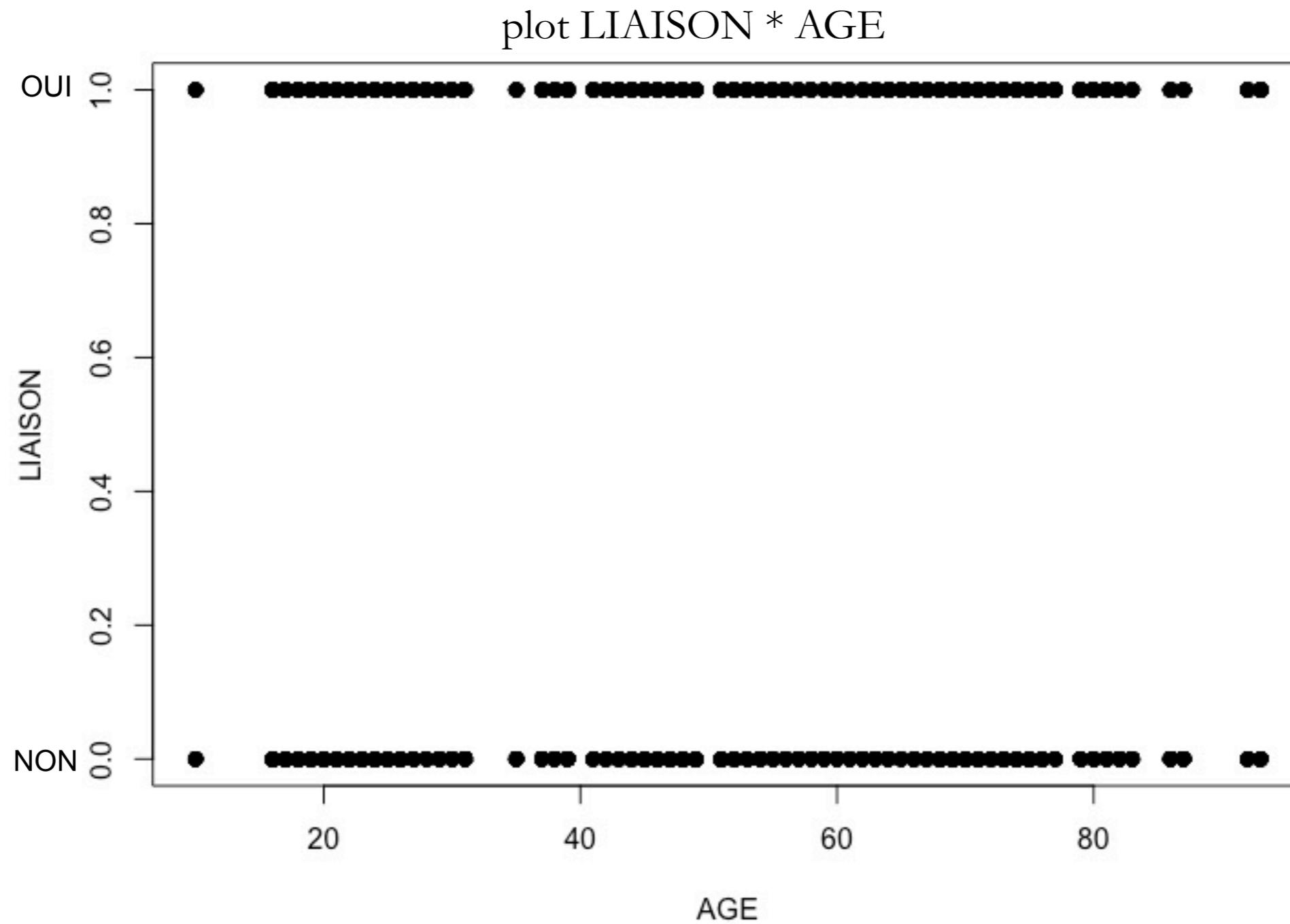
$$Y = \beta_0 + \beta_1 \alpha_1 + \epsilon$$

Les méthodes GAM s'appuient sur des transformations empiriques des variables explicatives par des techniques de *lissage*

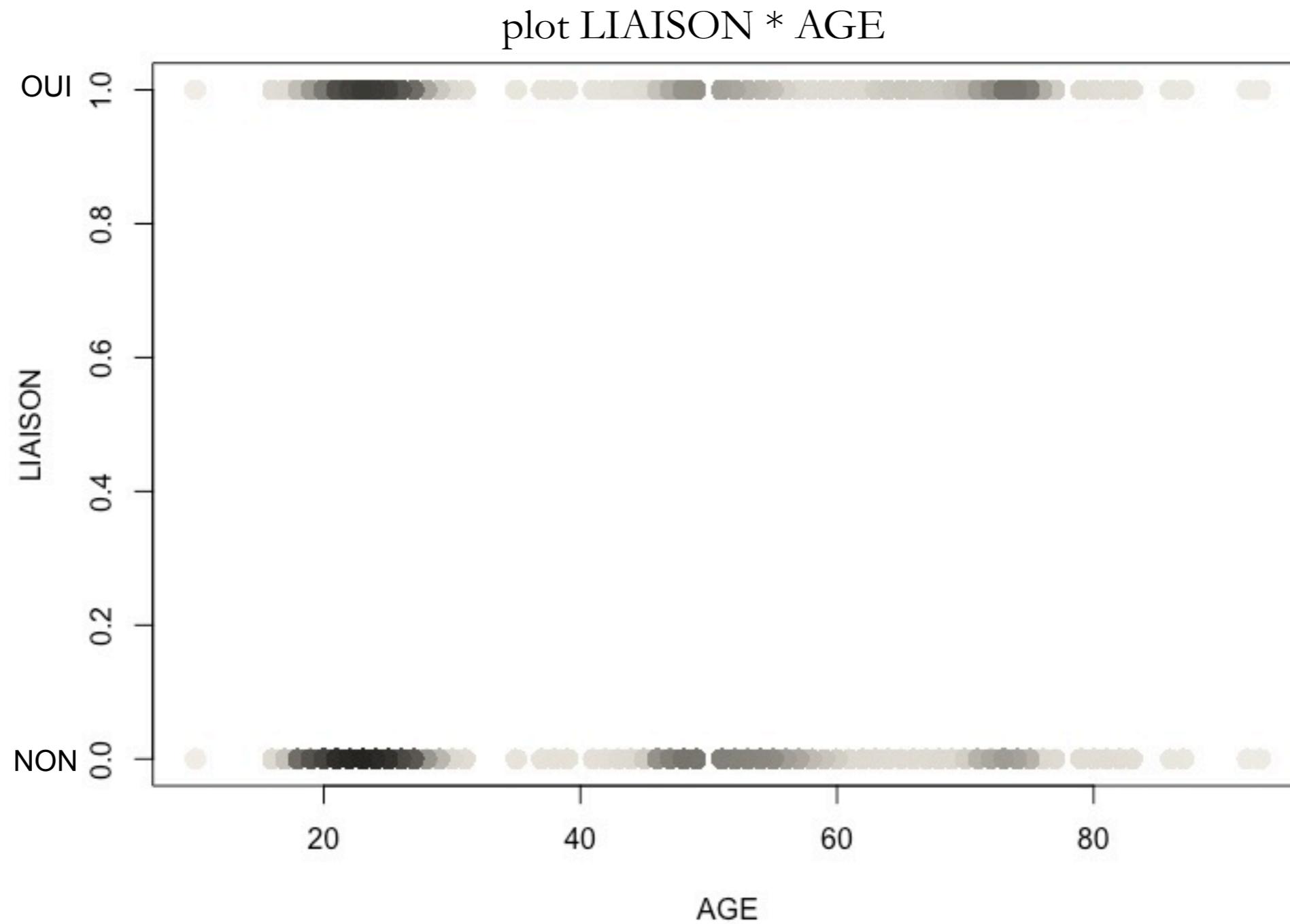
$$g(E(Y)) = \beta_0 + f_1(\alpha_1) + \epsilon$$



4.5 Modèle additif généralisé (GAM, Wood 2016)

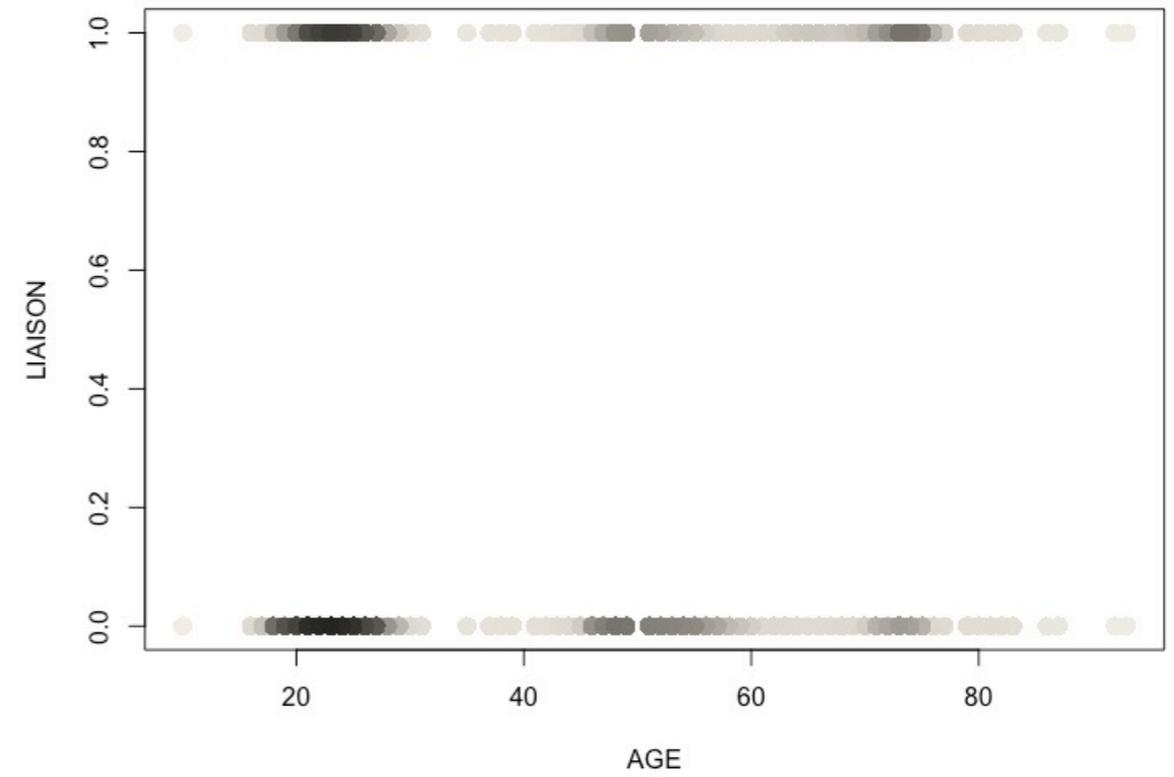


Visualisation de la densité



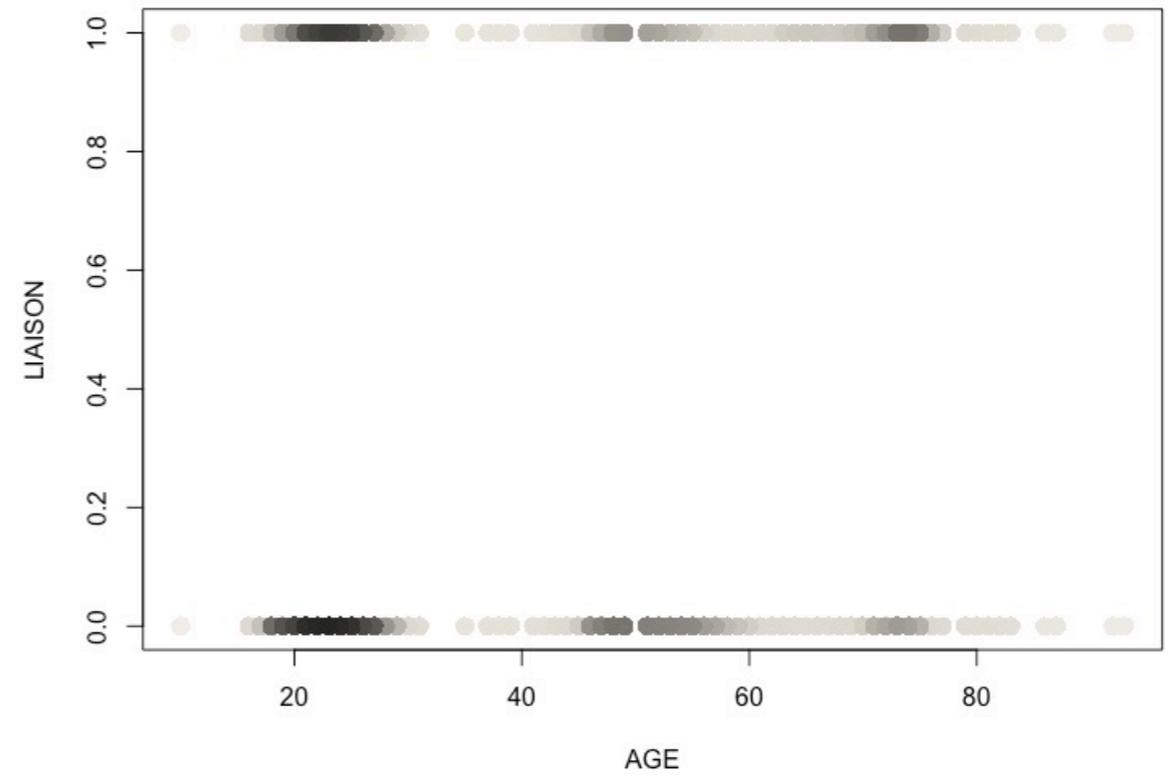
4.5 Modèle additif généralisé (GAM, Wood 2016)

Fonction plus adaptative à la distribution réelle des données

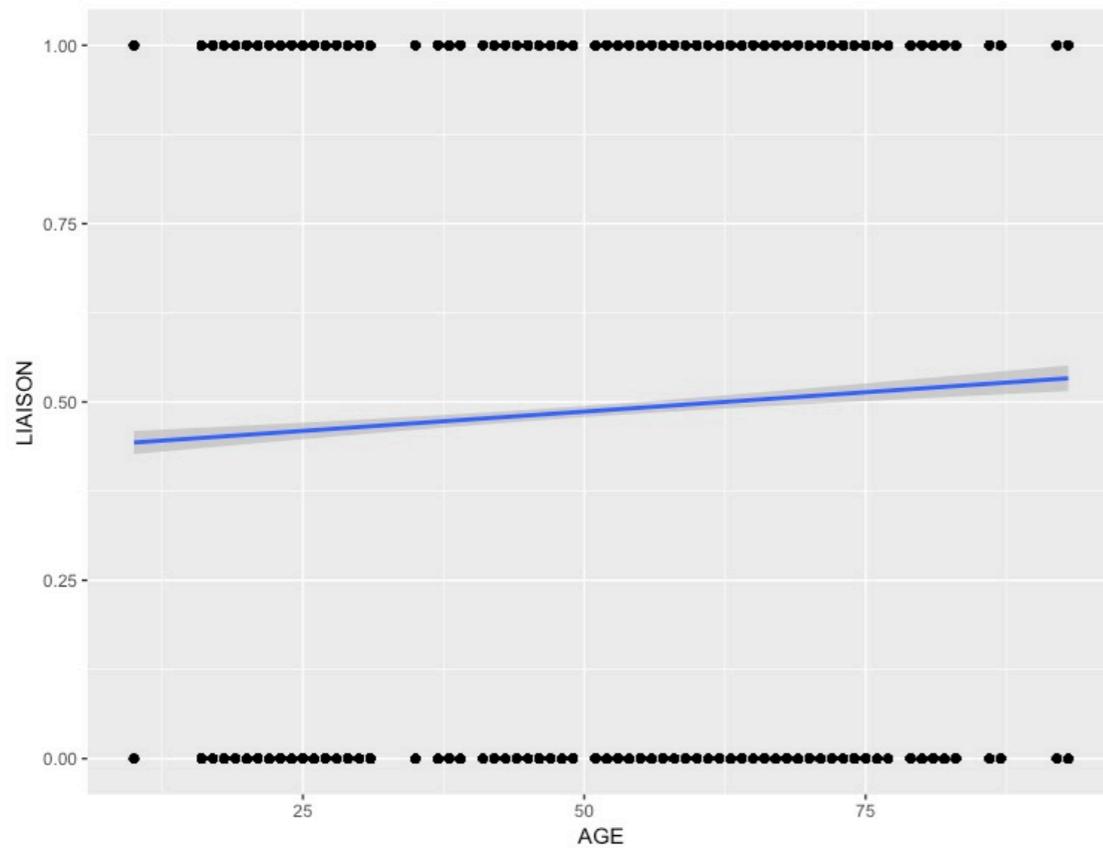


4.5 Modèle additif généralisé (GAM, Wood 2016)

Fonction plus adaptative à la distribution réelle des données



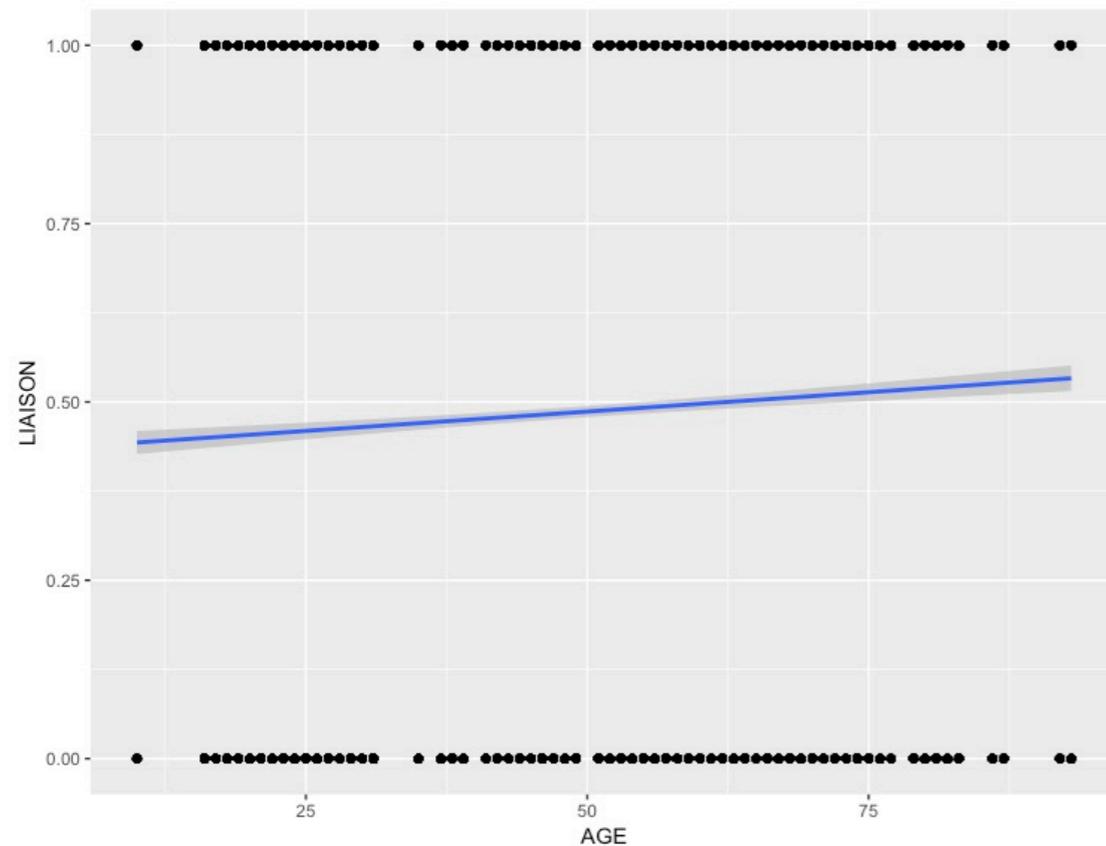
$$mod.AGE_{RL} = \beta_0 + \beta_1 AGE$$



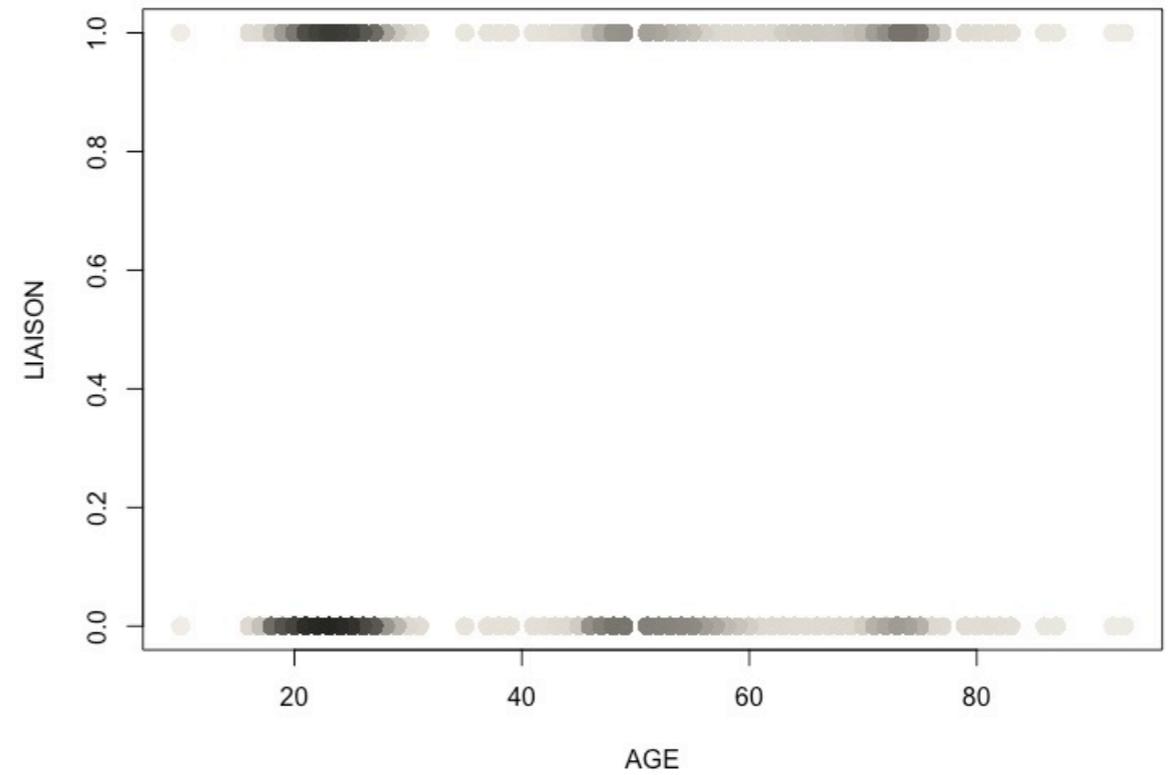
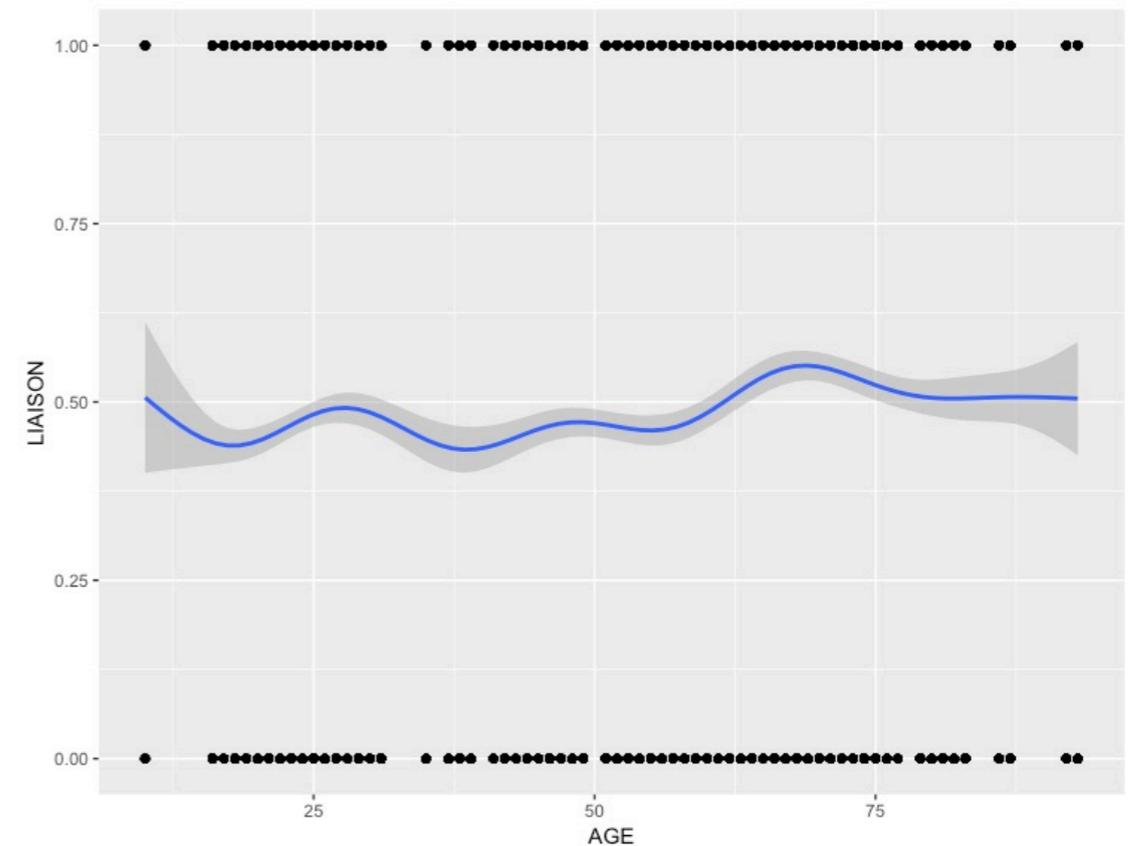
4.5 Modèle additif généralisé (GAM, Wood 2016)

Fonction plus adaptative à la distribution réelle des données

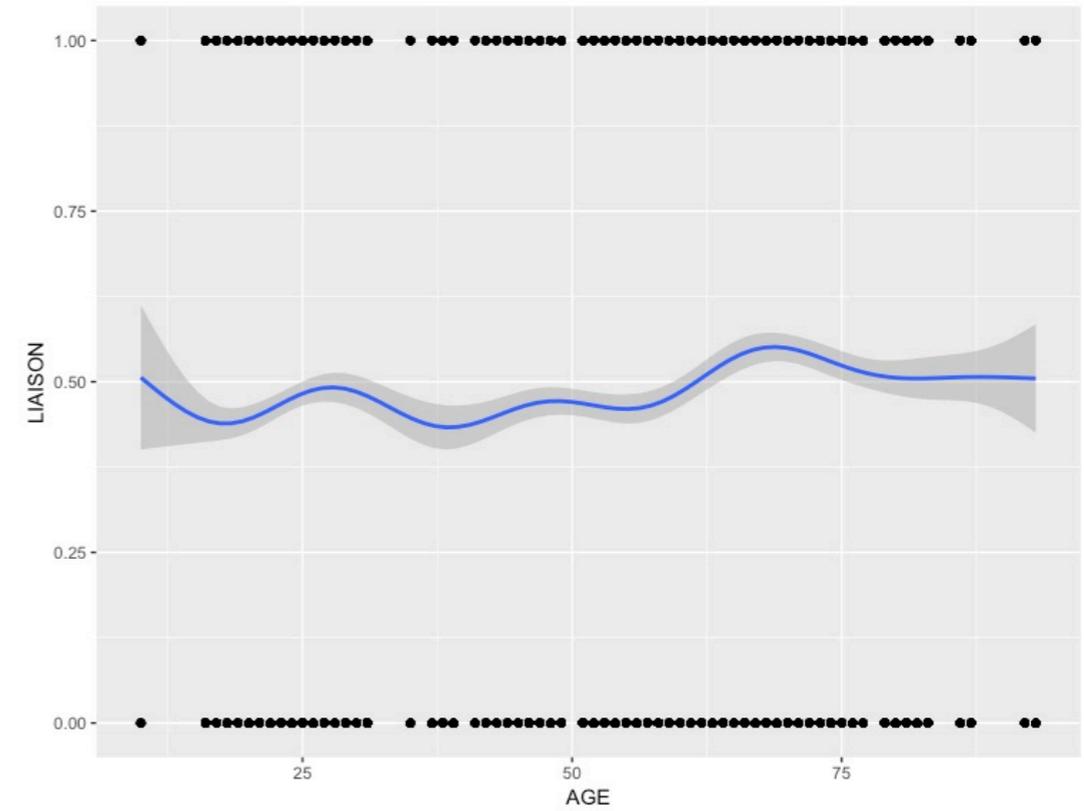
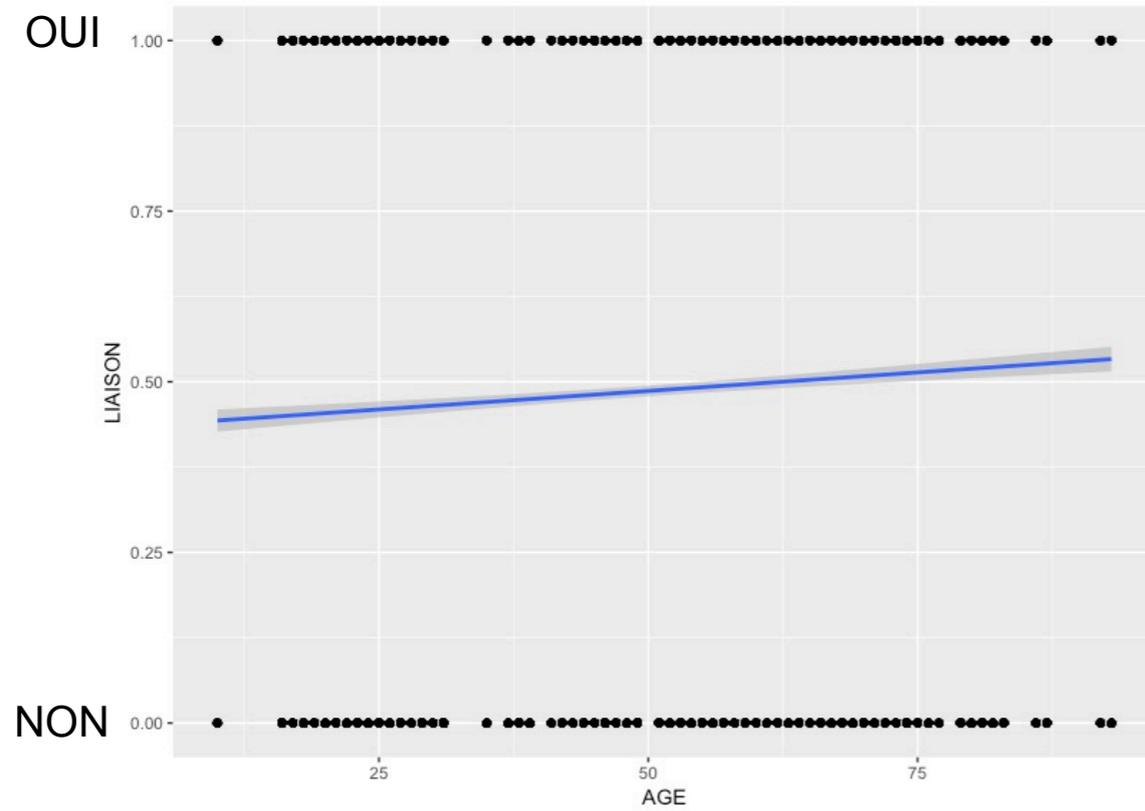
$$\text{mod.}AGE_{RL} = \beta_0 + \beta_1 AGE$$



$$\text{mod.}AGE_{GAM} = \beta_0 + s(AGE)$$

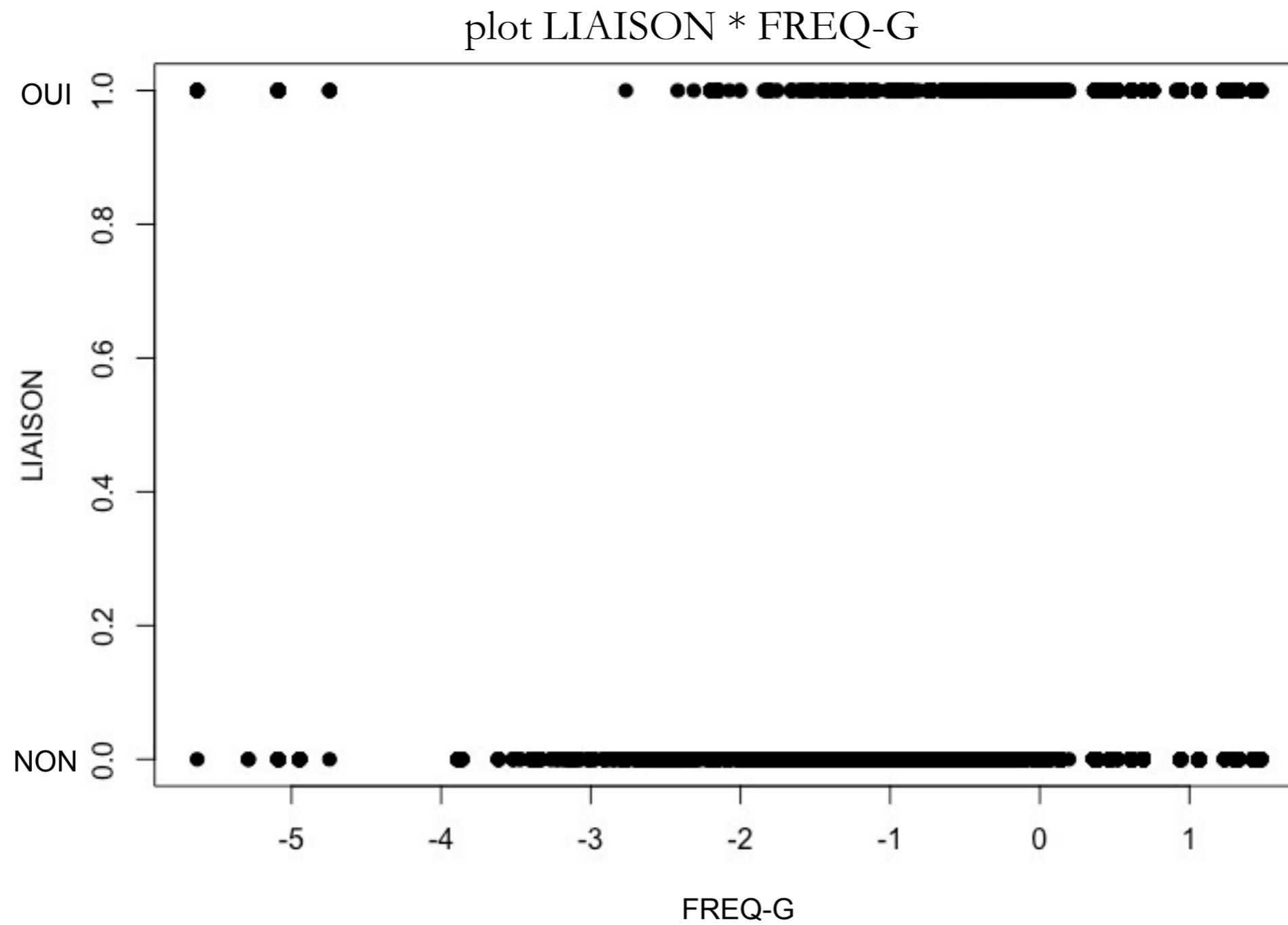


Comparaison des deux modèles

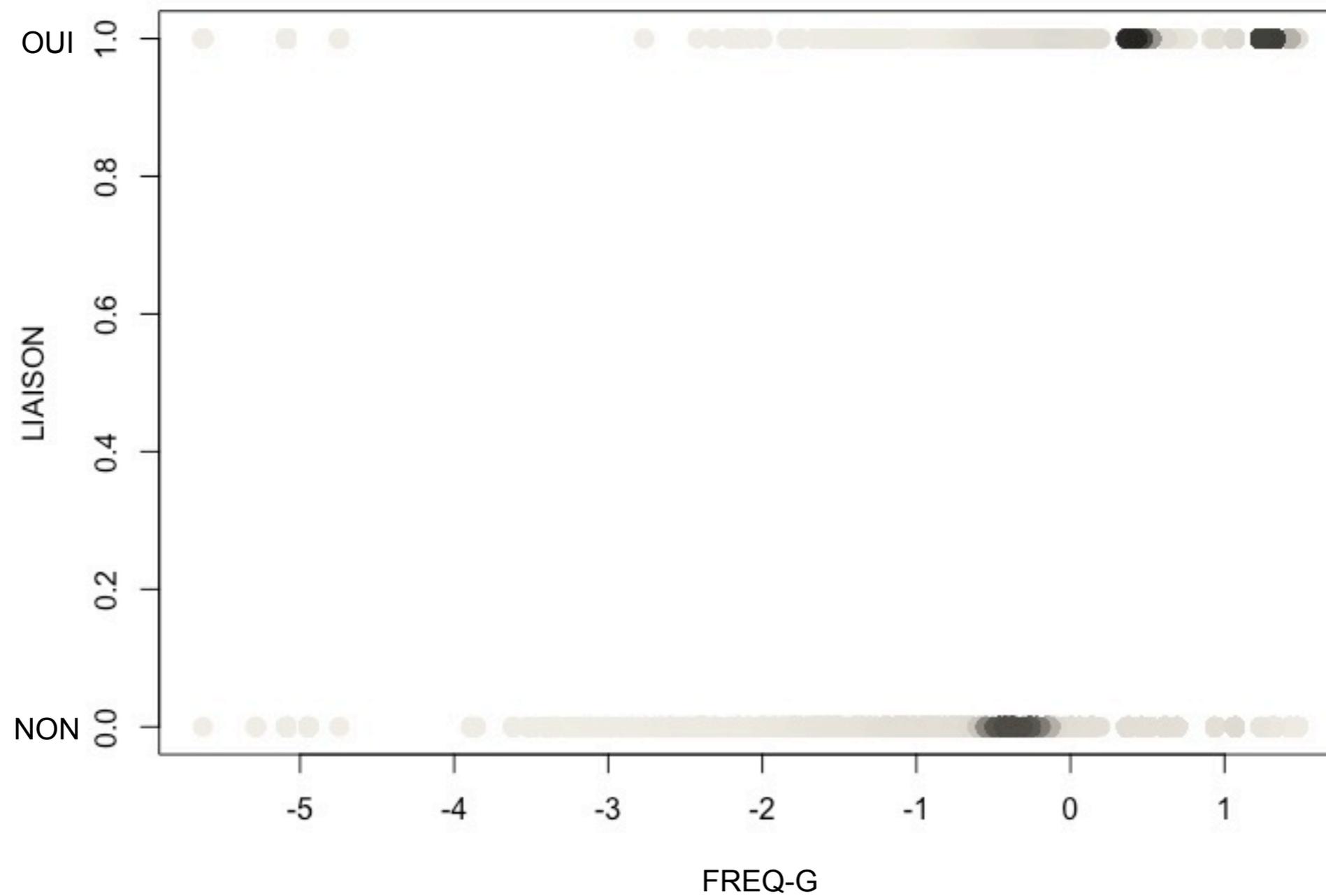


MODEL	AIC	%dev. expl.
mod.AGE _{LR}	21947.78	0.15
mod.AGE _{GAM}	21931.02	0.29

4.5 Modèle additif généralisé (GAM, Wood 2016)

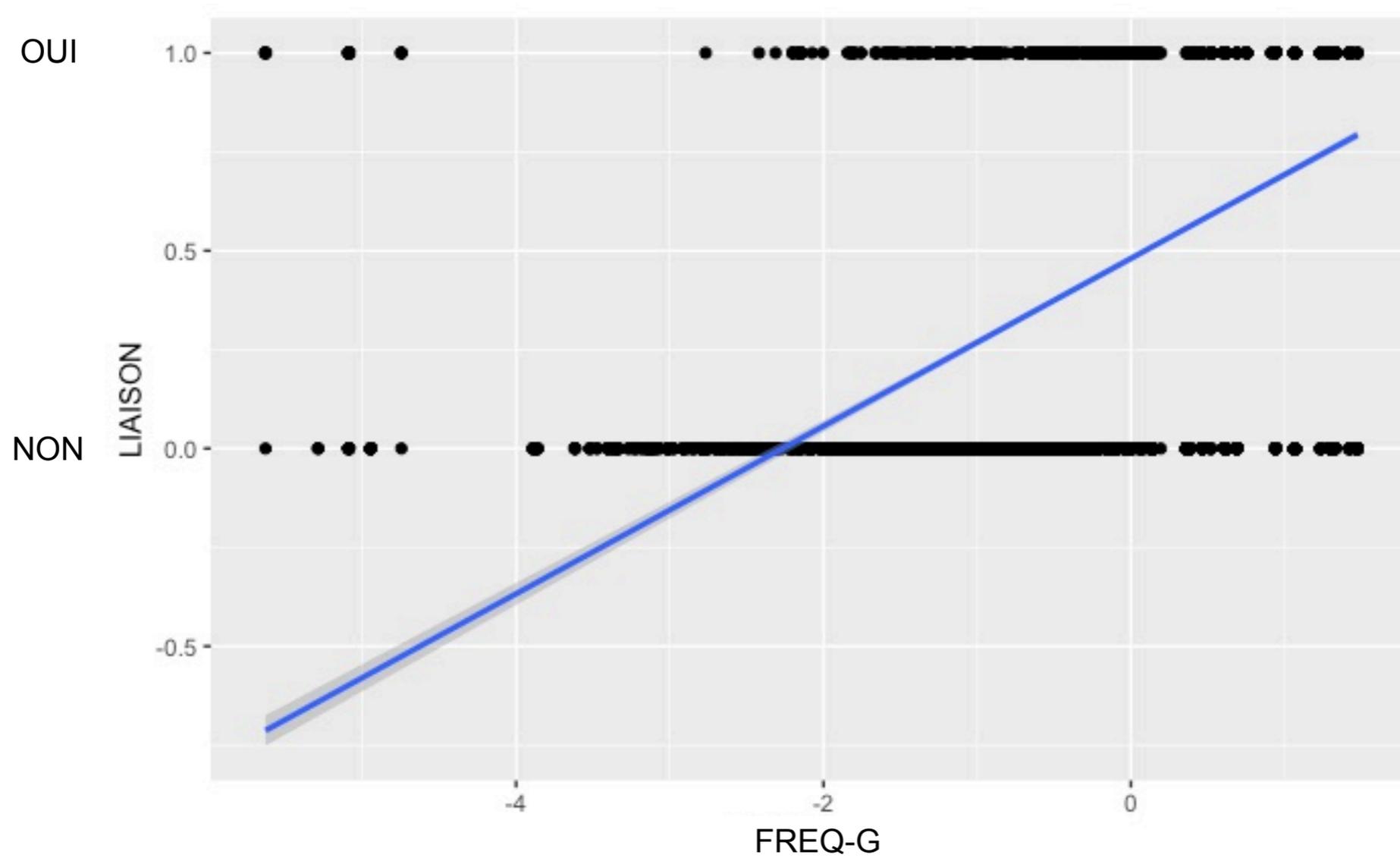


Visualisation de la densité
plot LIAISON * FREQ-G



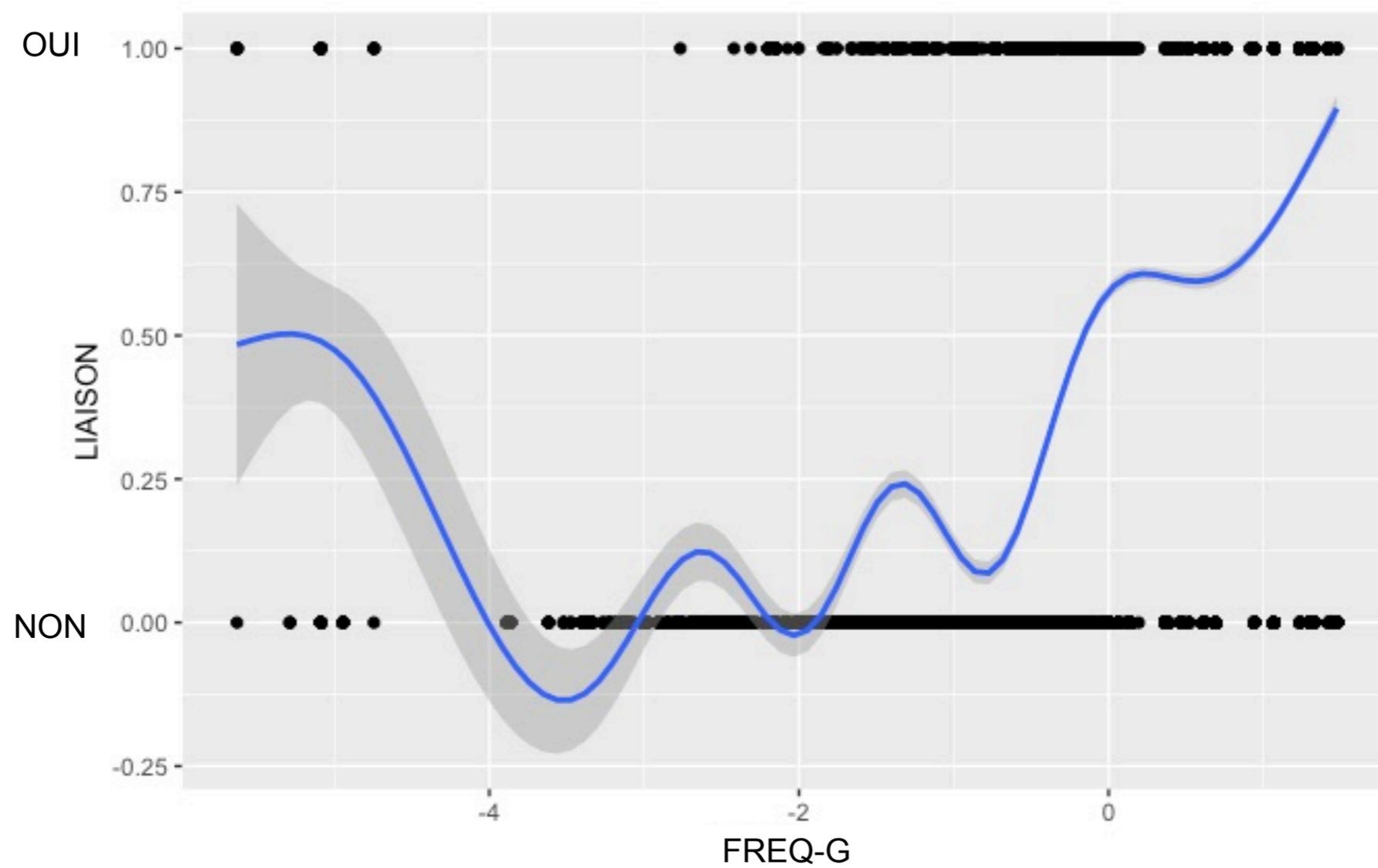
4.5 Modèle additif généralisé (GAM, Wood 2016)

$$\text{mod.}FREQ.G_{RL} = \beta_0 + \beta_1 FREQ.G$$



4.5 Modèle additif généralisé (GAM, Wood 2016)

$$\text{mod.FREQ.G}_{\text{GAM}} = \beta_0 + s(\text{FREQ.G})$$



Une fois bien équipé avec ces outils statistiques on essaie de créer un model de prédiction de la liaison

Une fois bien équipé avec ces outils statistiques on essaie de créer un modèle de prédiction de la liaison

Calcul des variables outre que les variables présentes dans PFC, comme par exemple la fréquence des mots à gauche et à droite (information tirée par FrWaC) ou la longueur du mot phonologique (information tirée par la ressource PsychoGLAFF)

4.6 Les variables de notre modèle

LIAISON	LOCUTEUR	ÉDUCATION	SEXE	AGE	VILLE	MOT_GAUCHE	MOT_DROITE	POS_GAUCHE	POS_DROITE
OUI	21aml1	20	F	65	Dijon	temps	en	NOM	PRP
NON	44ajs2	14	F	59	Nantes	est	en	VER	PRP
...
oui = 8350 no = 9036 TOT = 17386	niveaux = 192	ÉDUCATION = 15,41	F = 9100 H = 8286	AGE = 48,11	niveaux = 20	niveaux = 487	niveaux = 487	niveaux = 11	niveaux = 7

L-SYLL_G	L-SYLL_D	L-PHONO_G	L-PHONO_D	FREQ_G	FREQ_D	F.BIGRAM	TP.BIGRAM	LATITUDE	LONGITUDE
1	2	1	3	4233	7655	433	0.09	47.32	5.04
2	3	2	2	190	651	115	0.13	47.21	1.55
...
SYLL_G = 1,25	SYLL_D = 1,84	PHONO_G = 1,11	PHONO_D = 2,31	FREQ_G = 3213	FREQ_D = 4390	F.BIGRAM = 3691	TP.BIGRAM = 0.02		

4.6 Les variables de notre modèle

LIAISON	LOCUTEUR	ÉDUCATION	SEXE	AGE	VILLE	MOT_GAUCHE	MOT_DROITE	POS_GAUCHE	POS_DROITE
OUI	21aml1	20	F	65	Dijon	temps	en	NOM	PRP
NON	44ajs2	14	F	59	Nantes	est	en	VER	PRP
...
oui = 8350 no = 9036 TOT = 17386	niveaux = 192	ÉDUCATION = 15,41	F = 9100 H = 8286	AGE = 48,11	niveaux = 20	niveaux = 487	niveaux = 487	niveaux = 11	niveaux = 7

Toutes les variables quantitatives sont centrées réduites (c'est-à-dire : moyenne = 0 et écart-type = 1)

L-SYLL_G	L-SYLL_D	L-PHONO_G	L-PHONO_D	FREQ_G	FREQ_D	F.BIGRAM	TP.BIGRAM	LATITUDE	LONGITUDE
1	2	1	3	4233	7655	433	0.09	47.32	5.04
2	3	2	2	190	651	115	0.13	47.21	1.55
...
SYLL_G = 1,25	SYLL_D = 1,84	PHONO_G = 1,11	PHONO_D = 2,31	FREQ_G = 3213	FREQ_D = 4390	F.BIGRAM = 3691	TP.BIGRAM = 0.02		

4.6 Les variables de notre modèle

LIAISON	LOCUTEUR	ÉDUCATION	SEXE	AGE	VILLE	MOT_GAUCHE	MOT_DROITE	POS_GAUCHE	POS_DROITE
OUI	21aml1	20	F	65	Dijon	temps	en	NOM	PRP
NON	44ajs2	14	F	59	Nantes	est	en	VER	PRP
...
oui = 8350 no = 9036 TOT = 17386	niveaux = 192	ÉDUCATION = 15,41	F = 9100 H = 8286	AGE = 48,11	niveaux = 20	niveaux = 487	niveaux = 487	niveaux = 11	niveaux = 7

Toutes les variables quantitatives sont centrées réduites (c'est-à-dire : moyenne = 0 et écart-type = 1)

Toutes les variables *FREQ* sont centrées réduites en fonction *log*

L-SYLL_G	L-SYLL_D	L-PHONO_G	L-PHONO_D	FREQ_G	FREQ_D	F.BIGRAM	TP.BIGRAM	LATITUDE	LONGITUDE
1	2	1	3	4233	7655	433	0.09	47.32	5.04
2	3	2	2	190	651	115	0.13	47.21	1.55
...
SYLL_G = 1,25	SYLL_D = 1,84	PHONO_G = 1,11	PHONO_D = 2,31	FREQ_G = 3213	FREQ_D = 4390	F.BIGRAM = 3691	TP.BIGRAM = 0.02		

La méthode adoptée pour la sélection des variables définissant le modèle est la régression itérative (procédure *stepwise progressive*) qui inclut d'abord la variable qui propose la meilleur déviance expliquée.

Ensuite, celle qui améliore le plus la déviance et ainsi de suite.

La méthode adoptée pour la sélection des variables définissant le modèle est la régression itérative (procédure *stepwise progressive*) qui inclut d'abord la variable qui propose la meilleur déviance expliquée.

Ensuite, celle qui améliore le plus la déviance et ainsi de suite.

$$model1 = \beta_0 + s(MOT - G)$$

La méthode adoptée pour la sélection des variables définissant le modèle est la régression itérative (procédure *stepwise progressive*) qui inclut d'abord la variable qui propose la meilleur déviance expliquée.

Ensuite, celle qui améliore le plus la déviance et ainsi de suite.

$$model1 = \beta_0 + s(MOT - G)$$

$$model2 = \beta_0 + s(MOT - G) + s(FREQ - G)$$

La méthode adoptée pour la sélection des variables définissant le modèle est la régression itérative (procédure *stepwise progressive*) qui inclut d'abord la variable qui propose la meilleur déviance expliquée.

Ensuite, celle qui améliore le plus la déviance et ainsi de suite.

Création des modèles : calcul de AIC et (%) déviance expliquée

$$model1 = \beta_0 + s(MOT - G)$$

$$model2 = \beta_0 + s(MOT - G) + s(FREQ - G)$$

La méthode adoptée pour la sélection des variables définissant le modèle est la régression itérative (procédure *stepwise progressive*) qui inclut d'abord la variable qui propose la meilleur déviance expliquée.

Ensuite, celle qui améliore le plus la déviance et ainsi de suite.

Création des modèles : calcul de AIC et (%) déviance expliquée

$$model1 = \beta_0 + s(MOT - G)$$

$$model2 = \beta_0 + s(MOT - G) + s(FREQ - G)$$

$$model3 = \beta_0 + s(MOT - G) + s(FREQ - G) + s(POS - D)$$

$$model4 = \beta_0 + s(MOT - G) + s(FREQ - G) + s(POS - D) + s(TP.BIGRAM)$$

...

4.8 Puissance explicative des variables

Predictors	Type	%dev. expl.	AIC	t-test
locuteur	factor (192)	3.72	23479.67	$p < 0.001$ * **
age	numer.	0.29	21931.02	$p < 0.001$ * **
éducation	numer.	0.63	15388.94	$p < 0.001$ * **
sexe	factor (2)	0.02	24074.83	$p < 0.001$ * **
mot-gauche	factor (487)	70.80	7699.546	$p < 0.001$ * **
mot-droite	factor (1531)	48.41	14551.51	$p < 0.001$ * **
pos-gauche	factor (11)	49.91	11933.72	$p < 0.001$ * **
pos-droite	factor (7)	23.00	17260.84	$p < 0.001$ * **
long.syll-g	numer.	16.50	20081.17	$p < 0.001$ * **
long.syll-d	numer.	9.30	22932.11	$p < 0.001$ * **
long.phono-g	numer.	15.43	21341.09	$p < 0.001$ * **
long.phono-d	numer.	7.52	22231.40	$p < 0.001$ * **
freq-g	numer.	18.41	19662.86	$p < 0.001$ * **
freq-d	numer.	1.32	26731.56	$p < 0.001$ * **
f.bigram	numer.	7.97	20633.03	$p < 0.001$ * **
tp.bigram	numer.	5.20	21252.77	$p < 0.001$ * **
(longitude,latitude)	numer.	1.72	23923.42	$p < 0.001$ * **

4.8 Puissance explicative des variables

Predictors	Type	%dev. expl.	AIC	t-test
locuteur	factor (192)	3.72	23479.67	$p < 0.001$ * **
age	numer.	0.29	21931.02	$p < 0.001$ * **
éducation	numer.	0.63	15388.94	$p < 0.001$ * **
sexe	factor (2)	0.02	24074.83	$p < 0.001$ * **
mot-gauche	factor (487)	70.80	7699.546	$p < 0.001$ * **
mot-droite	factor (1531)	48.41	14551.51	$p < 0.001$ * **
pos-gauche	factor (11)	49.91	11933.72	$p < 0.001$ * **
pos-droite	factor (7)	23.00	17260.84	$p < 0.001$ * **
long.syll-g	numer.	16.50	20081.17	$p < 0.001$ * **
long.syll-d	numer.	9.30	22932.11	$p < 0.001$ * **
long.phono-g	numer.	15.43	21341.09	$p < 0.001$ * **
long.phono-d	numer.	7.52	22231.40	$p < 0.001$ * **
freq-g	numer.	18.41	19662.86	$p < 0.001$ * **
freq-d	numer.	1.32	26731.56	$p < 0.001$ * **
f.bigram	numer.	7.97	20633.03	$p < 0.001$ * **
tp.bigram	numer.	5.20	21252.77	$p < 0.001$ * **
(longitude,latitude)	numer.	1.72	23923.42	$p < 0.001$ * **

- Objectif du travail
- La liaison - une petite présentation
- Un corpus phonologique : PFC
- Les outils statistiques (RL et GAM)
- Le model final et le résultats
- Limites du travail et perspectives futures

6. Le modèle final

$$\begin{aligned} model_{\text{FINAL}} = & -1.915713 + s(\text{AGE}) + s(\text{FREQ}-G) + s(\text{SYLL}-G) + s(\text{TP.BIGRAM}) + s(\text{MOT}-G) + \\ & + s(\text{POS}-G) + s(\text{POS}-D) + s(\text{LOCUTEUR}) + s(\text{LONGITUDE}, \text{LATITUDE}) \end{aligned}$$

6. Le modèle final

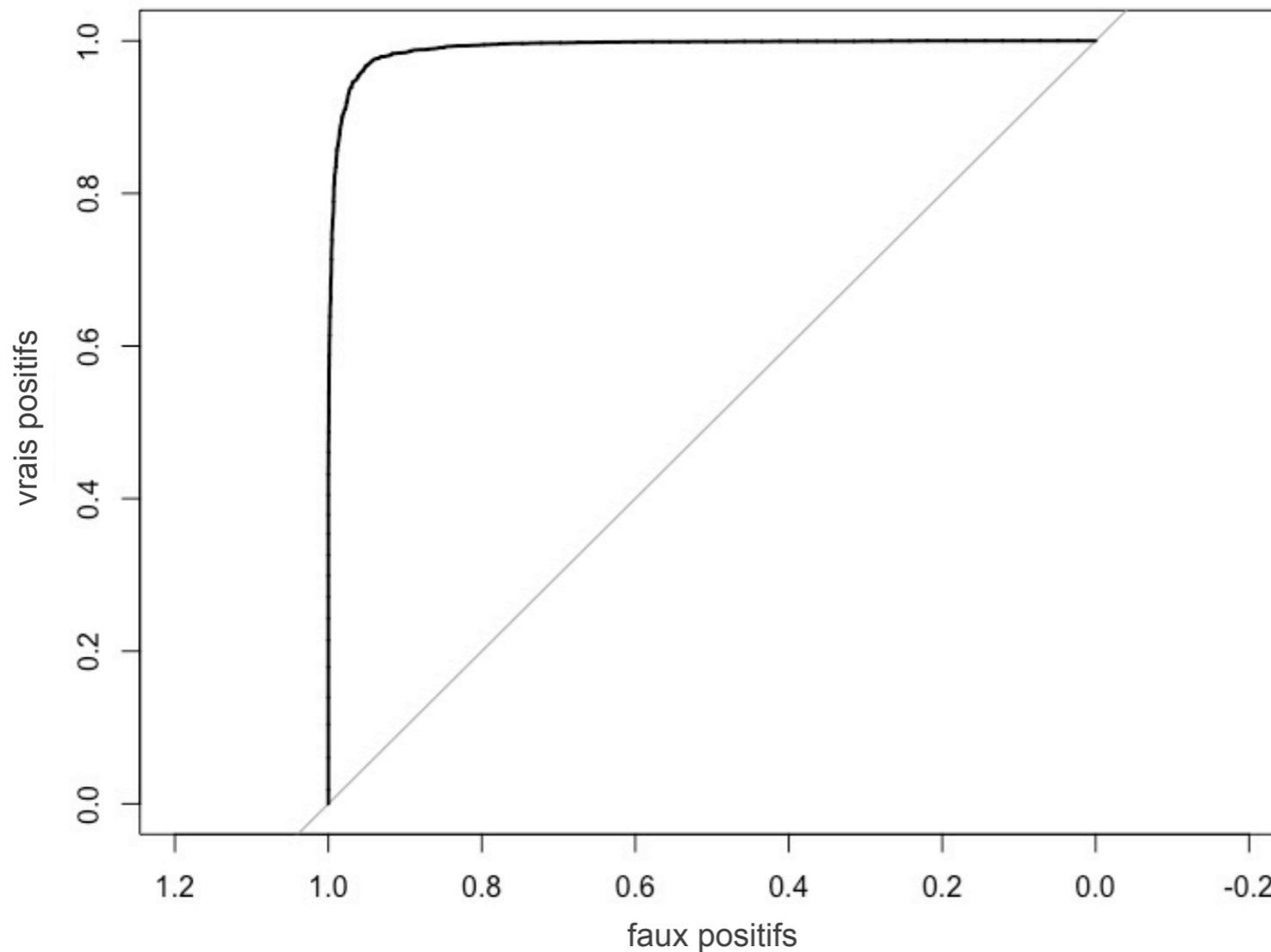
$$\begin{aligned} model_{\text{FINAL}} = & -1.915713 + s(\text{AGE}) + s(\text{FREQ}-G) + s(\text{SYLL}-G) + s(\text{TP.BIGRAM}) + s(\text{MOT}-G) + \\ & + s(\text{POS}-G) + s(\text{POS}-D) + s(\text{LOCUTEUR}) + s(\text{LONGITUDE}, \text{LATITUDE}) \end{aligned}$$

AIC	2999.090
% dev. expl.	83.214
AUC	0.991

6. Le modèle final

$$\begin{aligned} model_{FINAL} = & -1.915713 + s(AGE) + s(FREQ-G) + s(SYLL-G) + s(TP.BIGRAM) + s(MOT-G) + \\ & + s(POS-G) + s(POS-D) + s(LOCUTEUR) + s(LONGITUDE, LATITUDE) \end{aligned}$$

AIC	2999.090
% dev. expl.	83.214
AUC	0.991



6.1 Le modèle final - Les fonctions de lissage

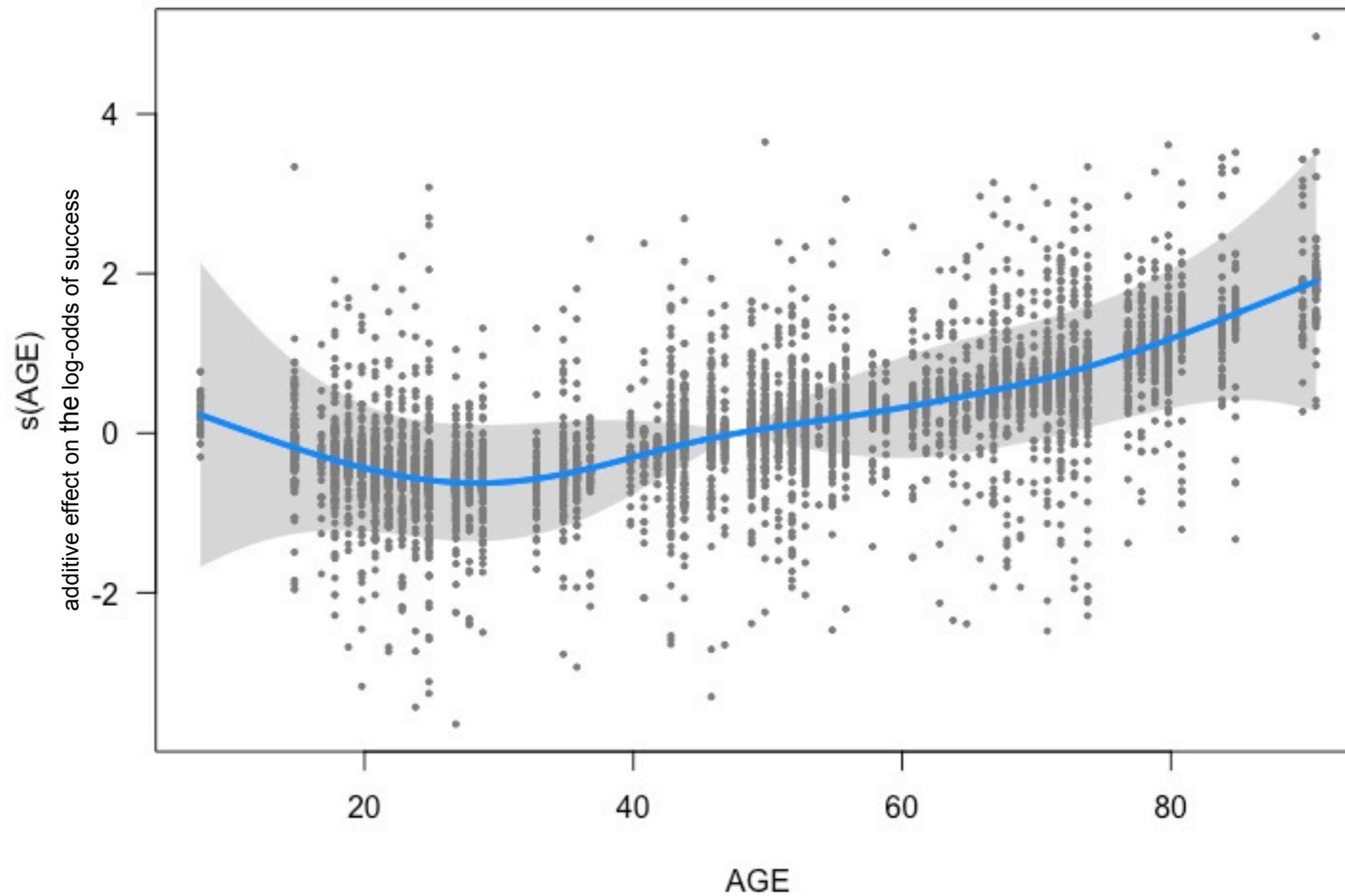
$$\begin{aligned} model_{\text{FINAL}} = & -1.915713 + s(\text{AGE}) + s(\text{FREQ}-G) + s(\text{SYLL}-G) + s(\text{TP.BIGRAM}) + s(\text{MOT}-G) + \\ & + s(\text{POS}-G) + s(\text{POS}-D) + s(\text{LOCUTEUR}) + s(\text{LONGITUDE}, \text{LATITUDE}) \end{aligned}$$

6.1 Le modèle final - Les fonctions de lissage

$$\begin{aligned} model_{\text{FINAL}} = & -1.915713 + s(\text{AGE}) + s(\text{FREQ}-G) + s(\text{SYLL}-G) + s(\text{TP.BIGRAM}) + s(\text{MOT}-G) + \\ & + s(\text{POS}-G) + s(\text{POS}-D) + s(\text{LOCUTEUR}) + s(\text{LONGITUDE}, \text{LATITUDE}) \end{aligned}$$

6.1 Le modèle final - Les fonctions de lissage

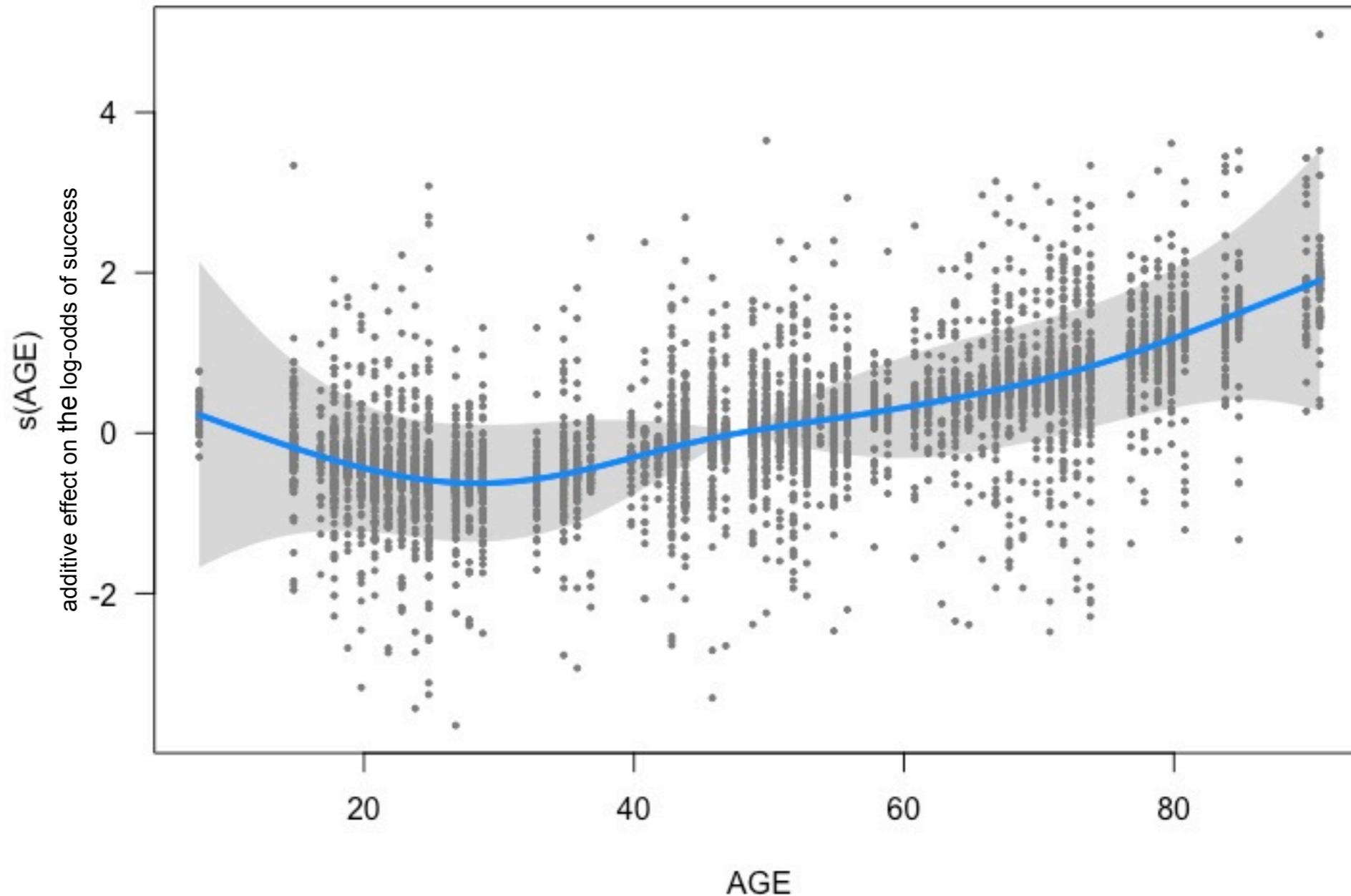
$$\begin{aligned} model_{\text{FINAL}} = & -1.915713 + s(\text{AGE}) + s(\text{FREQ}-G) + s(\text{SYLL}-G) + s(\text{TP.BIGRAM}) + s(\text{MOT}-G) + \\ & + s(\text{POS}-G) + s(\text{POS}-D) + s(\text{LOCUTEUR}) + s(\text{LONGITUDE}, \text{LATITUDE}) \end{aligned}$$



6.1 Le modèle final - Les fonctions de lissage

$$\begin{aligned} model_{FINAL} = & -1.915713 + s(AGE) + s(FREQ-G) + s(SYLL-G) + s(TP.BIGRAM) + s(MOT-G) + \\ & + s(POS-G) + s(POS-D) + s(LOCUTEUR) + s(LONGITUDE, LATITUDE) \end{aligned}$$

L'augmentation de l'AGE détermine une plus forte probabilité d'avoir une liaison réalisée



6.1 Le modèle final - Les fonctions de lissage

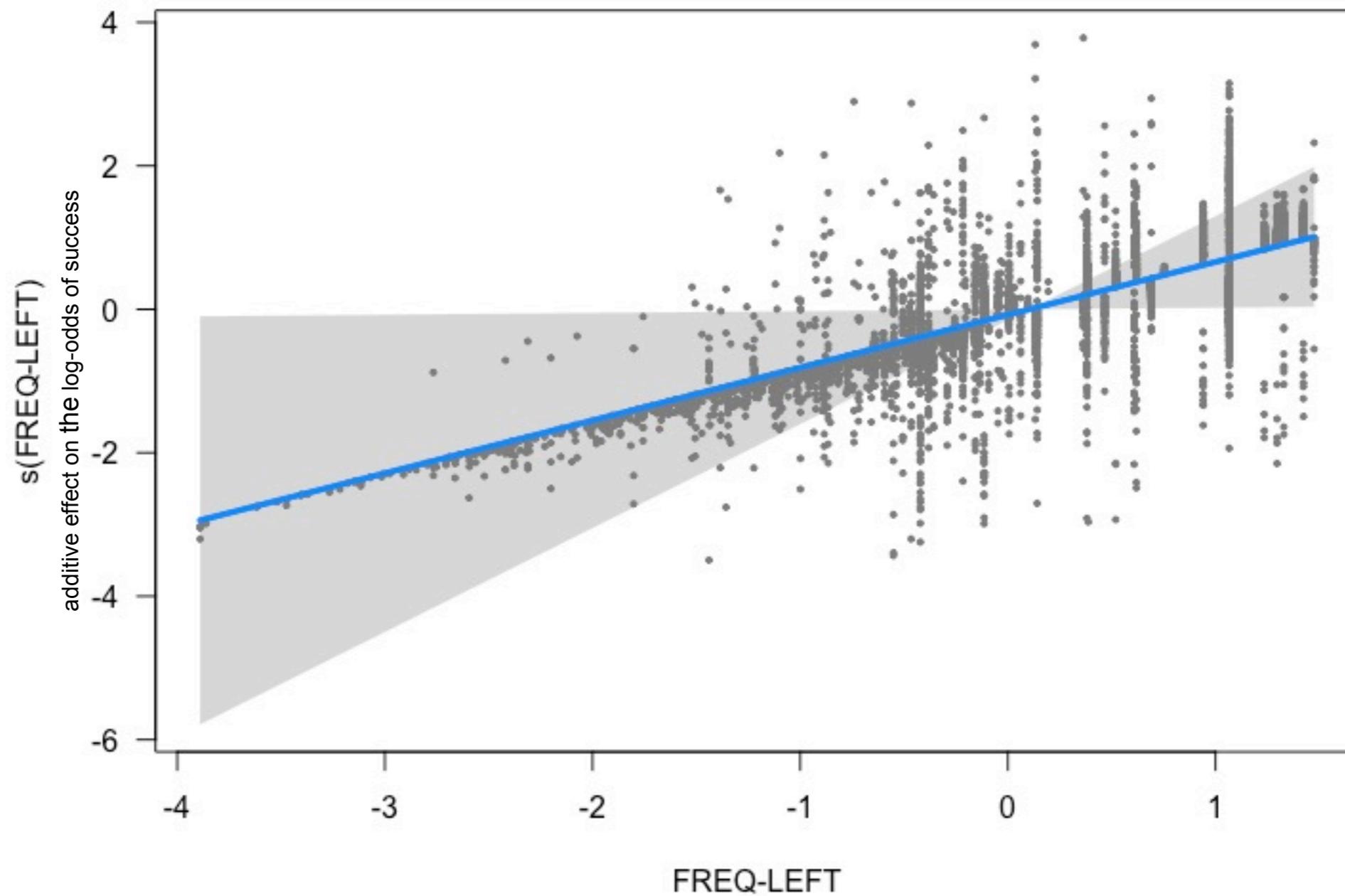
$$\begin{aligned} model_{\text{FINAL}} = & -1.915713 + s(\text{AGE}) + s(\text{FREQ}-G) + s(\text{SYLL}-G) + s(\text{TP.BIGRAM}) + s(\text{MOT}-G) + \\ & + s(\text{POS}-G) + s(\text{POS}-D) + s(\text{LOCUTEUR}) + s(\text{LONGITUDE}, \text{LATITUDE}) \end{aligned}$$

6.1 Le modèle final - Les fonctions de lissage

$$\begin{aligned} model_{\text{FINAL}} = & -1.915713 + s(\text{AGE}) + s(\text{FREQ-G}) + s(\text{SYLL-G}) + s(\text{TP.BIGRAM}) + s(\text{MOT-G}) + \\ & + s(\text{POS-G}) + s(\text{POS-D}) + s(\text{LOCUTEUR}) + s(\text{LONGITUDE, LATITUDE}) \end{aligned}$$

6.1 Le modèle final - Les fonctions de lissage

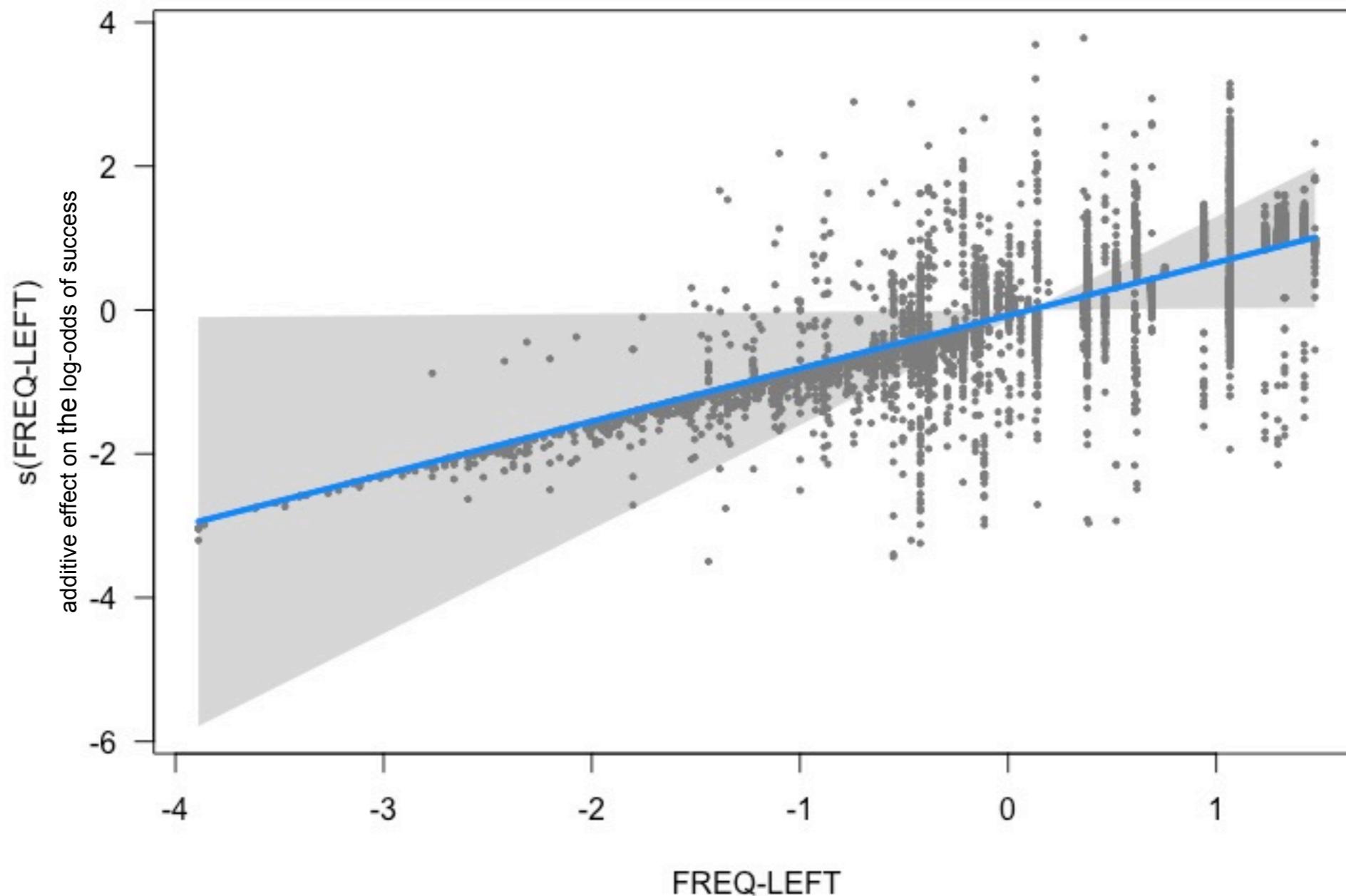
$$\begin{aligned} model_{FINAL} = & -1.915713 + s(AGE) + s(FREQ-G) + s(SYLL-G) + s(TP.BIGRAM) + s(MOT-G) + \\ & + s(POS-G) + s(POS-D) + s(LOCUTEUR) + s(LONGITUDE, LATITUDE) \end{aligned}$$



6.1 Le modèle final - Les fonctions de lissage

$$\begin{aligned} model_{FINAL} = & -1.915713 + s(AGE) + s(FREQ-G) + s(SYLL-G) + s(TP.BIGRAM) + s(MOT-G) + \\ & + s(POS-G) + s(POS-D) + s(LOCUTEUR) + s(LONGITUDE, LATITUDE) \end{aligned}$$

L'augmentation de FREQ-G détermine une plus forte probabilité d'avoir une liaison réalisée



6.1 Le modèle final - Les fonctions de lissage

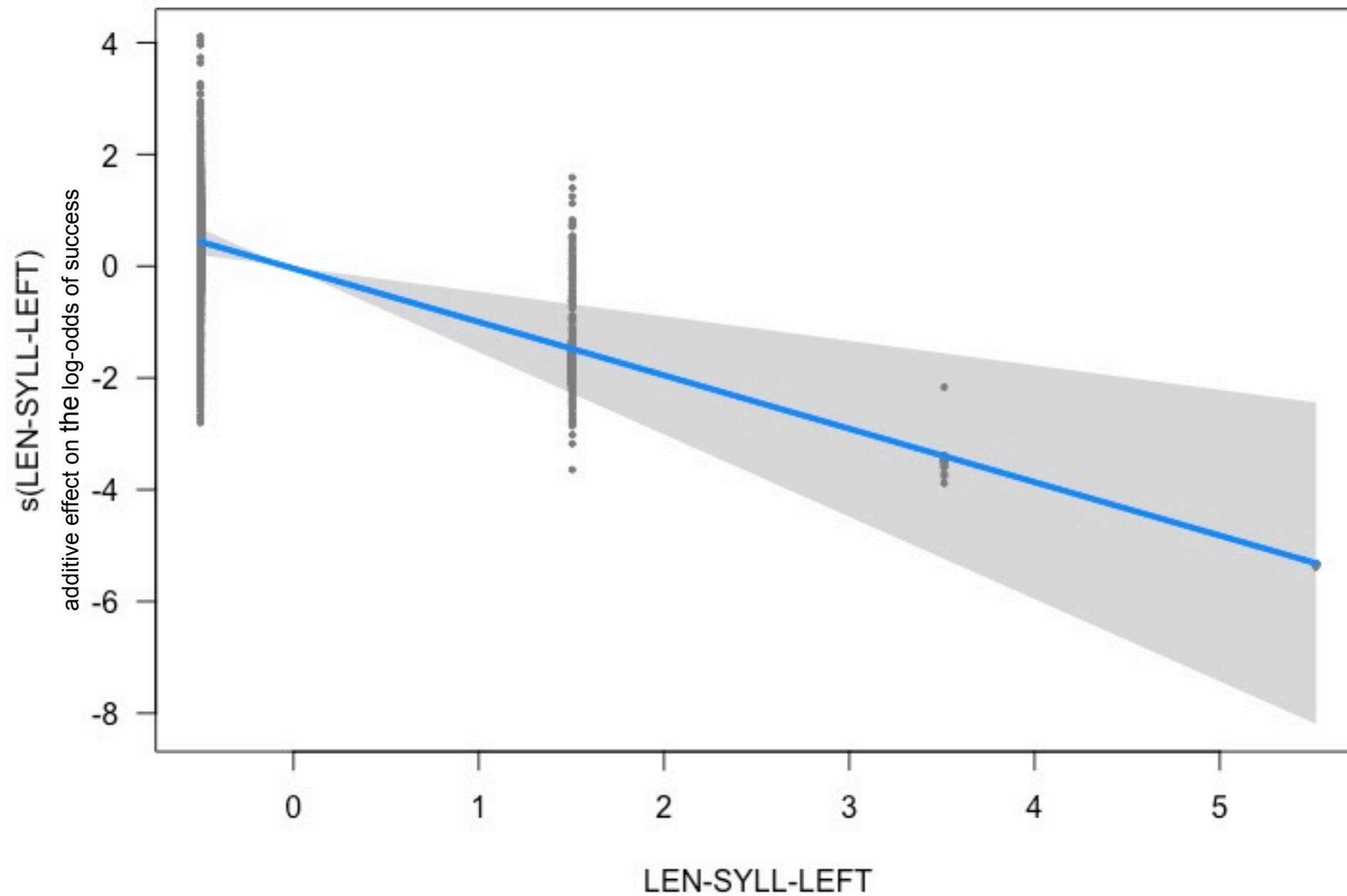
$$\begin{aligned} model_{\text{FINAL}} = & -1.915713 + s(\text{AGE}) + s(\text{FREQ}-G) + s(\text{SYLL}-G) + s(\text{TP.BIGRAM}) + s(\text{MOT}-G) + \\ & + s(\text{POS}-G) + s(\text{POS}-D) + s(\text{LOCUTEUR}) + s(\text{LONGITUDE}, \text{LATITUDE}) \end{aligned}$$

6.1 Le modèle final - Les fonctions de lissage

$$\begin{aligned} model_{\text{FINAL}} = & -1.915713 + s(\text{AGE}) + s(\text{FREQ}-G) + s(\text{SYLL}-G) + s(\text{TP.BIGRAM}) + s(\text{MOT}-G) + \\ & + s(\text{POS}-G) + s(\text{POS}-D) + s(\text{LOCUTEUR}) + s(\text{LONGITUDE}, \text{LATITUDE}) \end{aligned}$$

6.1 Le modèle final - Les fonctions de lissage

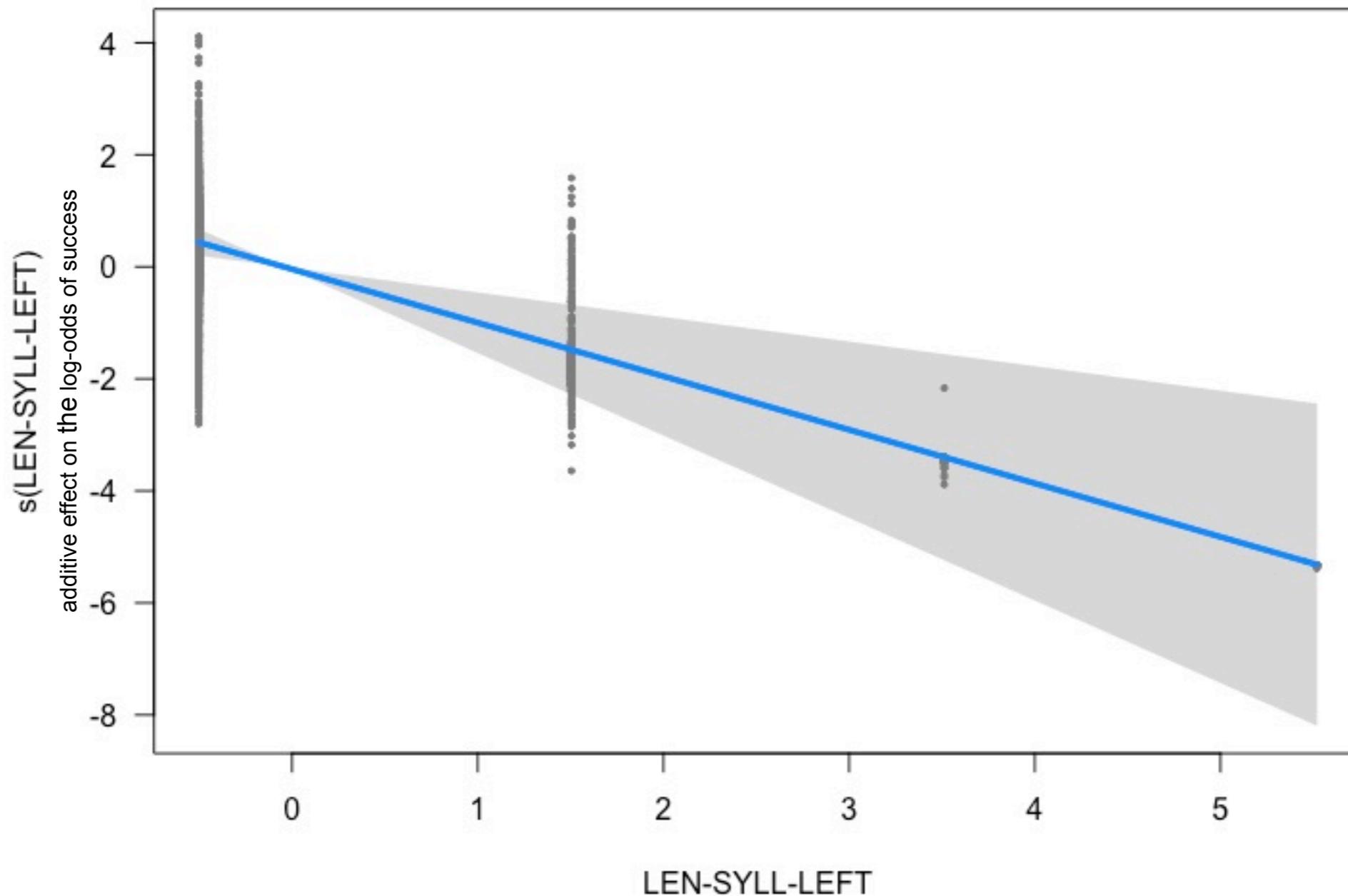
$$\begin{aligned} model_{FINAL} = & -1.915713 + s(AGE) + s(FREQ-G) + s(SYLL-G) + s(TP.BIGRAM) + s(MOT-G) + \\ & + s(POS-G) + s(POS-D) + s(LOCUTEUR) + s(LONGITUDE, LATITUDE) \end{aligned}$$



6.1 Le modèle final - Les fonctions de lissage

$$\begin{aligned} model_{FINAL} = & -1.915713 + s(AGE) + s(FREQ-G) + s(SYLL-G) + s(TP.BIGRAM) + s(MOT-G) + \\ & + s(POS-G) + s(POS-D) + s(LOCUTEUR) + s(LONGITUDE, LATITUDE) \end{aligned}$$

L'augmentation de SYLL-G détermine une probabilité plus faible d'avoir une liaison réalisée



6.1 Le modèle final - Les fonctions de lissage

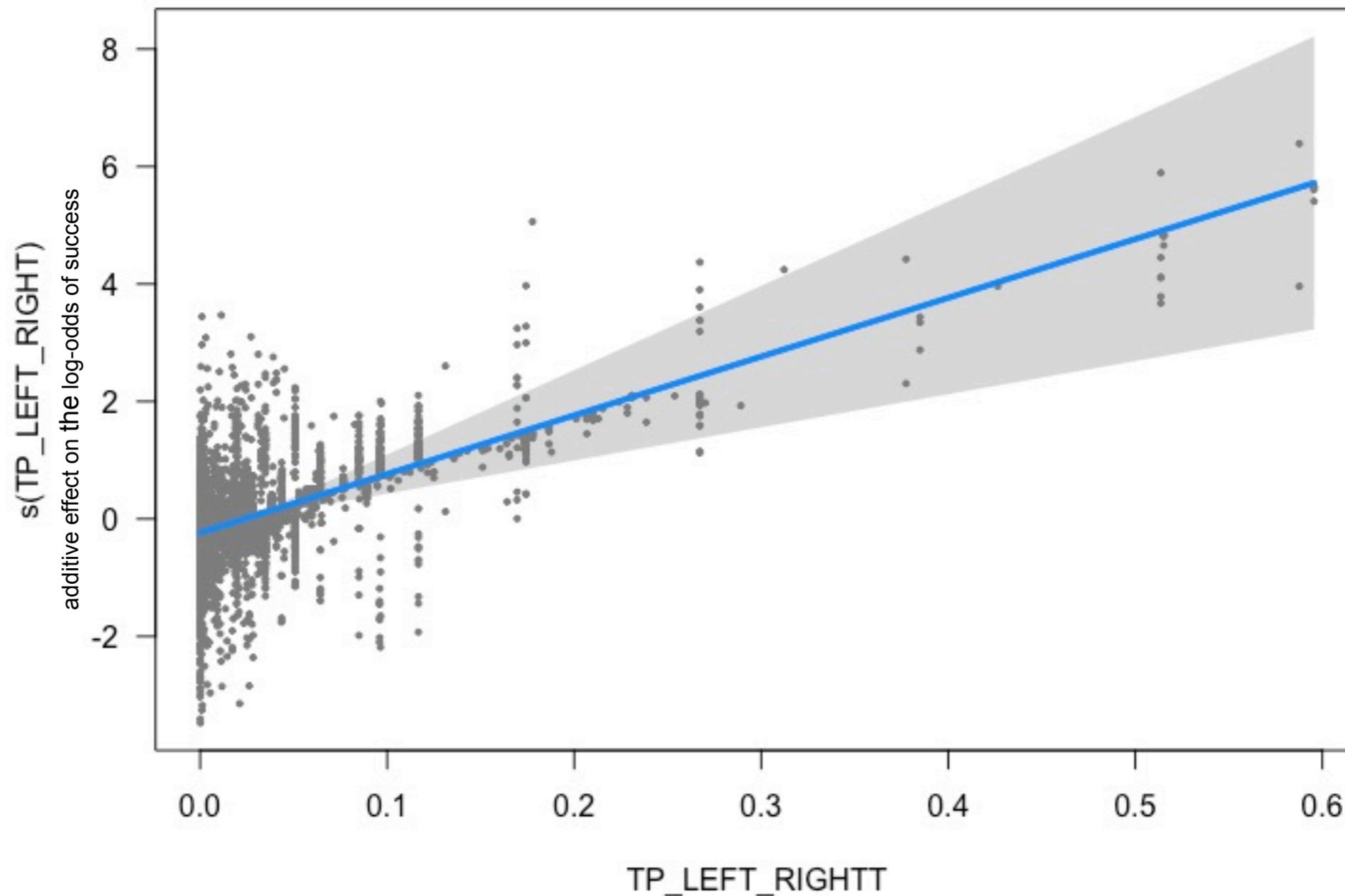
$$\begin{aligned} model_{\text{FINAL}} = & -1.915713 + s(\text{AGE}) + s(\text{FREQ}-G) + s(\text{SYLL}-G) + s(\text{TP.BIGRAM}) + s(\text{MOT}-G) + \\ & + s(\text{POS}-G) + s(\text{POS}-D) + s(\text{LOCUTEUR}) + s(\text{LONGITUDE}, \text{LATITUDE}) \end{aligned}$$

6.1 Le modèle final - Les fonctions de lissage

$$\begin{aligned} model_{\text{FINAL}} = & -1.915713 + s(\text{AGE}) + s(\text{FREQ}-G) + s(\text{SYLL}-G) + s(\text{TP.BIGRAM}) + s(\text{MOT}-G) + \\ & + s(\text{POS}-G) + s(\text{POS}-D) + s(\text{LOCUTEUR}) + s(\text{LONGITUDE}, \text{LATITUDE}) \end{aligned}$$

6.1 Le modèle final - Les fonctions de lissage

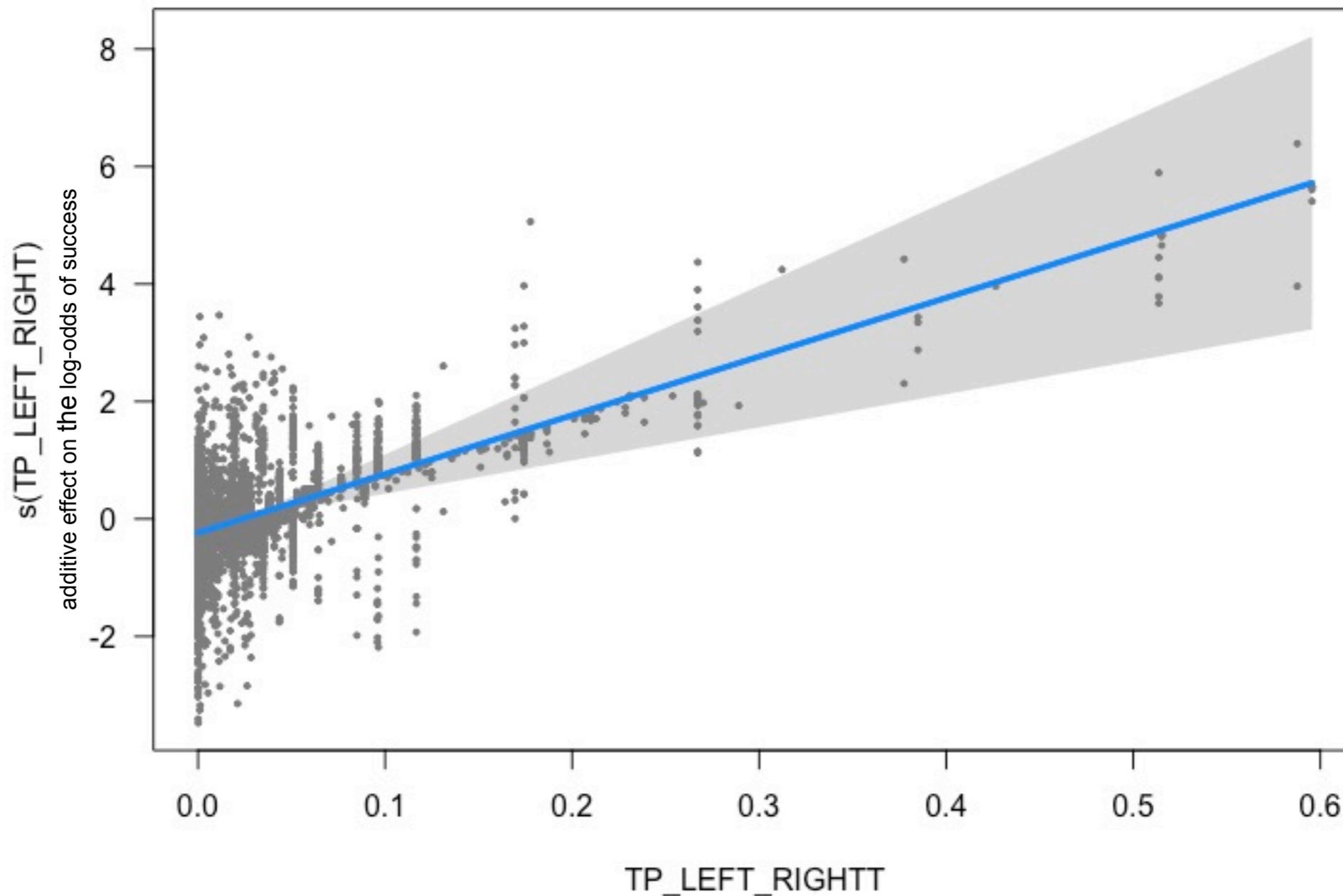
$$\begin{aligned} model_{FINAL} = & -1.915713 + s(AGE) + s(FREQ-G) + s(SYLL-G) + s(TP.BIGRAM) + s(MOT-G) + \\ & + s(POS-G) + s(POS-D) + s(LOCUTEUR) + s(LONGITUDE, LATITUDE) \end{aligned}$$



6.1 Le modèle final - Les fonctions de lissage

$$\begin{aligned} model_{FINAL} = & -1.915713 + s(AGE) + s(FREQ-G) + s(SYLL-G) + s(TP.BIGRAM) + s(MOT-G) + \\ & + s(POS-G) + s(POS-D) + s(LOCUTEUR) + s(LONGITUDE, LATITUDE) \end{aligned}$$

L'augmentation de TP.BIGRAM détermine une plus forte probabilité d'avoir une liaison réalisée



6.1 Le modèle final - Les fonctions de lissage

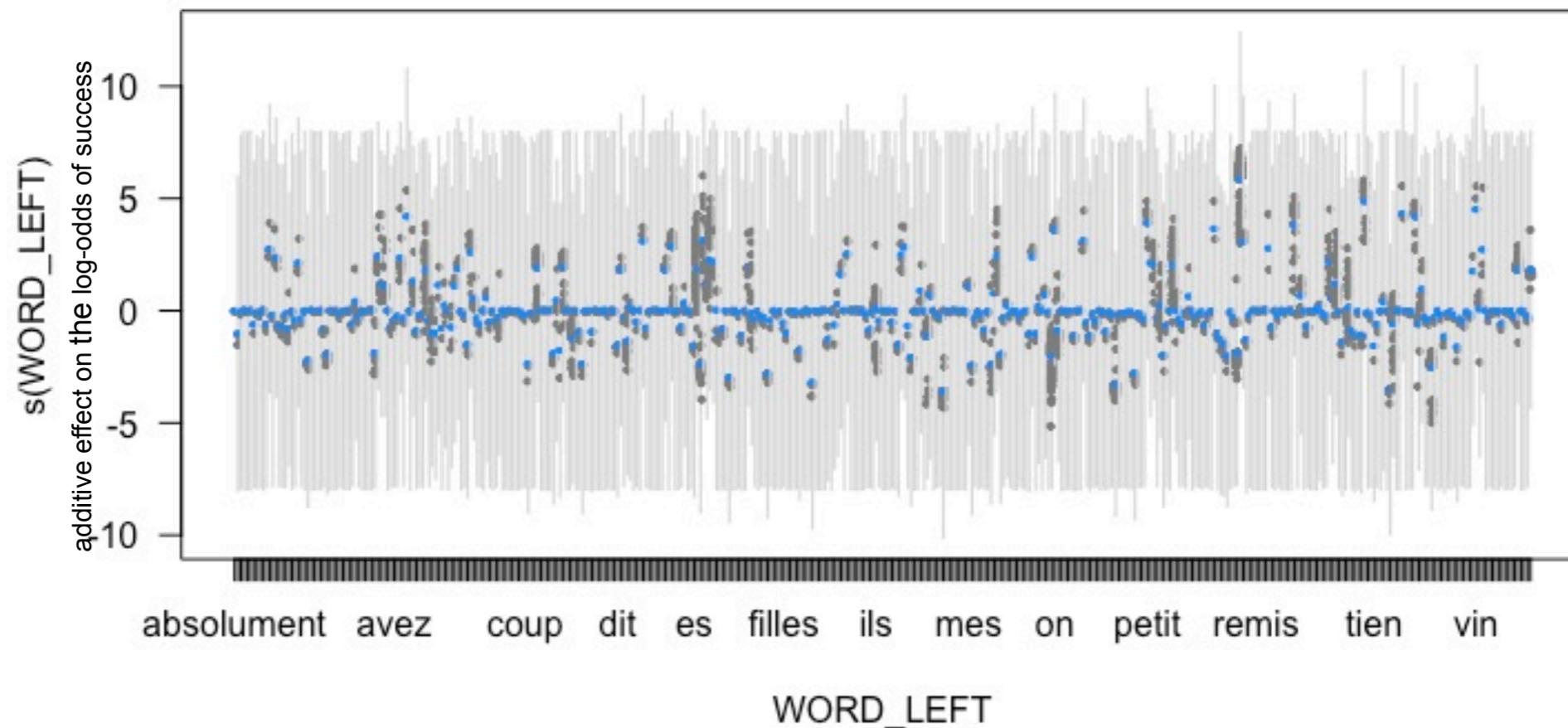
$$\begin{aligned} model_{\text{FINAL}} = & -1.915713 + s(\text{AGE}) + s(\text{FREQ}-G) + s(\text{SYLL}-G) + s(\text{TP.BIGRAM}) + s(\text{MOT}-G) + \\ & + s(\text{POS}-G) + s(\text{POS}-D) + s(\text{LOCUTEUR}) + s(\text{LONGITUDE}, \text{LATITUDE}) \end{aligned}$$

6.1 Le modèle final - Les fonctions de lissage

$$\begin{aligned} model_{\text{FINAL}} = & -1.915713 + s(\text{AGE}) + s(\text{FREQ}-G) + s(\text{SYLL}-G) + s(\text{TP.BIGRAM}) + s(\text{MOT}-G) + \\ & + s(\text{POS}-G) + s(\text{POS}-D) + s(\text{LOCUTEUR}) + s(\text{LONGITUDE}, \text{LATITUDE}) \end{aligned}$$

6.1 Le modèle final - Les fonctions de lissage

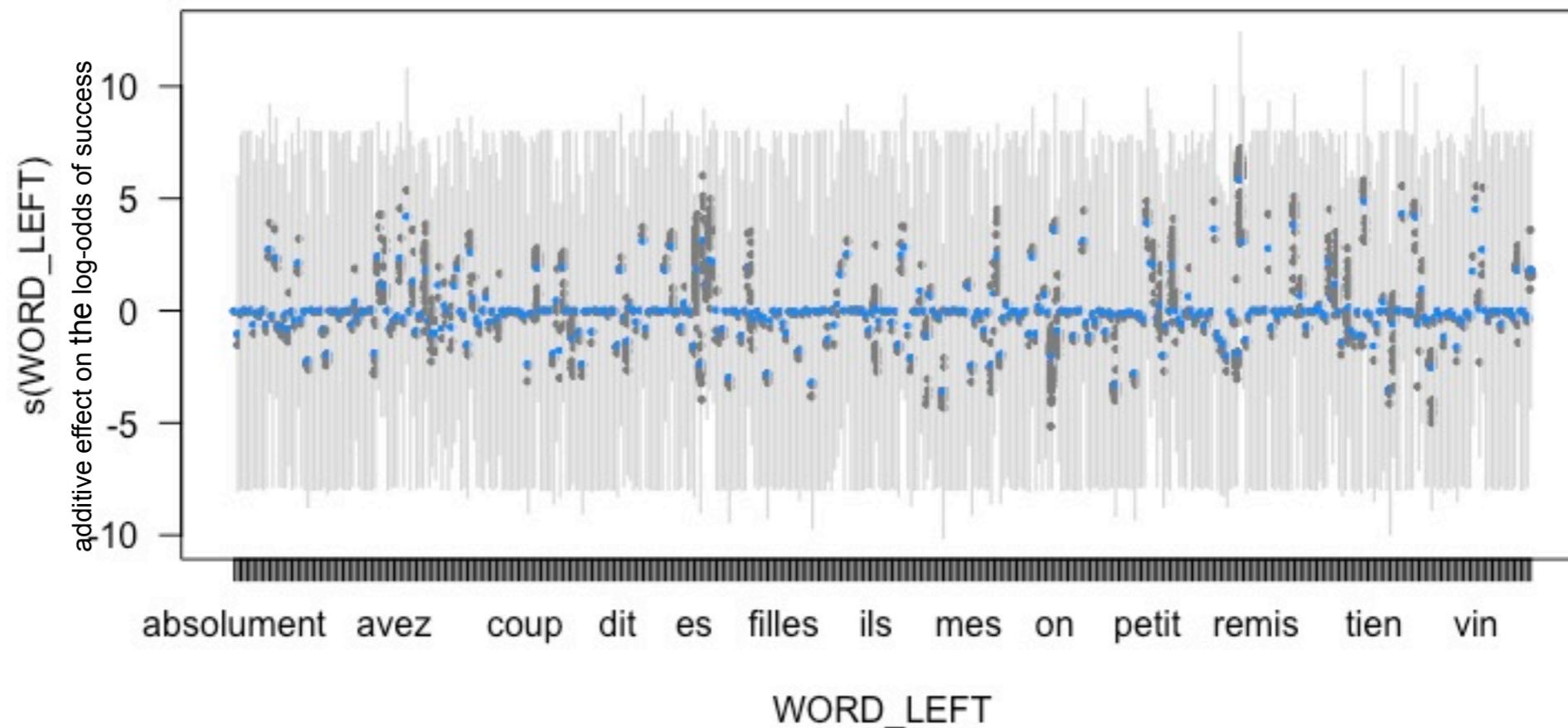
$$\begin{aligned} model_{FINAL} = & -1.915713 + s(AGE) + s(FREQ-G) + s(SYLL-G) + s(TP.BIGRAM) + s(MOT-G) + \\ & + s(POS-G) + s(POS-D) + s(LOCUTEUR) + s(LONGITUDE, LATITUDE) \end{aligned}$$



6.1 Le modèle final - Les fonctions de lissage

$$\begin{aligned} model_{FINAL} = & -1.915713 + s(AGE) + s(FREQ-G) + s(SYLL-G) + s(TP.BIGRAM) + s(MOT-G) + \\ & + s(POS-G) + s(POS-D) + s(LOCUTEUR) + s(LONGITUDE, LATITUDE) \end{aligned}$$

La probabilité de réalisation d'une liaison est spécifique à chaque valeur de MOT-G



6.1 Le modèle final - Les fonctions de lissage

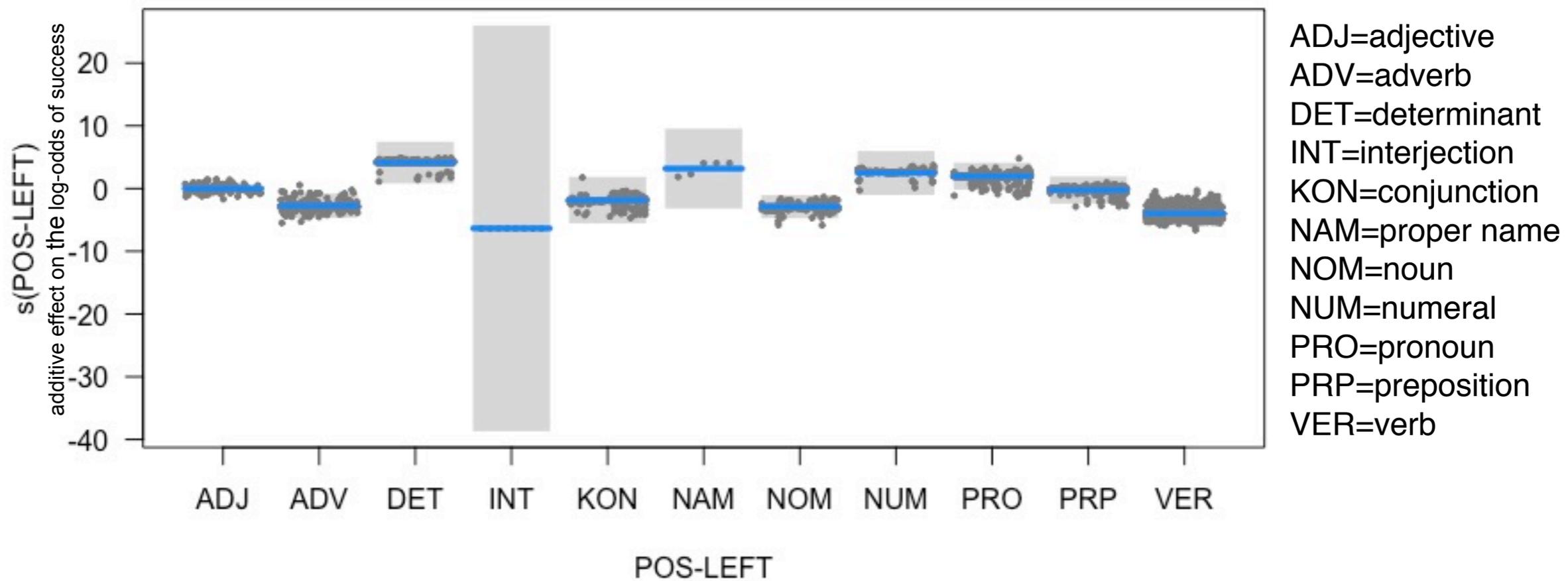
$$\begin{aligned} model_{\text{FINAL}} = & -1.915713 + s(\text{AGE}) + s(\text{FREQ}-G) + s(\text{SYLL}-G) + s(\text{TP.BIGRAM}) + s(\text{MOT}-G) + \\ & + s(\text{POS}-G) + s(\text{POS}-D) + s(\text{LOCUTEUR}) + s(\text{LONGITUDE}, \text{LATITUDE}) \end{aligned}$$

6.1 Le modèle final - Les fonctions de lissage

$$\begin{aligned} model_{\text{FINAL}} = & -1.915713 + s(\text{AGE}) + s(\text{FREQ}-G) + s(\text{SYLL}-G) + s(\text{TP.BIGRAM}) + s(\text{MOT}-G) + \\ & + s(\text{POS}-G) + s(\text{POS}-D) + s(\text{LOCUTEUR}) + s(\text{LONGITUDE}, \text{LATITUDE}) \end{aligned}$$

6.1 Le modèle final - Les fonctions de lissage

$$\text{model}_{\text{FINAL}} = -1.915713 + s(\text{AGE}) + s(\text{FREQ}-G) + s(\text{SYLL}-G) + s(\text{TP.BIGRAM}) + s(\text{MOT}-G) + \\ + s(\text{POS}-G) + s(\text{POS}-D) + s(\text{LOCUTEUR}) + s(\text{LONGITUDE}, \text{LATITUDE})$$



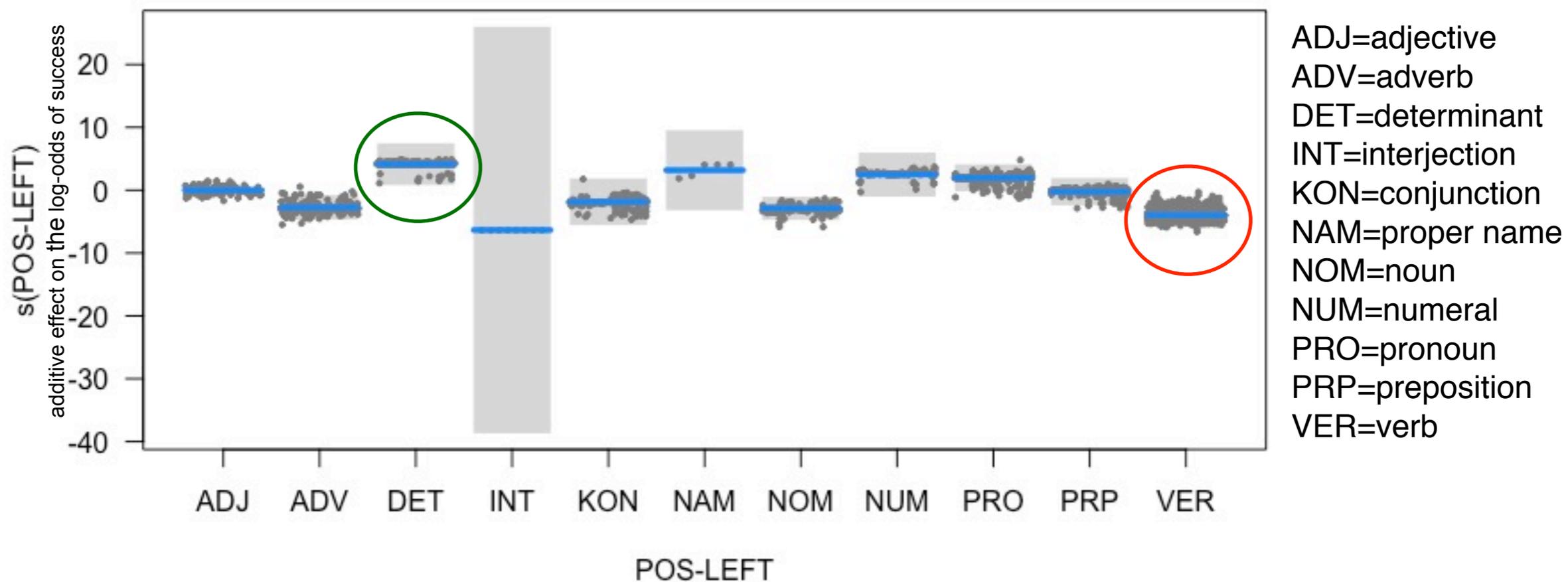
6.1 Le modèle final - Les fonctions de lissage

$$\text{model}_{\text{FINAL}} = -1.915713 + s(\text{AGE}) + s(\text{FREQ}-G) + s(\text{SYLL}-G) + s(\text{TP.BIGRAM}) + s(\text{MOT}-G) + \\ + s(\text{POS}-G) + s(\text{POS}-D) + s(\text{LOCUTEUR}) + s(\text{LONGITUDE}, \text{LATITUDE})$$

La probabilité de réalisation d'une liaison est liée à chaque valeur de POS-G

POS-G qui favorise mieux la réalisation de la liaison : **DET**

POS-G qui favorise moins la réalisation de la liaison : **VER**



6.1 Le modèle final - Les fonctions de lissage

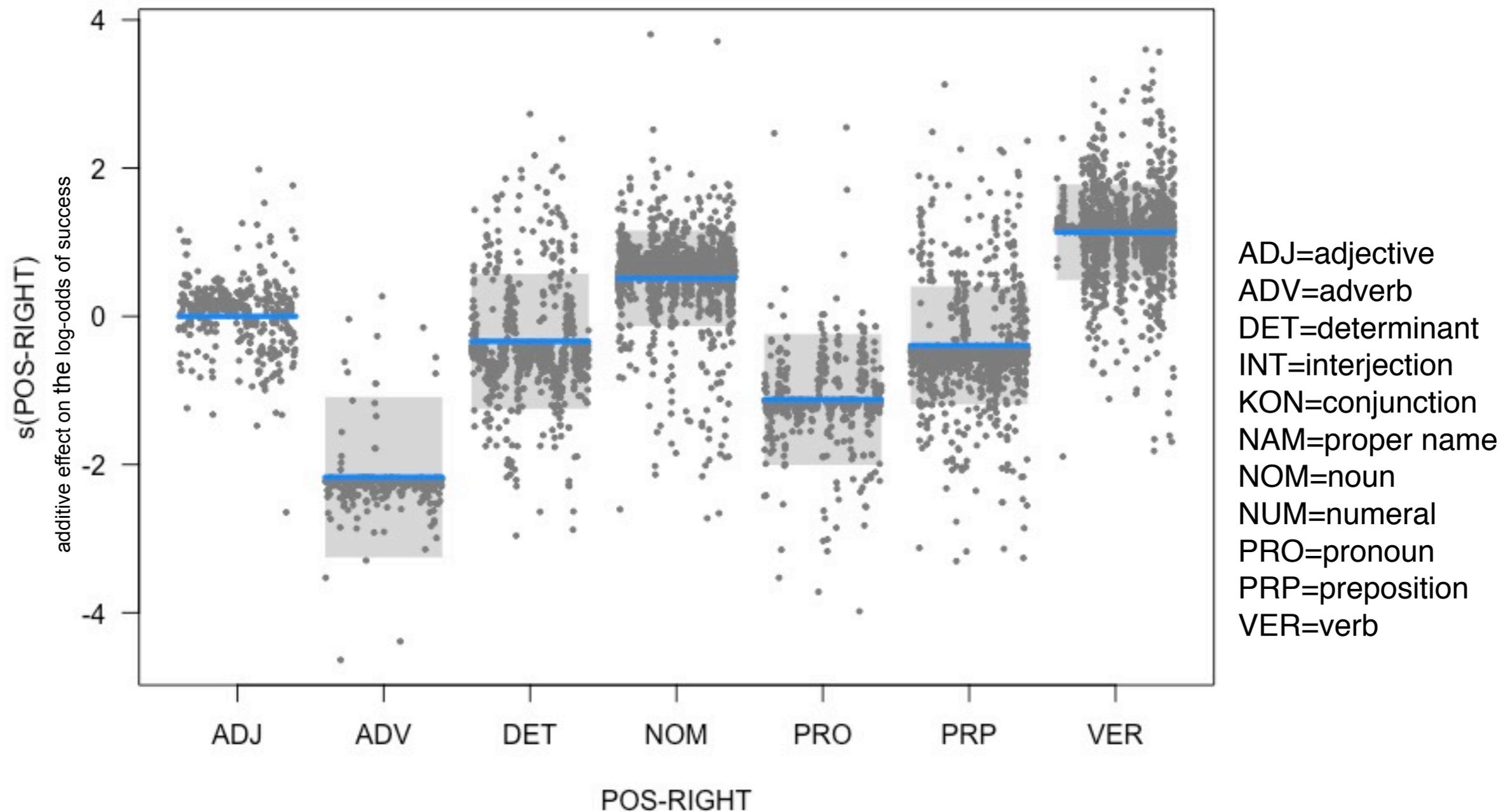
$$\begin{aligned} model_{\text{FINAL}} = & -1.915713 + s(\text{AGE}) + s(\text{FREQ}-G) + s(\text{SYLL}-G) + s(\text{TP.BIGRAM}) + s(\text{MOT}-G) + \\ & + s(\text{POS}-G) + s(\text{POS}-D) + s(\text{LOCUTEUR}) + s(\text{LONGITUDE}, \text{LATITUDE}) \end{aligned}$$

6.1 Le modèle final - Les fonctions de lissage

$$\begin{aligned} model_{\text{FINAL}} = & -1.915713 + s(\text{AGE}) + s(\text{FREQ}-G) + s(\text{SYLL}-G) + s(\text{TP.BIGRAM}) + s(\text{MOT}-G) + \\ & + s(\text{POS}-G) - s(\text{POS}-D) + s(\text{LOCUTEUR}) + s(\text{LONGITUDE}, \text{LATITUDE}) \end{aligned}$$

6.1 Le modèle final - Les fonctions de lissage

$$\begin{aligned} model_{FINAL} = & -1.915713 + s(AGE) + s(FREQ-G) + s(SYLL-G) + s(TP.BIGRAM) + s(MOT-G) + \\ & + s(POS-G) - s(POS-D) + s(LOCUTEUR) + s(LONGITUDE, LATITUDE) \end{aligned}$$



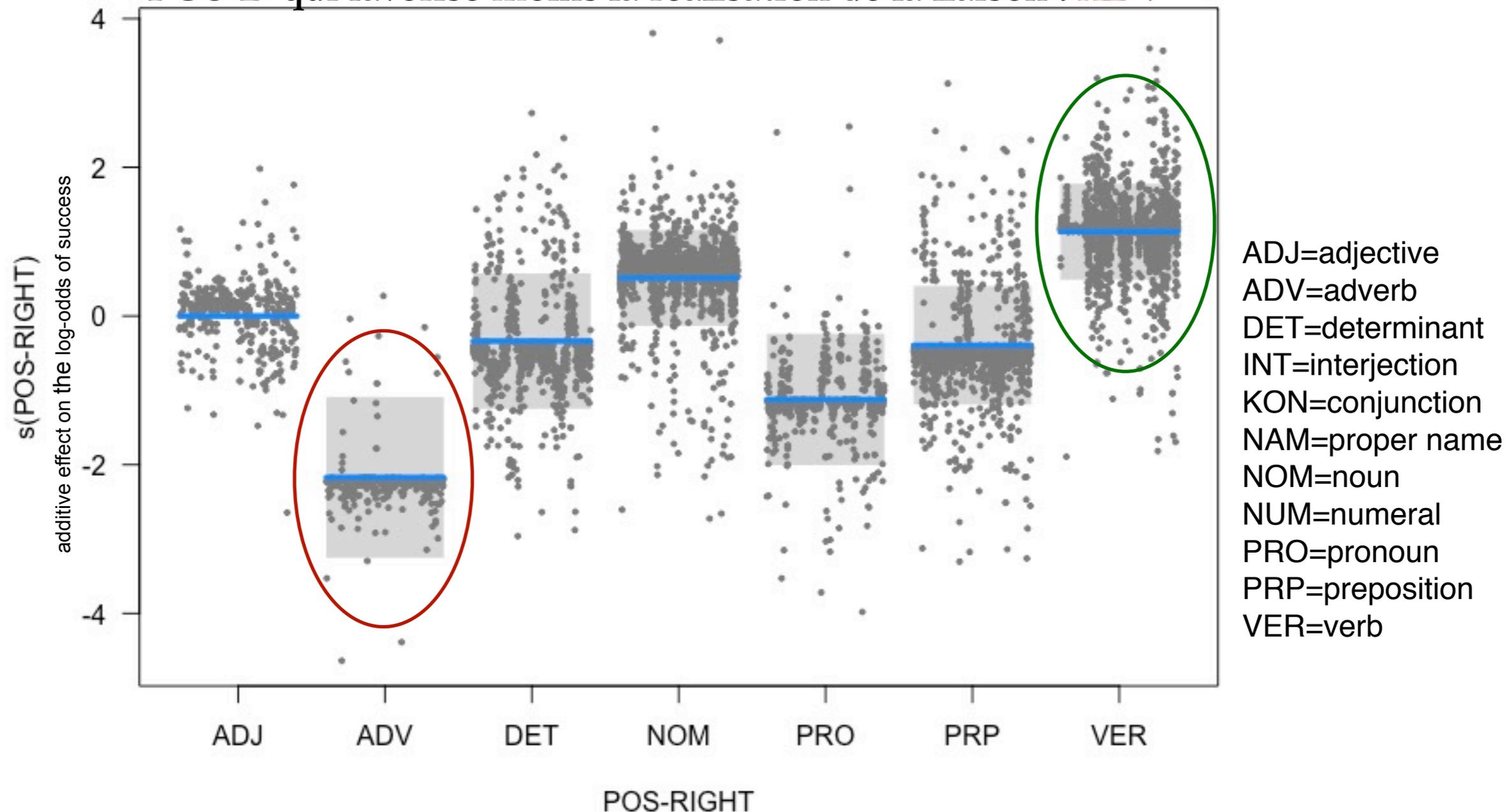
6.1 Le modèle final - Les fonctions de lissage

$$\begin{aligned} model_{FINAL} = & -1.915713 + s(AGE) + s(FREQ-G) + s(SYLL-G) + s(TP.BIGRAM) + s(MOT-G) + \\ & + s(POS-G) - s(POS-D) + s(LOCUTEUR) + s(LONGITUDE, LATITUDE) \end{aligned}$$

La probabilité de réalisation d'une liaison est liée à chaque valeur de POS-D

POS-D qui favorise mieux la réalisation de la liaison : **VER**

POS-D qui favorise moins la réalisation de la liaison : **ADV**



6.1 Le modèle final - Les fonctions de lissage

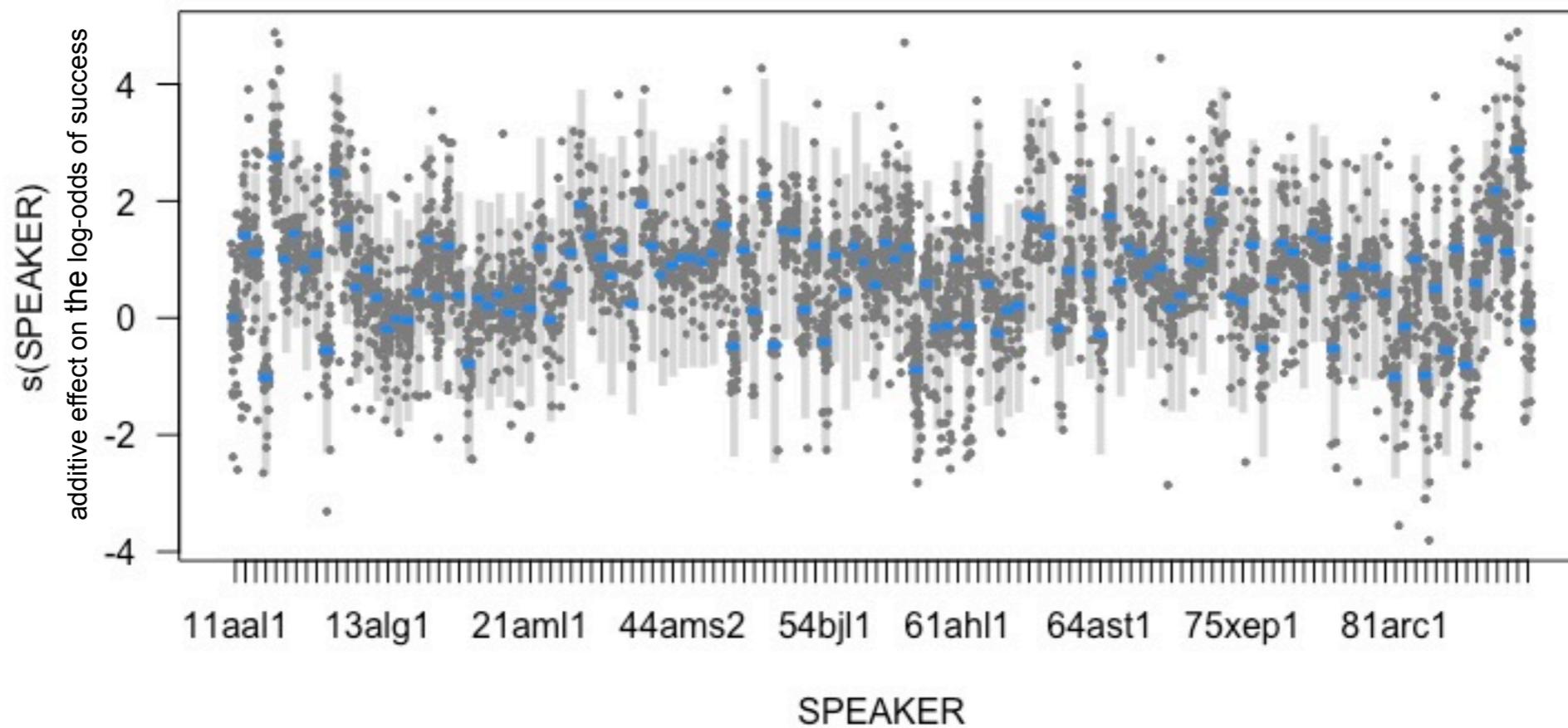
$$\begin{aligned} model_{\text{FINAL}} = & -1.915713 + s(\text{AGE}) + s(\text{FREQ}-G) + s(\text{SYLL}-G) + s(\text{TP.BIGRAM}) + s(\text{MOT}-G) + \\ & + s(\text{POS}-G) + s(\text{POS}-D) + s(\text{LOCUTEUR}) + s(\text{LONGITUDE}, \text{LATITUDE}) \end{aligned}$$

6.1 Le modèle final - Les fonctions de lissage

$$\begin{aligned} model_{\text{FINAL}} = & -1.915713 + s(\text{AGE}) + s(\text{FREQ}-G) + s(\text{SYLL}-G) + s(\text{TP.BIGRAM}) + s(\text{MOT}-G) + \\ & + s(\text{POS}-G) + s(\text{POS}-D) - s(\text{LOCUTEUR}) - s(\text{LONGITUDE}, \text{LATITUDE}) \end{aligned}$$

6.1 Le modèle final - Les fonctions de lissage

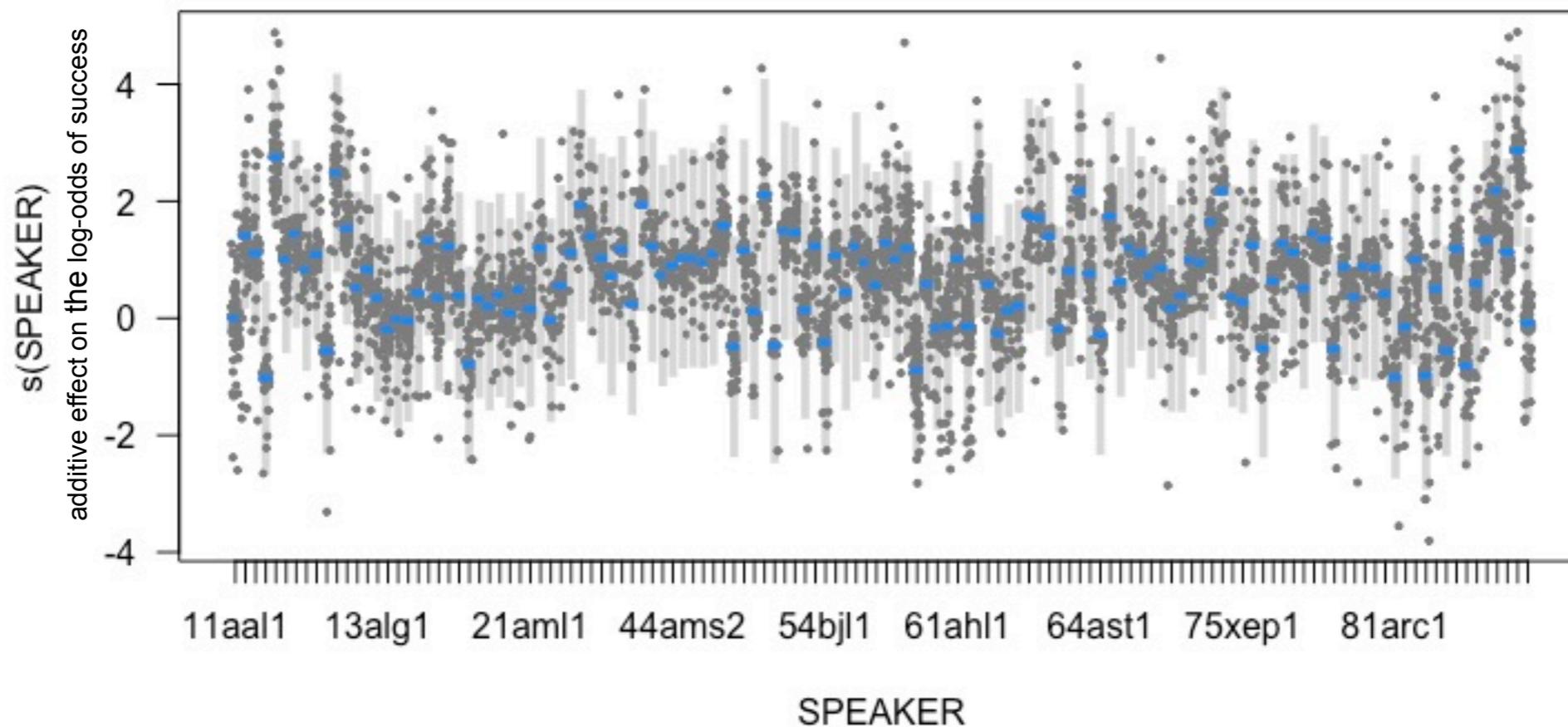
$$\begin{aligned} model_{FINAL} = & -1.915713 + s(AGE) + s(FREQ-G) + s(SYLL-G) + s(TP.BIGRAM) + s(MOT-G) + \\ & + s(POS-G) + s(POS-D) - s(LOCUTEUR) + s(LONGITUDE, LATITUDE) \end{aligned}$$



6.1 Le modèle final - Les fonctions de lissage

$$\begin{aligned} model_{FINAL} = & -1.915713 + s(AGE) + s(FREQ-G) + s(SYLL-G) + s(TP.BIGRAM) + s(MOT-G) + \\ & + s(POS-G) + s(POS-D) - s(LOCUTEUR) + s(LONGITUDE, LATITUDE) \end{aligned}$$

La probabilité de réalisation d'une liaison est spécifique à chaque valeur de LOCUTEUR



6.1 Le modèle final - Les fonctions de lissage

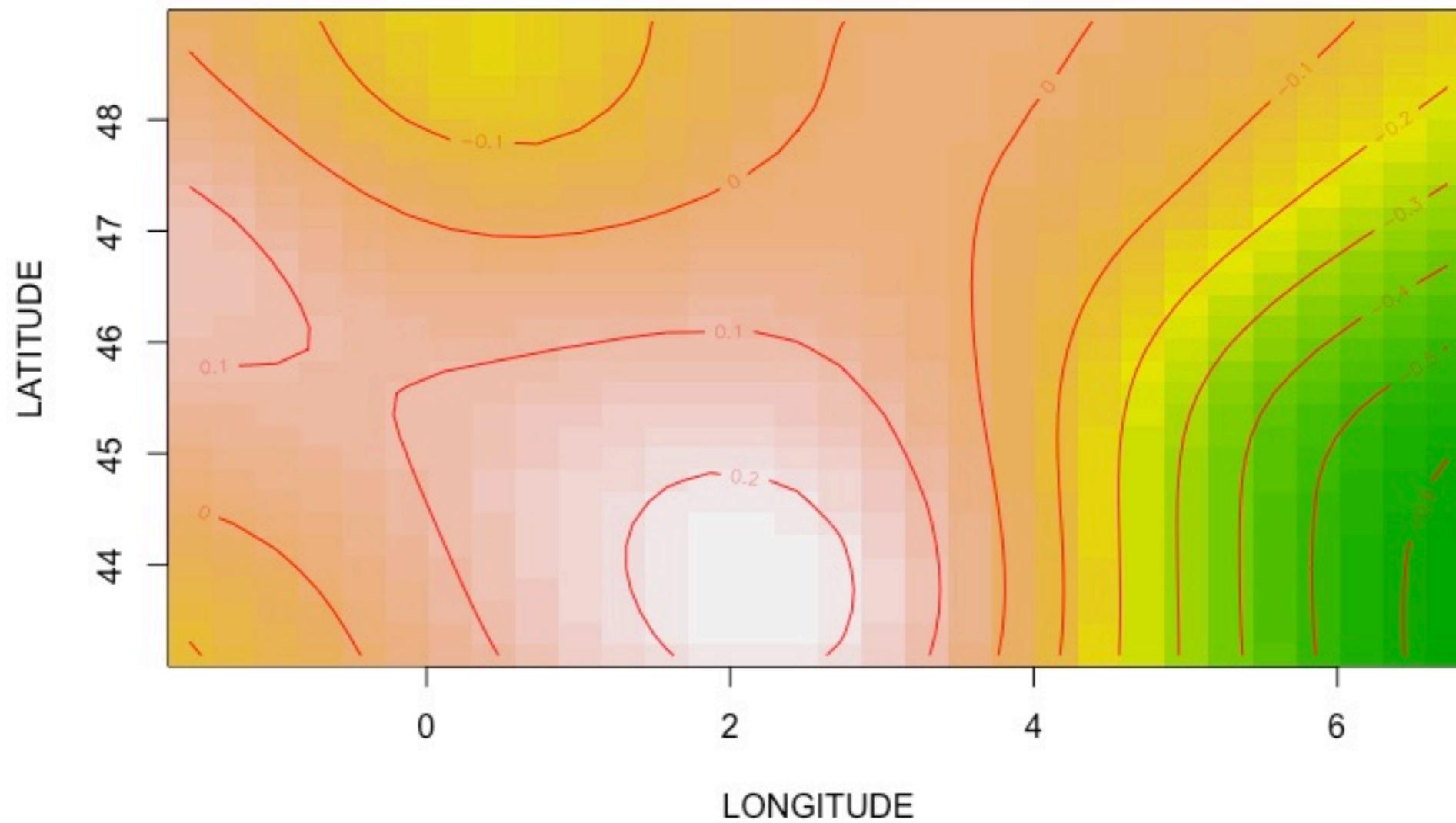
$$\begin{aligned} model_{\text{FINAL}} = & -1.915713 + s(\text{AGE}) + s(\text{FREQ}-G) + s(\text{SYLL}-G) + s(\text{TP.BIGRAM}) + s(\text{MOT}-G) + \\ & + s(\text{POS}-G) + s(\text{POS}-D) + s(\text{LOCUTEUR}) + s(\text{LONGITUDE}, \text{LATITUDE}) \end{aligned}$$

6.1 Le modèle final - Les fonctions de lissage

$$\begin{aligned} model_{\text{FINAL}} = & -1.915713 + s(\text{AGE}) + s(\text{FREQ}-G) + s(\text{SYLL}-G) + s(\text{TP.BIGRAM}) + s(\text{MOT}-G) + \\ & + s(\text{POS}-G) + s(\text{POS}-D) + s(\text{LOCUTEUR}) + s(\text{LONGITUDE}, \text{LATITUDE}) \end{aligned}$$

6.1 Le modèle final - Les fonctions de lissage

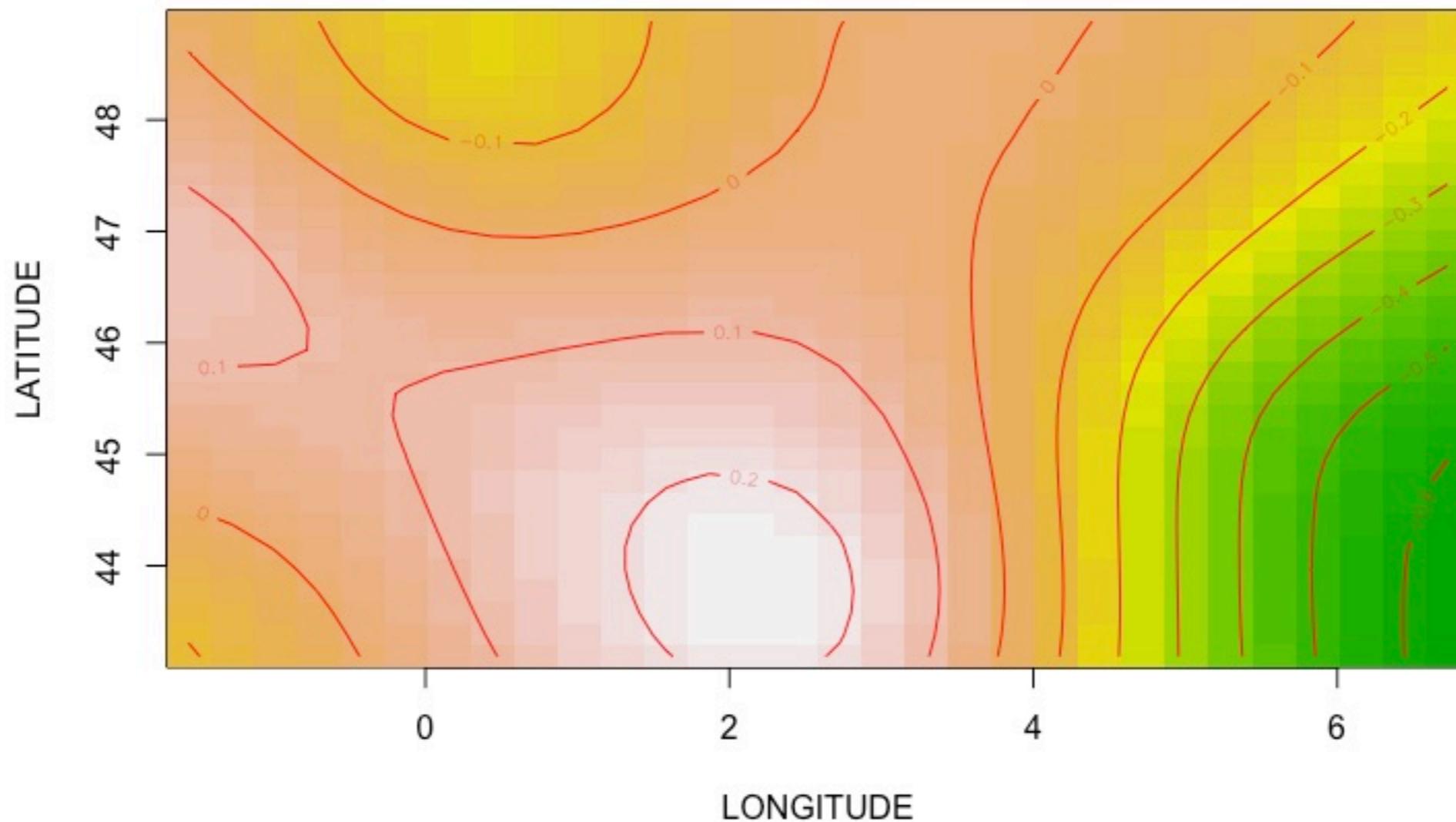
$$\begin{aligned} model_{FINAL} = & -1.915713 + s(AGE) + s(FREQ-G) + s(SYLL-G) + s(TP.BIGRAM) + s(MOT-G) + \\ & + s(POS-G) + s(POS-D) + s(LOCUTEUR) + s(LONGITUDE, LATITUDE) \end{aligned}$$



6.1 Le modèle final - Les fonctions de lissage

$$\begin{aligned} model_{FINAL} = & -1.915713 + s(AGE) + s(FREQ-G) + s(SYLL-G) + s(TP.BIGRAM) + s(MOT-G) + \\ & + s(POS-G) + s(POS-D) + s(LOCUTEUR) + s(LONGITUDE, LATITUDE) \end{aligned}$$

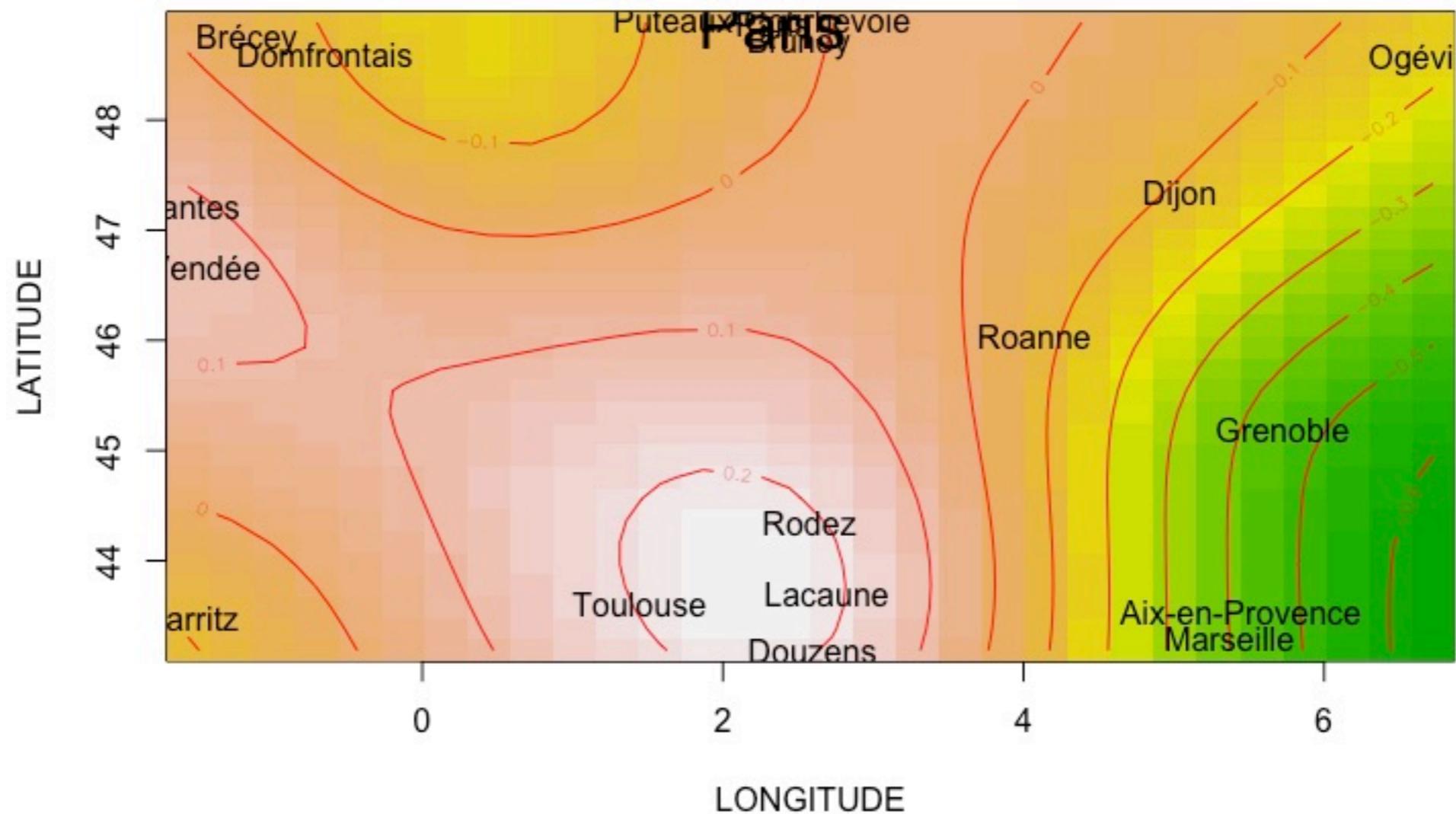
La probabilité de réalisation d'une liaison par rapport à la **LONGITUDE** et **LATITUDE** du locuteur



6.1 Le modèle final - Les fonctions de lissage

$$\begin{aligned} model_{FINAL} = & -1.915713 + s(AGE) + s(FREQ-G) + s(SYLL-G) + s(TP.BIGRAM) + s(MOT-G) + \\ & + s(POS-G) + s(POS-D) + s(LOCUTEUR) + s(LONGITUDE, LATITUDE) \end{aligned}$$

La probabilité de réalisation d'une liaison par rapport à la **LONGITUDE** et **LATITUDE** du locuteur



- Objectif du travail
- La liaison - une petite présentation
- Un corpus phonologique: PFC
- Les outils statistiques (RL et GAM)
- Le model final et le résultats
- Limites du travail et perspectives futures

7. Conclusions

7. Conclusions

- Nous avons créé un modèle prédictif qui explique plus que 83% de la variabilité de la liaison en français

7. Conclusions

- Nous avons créé un modèle prédictif qui explique plus que 83% de la variabilité de la liaison en français
- Parmi les variables explicatives, le mot à gauche (et son POS) semble être la variable la plus explicative du phénomène liaison, suivi par des variables distributionnelles comme la fréquence du mot à gauche et des variables socio-linguistiques comme l'âge du locuteur

7. Conclusions

- Nous avons créé un modèle prédictif qui explique plus que 83% de la variabilité de la liaison en français
- Parmi les variables explicatives, le mot à gauche (et son POS) semble être la variable la plus explicative du phénomène liaison, suivi par des variables distributionnelles comme la fréquence du mot à gauche et des variables socio-linguistiques comme l'âge du locuteur
- Le modèle est prédictif. Il est possible prédire la probabilité de liaison de données inconnues : $P[(Liaison=oui) \mid \text{M. Dupont de Bordeaux de 70 ans face à la construction } \langle \text{pas encore} \rangle]$

7. Conclusions

- Nous avons créé un modèle prédictif qui explique plus que 83% de la variabilité de la liaison en français
- Parmi les variables explicatives, le mot à gauche (et son POS) semble être la variable la plus explicative du phénomène liaison, suivi par des variables distributionnelles comme la fréquence du mot à gauche et des variables socio-linguistiques comme l'âge du locuteur
- Le modèle est prédictif. Il est possible prédire la probabilité de liaison de données inconnues : $P[(Liaison=oui) \mid \text{M. Dupont de Bordeaux de 70 ans face à la construction } \langle \text{pas encore} \rangle]$
- Le modèle incorpore aussi des **effets aléatoires** (pas illustrés aujourd'hui) pour les variables liés aux **locuteurs** et aux **mots** (Baayen, 2014).

7. Conclusions

- Nous avons créé un modèle prédictif qui explique plus que 83% de la variabilité de la liaison en français
- Parmi les variables explicatives, le mot à gauche (et son POS) semble être la variable la plus explicative du phénomène liaison, suivi par des variables distributionnelles comme la fréquence du mot à gauche et des variables socio-linguistiques comme l'âge du locuteur
- Le modèle est prédictif. Il est possible prédire la probabilité de liaison de données inconnues : $P[(Liaison=oui) \mid \text{M. Dupont de Bordeaux de 70 ans face à la construction } \langle \text{pas encore} \rangle]$
- Le modèle incorpore aussi des **effets aléatoires** (pas illustrés aujourd'hui) pour les variables liés aux **locuteurs** et aux **mots** (Baayen, 2014).
 - Nous n'avons pas enregistré le comportement de liaison chez tous les locuteurs francophones, mais seulement chez un nombre restreint (192 locuteurs), pris aléatoirement.

7. Conclusions

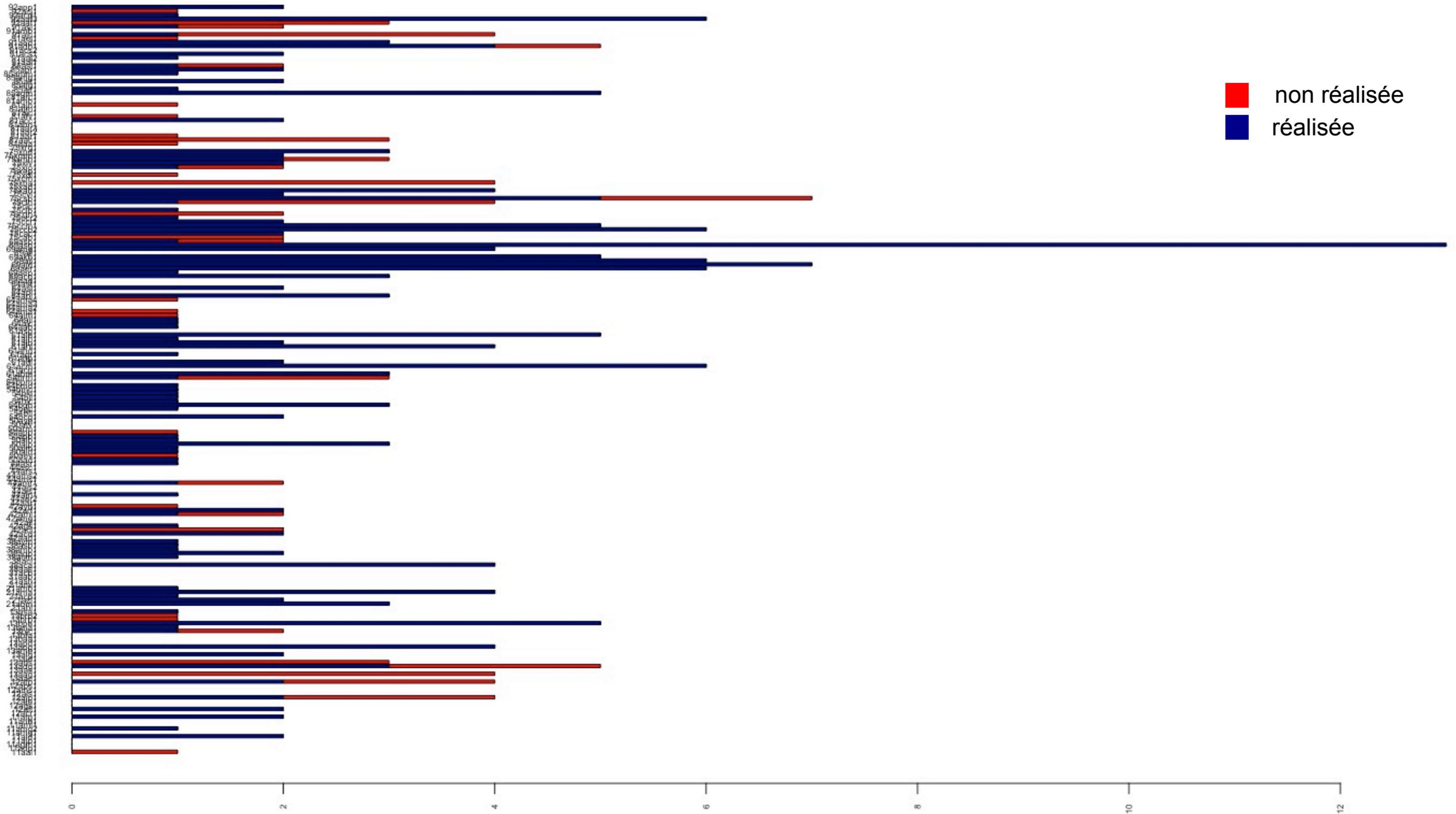
- Nous avons créé un modèle prédictif qui explique plus que 83% de la variabilité de la liaison en français
- Parmi les variables explicatives, le mot à gauche (et son POS) semble être la variable la plus explicative du phénomène liaison, suivi par des variables distributionnelles comme la fréquence du mot à gauche et des variables socio-linguistiques comme l'âge du locuteur
- Le modèle est prédictif. Il est possible prédire la probabilité de liaison de données inconnues : $P[(Liaison=oui) \mid \text{M. Dupont de Bordeaux de 70 ans face à la construction } \langle \text{pas encore} \rangle]$
- Le modèle incorpore aussi des **effets aléatoires** (pas illustrés aujourd'hui) pour les variables liés aux **locuteurs** et aux **mots** (Baayen, 2014).
 - Nous n'avons pas enregistré le comportement de liaison chez tous les locuteurs francophones, mais seulement chez un nombre restreint (192 locuteurs), pris aléatoirement.
 - Il y aura sûrement des différences de comportement de liaison entre un locuteur et l'autre et même chez le même locuteur. Cette variation est due à des facteurs incontrôlables (sujets plus sensibles à la langue, tension/émotion lors de la registration, etc).

7. Conclusions

- Nous avons créé un modèle prédictif qui explique plus que 83% de la variabilité de la liaison en français
- Parmi les variables explicatives, le mot à gauche (et son POS) semble être la variable la plus explicative du phénomène liaison, suivi par des variables distributionnelles comme la fréquence du mot à gauche et des variables socio-linguistiques comme l'âge du locuteur
- Le modèle est prédictif. Il est possible prédire la probabilité de liaison de données inconnues : $P[(Liaison=oui) \mid \text{M. Dupont de Bordeaux de 70 ans face à la construction } \langle \text{pas encore} \rangle]$
- Le modèle incorpore aussi des **effets aléatoires** (pas illustrés aujourd'hui) pour les variables liés aux **locuteurs** et aux **mots** (Baayen, 2014).
 - Nous n'avons pas enregistré le comportement de liaison chez tous les locuteurs francophones, mais seulement chez un nombre restreint (192 locuteurs), pris aléatoirement.
 - Il y aura sûrement des différences de comportement de liaison entre un locuteur et l'autre et même chez le même locuteur. Cette variation est due à des facteurs incontrôlables (sujets plus sensibles à la langue, tension/émotion lors de la registration, etc).
 - Grâce à l'introduction d'un effet aléatoire spécifique à chaque sujet (et à chaque mot) le modèle a pris en compte ces différences.

7. Conclusions

Distribution de <est un> chez les locuteurs



7. Limites du travail et perspectives futures

- Une limite de ce modèle est le fait de ne pas prendre en compte la dimension syntaxique (bien que l'on considère les informations de POS des mots à gauche et à droite, on ne possède pas d'information sur les syntagmes dans lesquels ces mots sont impliqués)

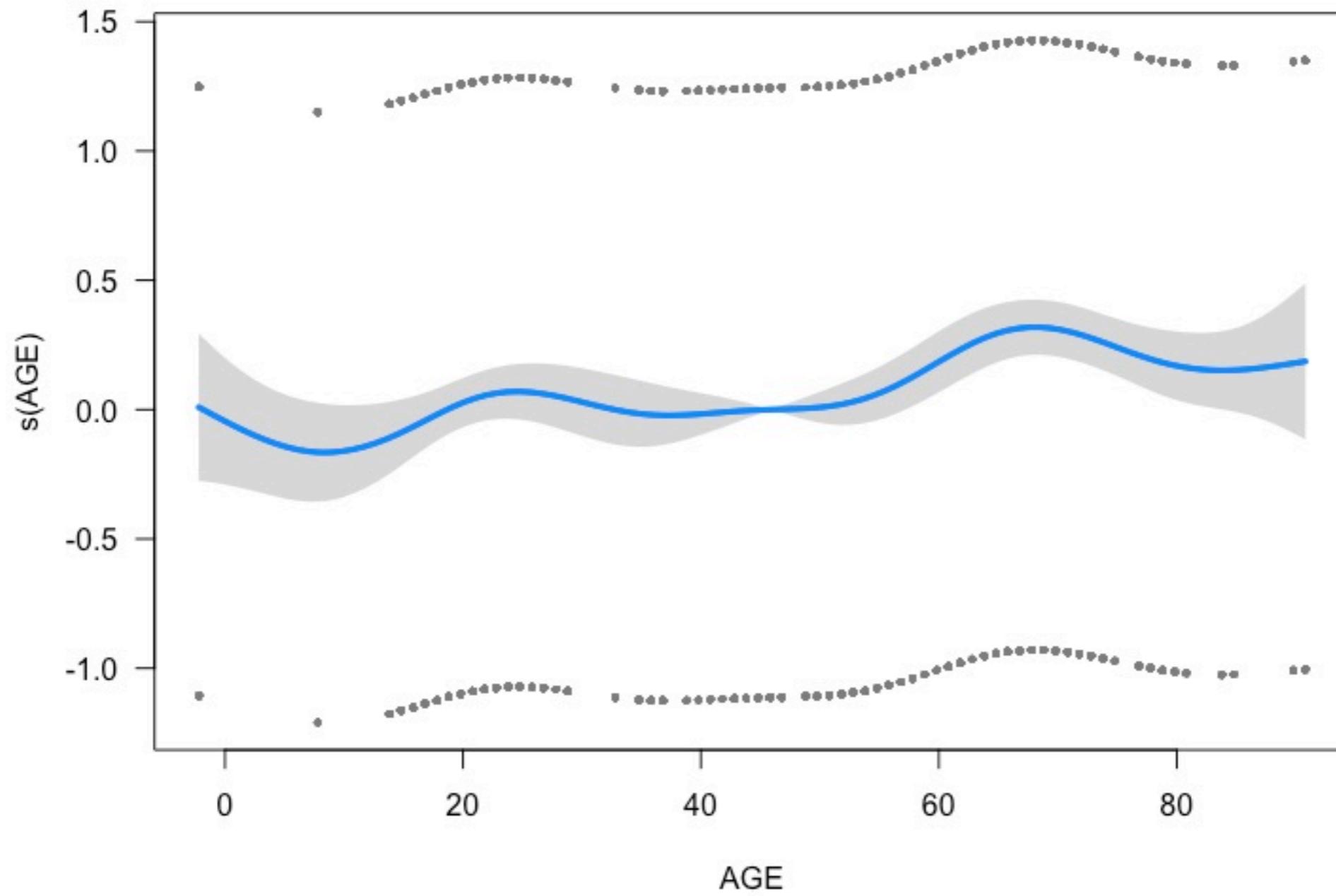
- Une limite de ce modèle est le fait de ne pas prendre en compte la dimension syntaxique (bien que l'on considère les informations de POS des mots à gauche et à droite, on ne possède pas d'information sur les syntagmes dans lesquels ces mots sont impliqués)
- Une limite ultérieure concerne les aspects phonotactiques : les informations sur les phonèmes à gauche et à droite pourraient en effet enrichir le modèle

7. Limites du travail et perspectives futures

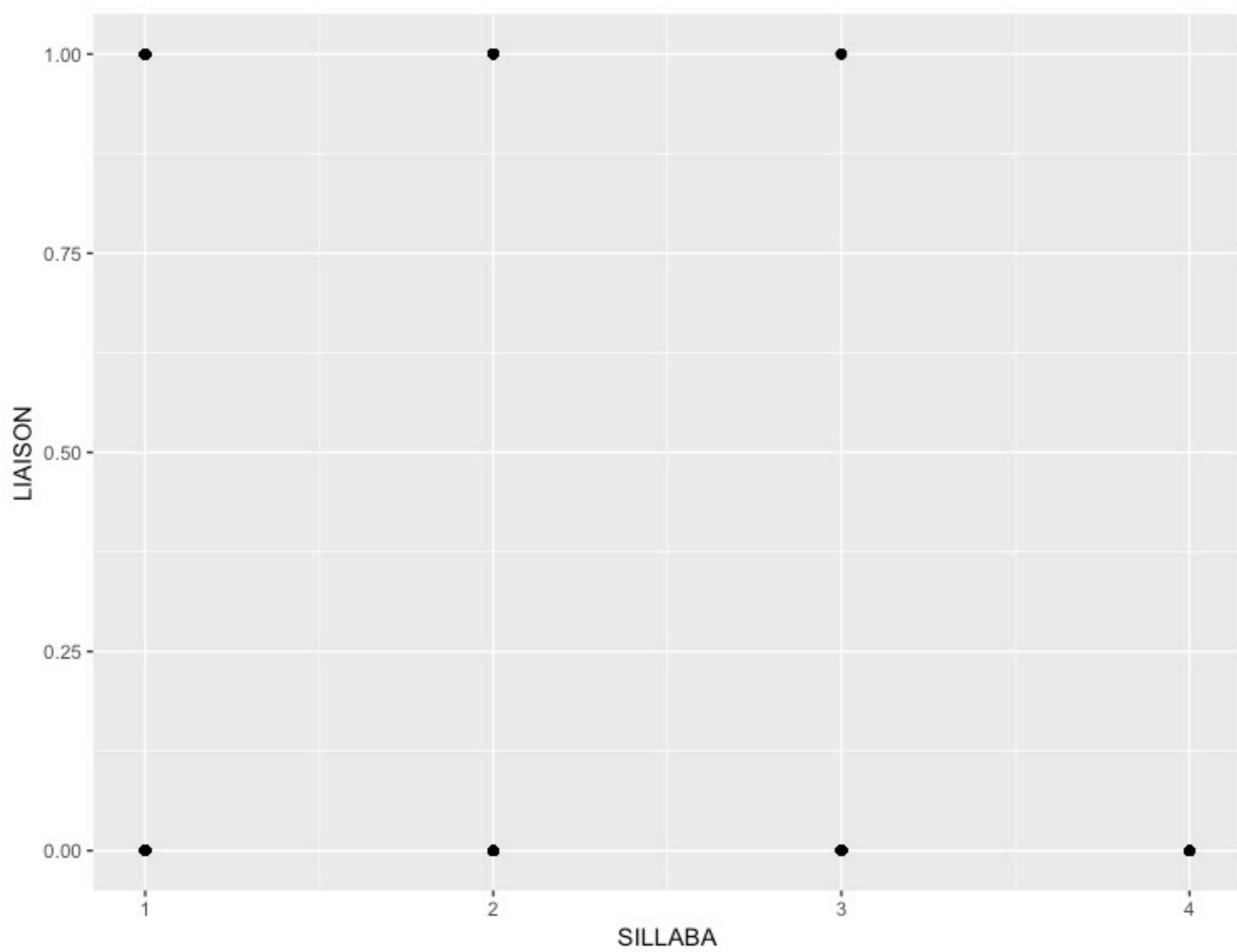
- Une limite de ce modèle est le fait de ne pas prendre en compte la dimension syntaxique (bien que l'on considère les informations de POS des mots à gauche et à droite, on ne possède pas d'information sur les syntagmes dans lesquels ces mots sont impliqués)
- Une limite ultérieure concerne les aspects phonotactiques : les informations sur les phonèmes à gauche et à droite pourraient en effet enrichir le modèle
- Considérer aussi la dimension holophrastique d'occurrence de la liaison

MERCI

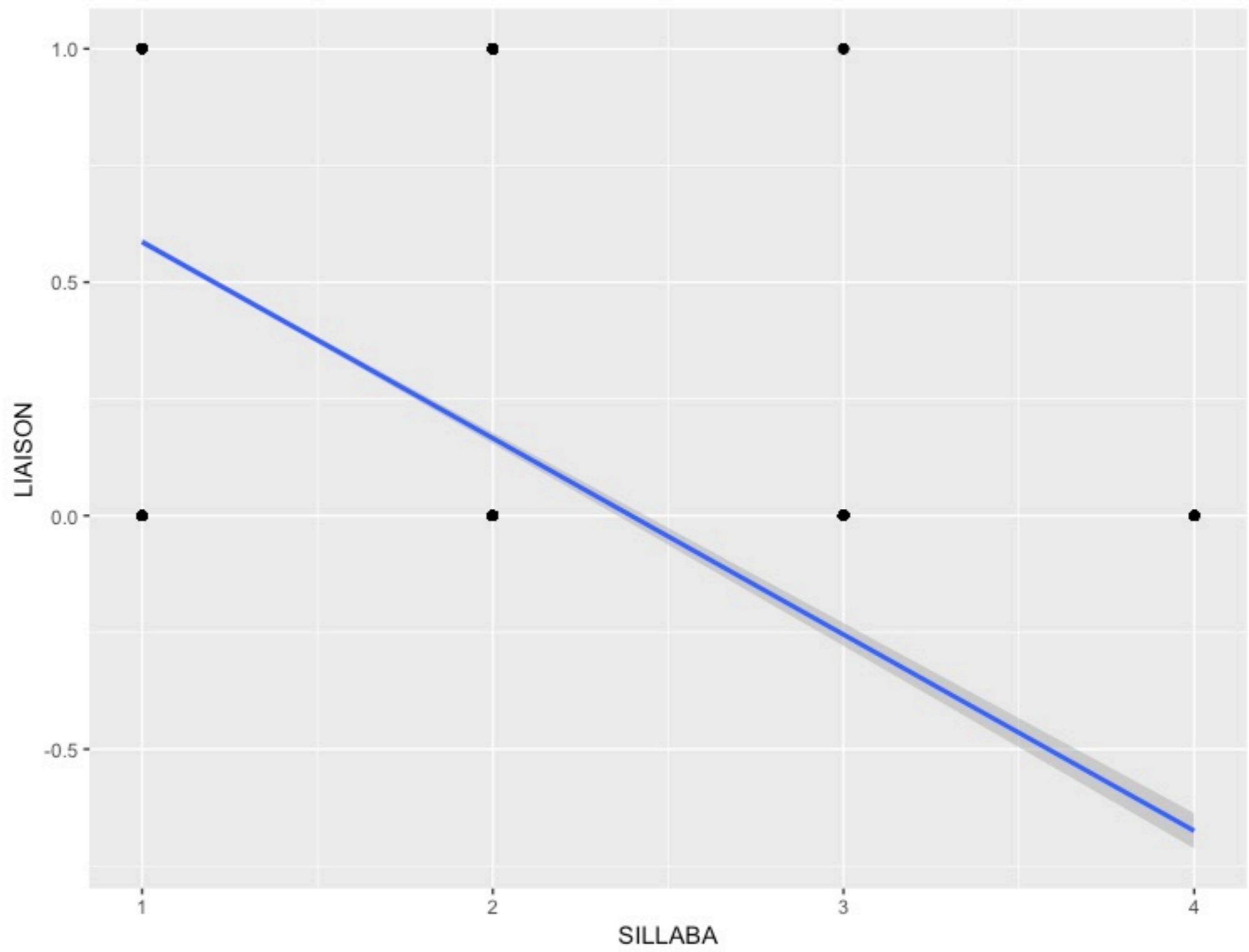
1. La liaison en français



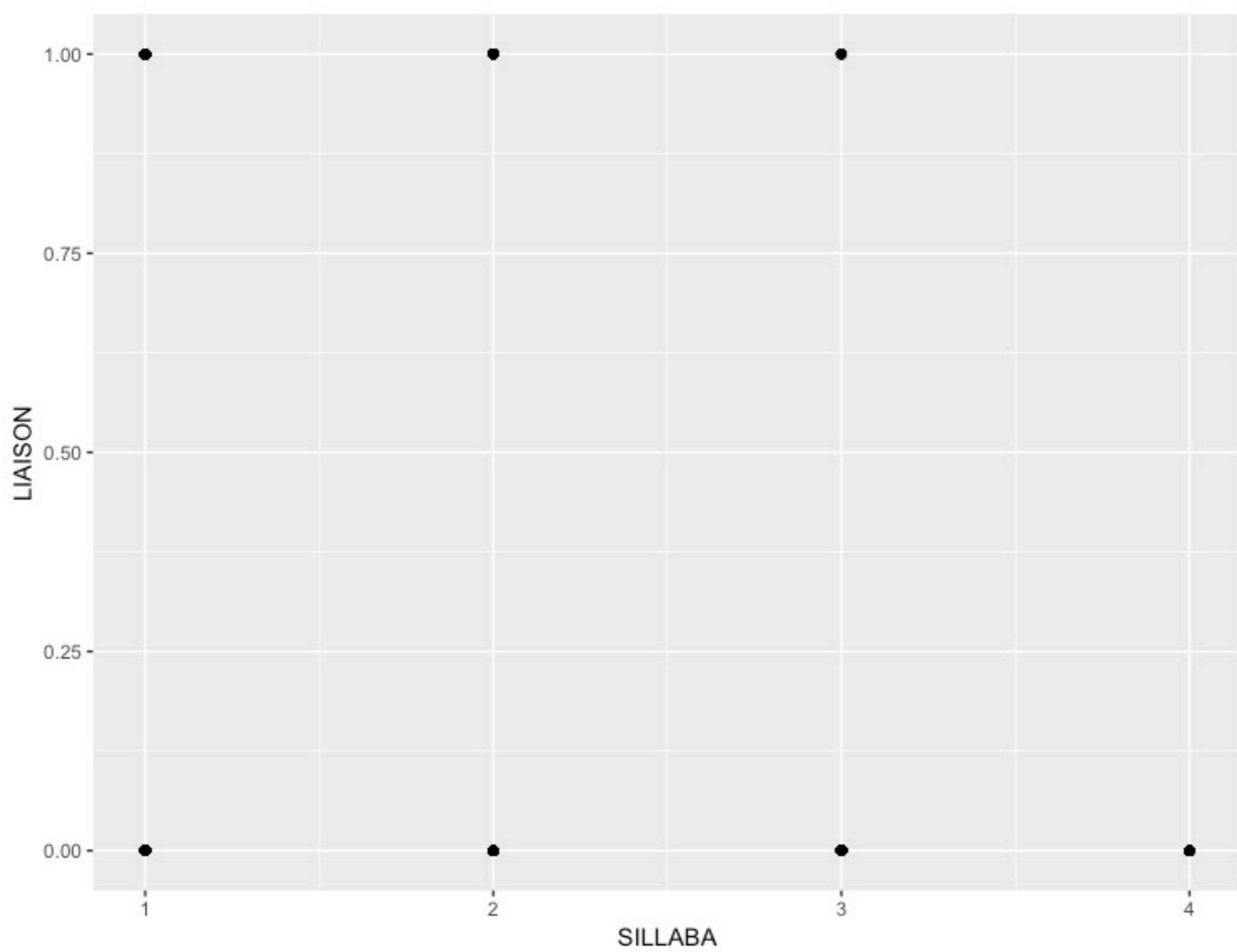
1. La liaison en français



1. La liaison en français



1. La liaison en français



1. La liaison en français

