

Introduction aux GLM

L. Rouvière
laurent.rouviere@univ-rennes2.fr

JANVIER 2015

1 Modèle statistique

- Modèle de densité
- Modèle de régression
- Rappels sur le modèle de régression linéaire

2 Introduction au modèle de régression logistique

- Exemples
- Régression logistique simple

3 Le modèle linéaire généralisé

- Introduction
- Définitions
- Modèle de Poisson

Bibliographie

Modèle statistique

Qu'est ce qu'un modèle ?

Mathématiquement, un modèle est un triplet $(\mathcal{H}, \mathcal{A}, \{P, P \in \mathcal{P}\})$ avec

- \mathcal{H} est l'espace des observations (l'ensemble de tous les résultats possibles de l'expérience) ;
- \mathcal{A} est une tribu sur \mathcal{H} ;
- \mathcal{P} est une famille de probabilités définie sur $(\mathcal{H}, \mathcal{A})$.

A quoi sert un modèle ?

Expliquer, décrire les mécanismes du phénomène considéré.

- **Question** : quel est le lien entre la définition mathématique et l'utilité du phénomène ?

Qu'est ce qu'un modèle ?

Mathématiquement, un modèle est un triplet $(\mathcal{H}, \mathcal{A}, \{P, P \in \mathcal{P}\})$ avec

- \mathcal{H} est l'espace des observations (l'ensemble de tous les résultats possibles de l'expérience) ;
- \mathcal{A} est une tribu sur \mathcal{H} ;
- \mathcal{P} est une famille de probabilités définie sur $(\mathcal{H}, \mathcal{A})$.

A quoi sert un modèle ?

Expliquer, décrire les mécanismes du phénomène considéré.

- **Question** : quel est le lien entre la définition mathématique et l'utilité du phénomène ?

Qu'est ce qu'un modèle ?

Mathématiquement, un modèle est un triplet $(\mathcal{H}, \mathcal{A}, \{P, P \in \mathcal{P}\})$ avec

- \mathcal{H} est l'espace des observations (l'ensemble de tous les résultats possibles de l'expérience) ;
- \mathcal{A} est une tribu sur \mathcal{H} ;
- \mathcal{P} est une famille de probabilités définie sur $(\mathcal{H}, \mathcal{A})$.

A quoi sert un modèle ?

Expliquer, décrire les mécanismes du phénomène considéré.

- **Question** : quel est le lien entre la définition mathématique et l'utilité du phénomène ?

Qu'est ce qu'un modèle ?

Mathématiquement, un modèle est un triplet $(\mathcal{H}, \mathcal{A}, \{P, P \in \mathcal{P}\})$ avec

- \mathcal{H} est l'espace des observations (l'ensemble de tous les résultats possibles de l'expérience) ;
- \mathcal{A} est une tribu sur \mathcal{H} ;
- \mathcal{P} est une famille de probabilités définie sur $(\mathcal{H}, \mathcal{A})$.

A quoi sert un modèle ?

Expliquer, décrire les mécanismes du phénomène considéré.

- **Question** : quel est le lien entre la définition mathématique et l'utilité du phénomène ?

Qu'est ce qu'un modèle ?

Mathématiquement, un modèle est un triplet $(\mathcal{H}, \mathcal{A}, \{P, P \in \mathcal{P}\})$ avec

- \mathcal{H} est l'espace des observations (l'ensemble de tous les résultats possibles de l'expérience) ;
- \mathcal{A} est une tribu sur \mathcal{H} ;
- \mathcal{P} est une famille de probabilités définie sur $(\mathcal{H}, \mathcal{A})$.

A quoi sert un modèle ?

Expliquer, décrire les mécanismes du phénomène considéré.

- **Question** : quel est le lien entre la définition mathématique et l'utilité du phénomène ?

Modèle de densité

Exemple 1

- On souhaite tester l'efficacité d'un nouveau traitement à l'aide d'un essai clinique.
- On traite $n = 100$ patients atteints de la pathologie.
- A l'issue de l'étude, 72 patients sont guéris.
- Soit p_0 la probabilité de guérison suite au traitement en question.
- On est tentés de conclure $p_0 \approx 0.72$.

Un tel résultat n'a cependant guère d'intérêt si on n'est pas capable de préciser l'erreur susceptible d'être commise par cette estimation.

Exemple 1

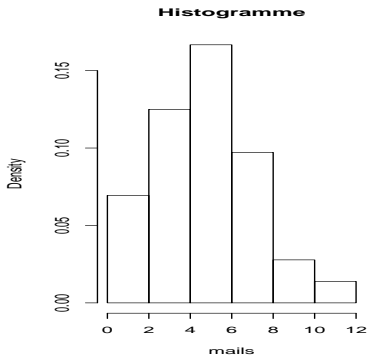
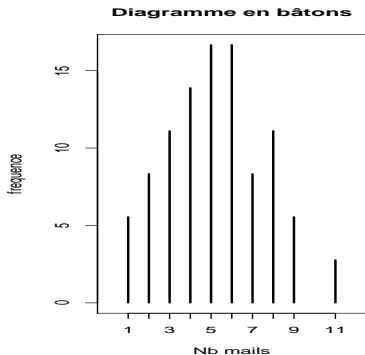
- On souhaite tester l'efficacité d'un nouveau traitement à l'aide d'un essai clinique.
- On traite $n = 100$ patients atteints de la pathologie.
- A l'issue de l'étude, 72 patients sont guéris.

- Soit p_0 la probabilité de guérison suite au traitement en question.
- On est tentés de conclure $p_0 \approx 0.72$.

Un tel résultat n'a cependant guère d'intérêt si on n'est pas capable de préciser l'erreur susceptible d'être commise par cette estimation.

Exemple 2

- On s'intéresse au nombre de mails reçus par jour par un utilisateur pendant 36 journées.
- $\bar{x} = 5.22$, $S_n^2 = 5.72$.

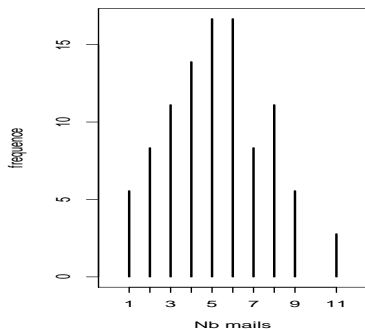


Quelle est la probabilité de recevoir plus de 5 mails dans une journée ?

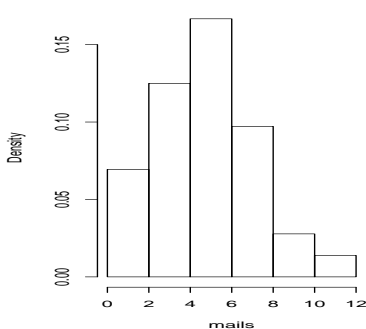
Exemple 2

- On s'intéresse au nombre de mails reçus par jour par un utilisateur pendant 36 journées.
- $\bar{x} = 5.22$, $S_n^2 = 5.72$.

Diagramme en bâtons



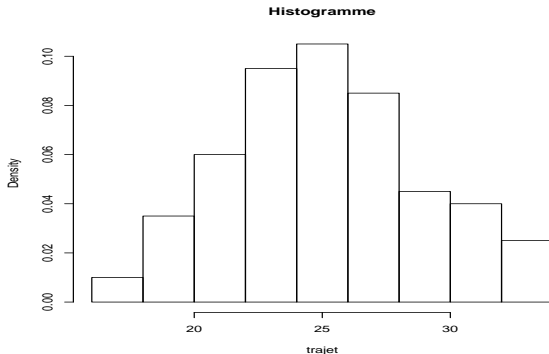
Histogramme



Quelle est la probabilité de recevoir plus de 5 mails dans une journée ?

Exemple 3

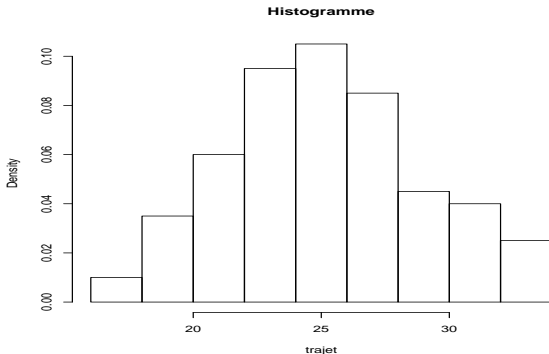
- Durée de trajet domicile-travail.
- On dispose de $n = 100$ mesures : $\bar{x} = 25.1$, $S_n^2 = 14.46$.



J'ai une réunion à 8h30, quelle est la probabilité que j'arrive en retard si je pars de chez moi à 7h55 ?

Exemple 3

- Durée de trajet domicile-travail.
- On dispose de $n = 100$ mesures : $\bar{x} = 25.1$, $S_n^2 = 14.46$.



J'ai une réunion à 8h30, quelle est la probabilité que j'arrive en retard si je pars de chez moi à 7h55 ?

Problème

- Nécessité de se dégager des observations x_1, \dots, x_n pour répondre à de telles questions.
- Si on mesure la durée du trajet pendant 100 nouveaux jours, on peut en effet penser que les nouvelles observations ne seront pas exactement les mêmes que les anciennes.

Idée

Considérer que les n valeurs observées x_1, \dots, x_n sont des réalisations de variables aléatoires X_1, \dots, X_n .

Attention

X_i est une variable aléatoire et x_i est une réalisation de cette variable, c'est-à-dire un nombre !

Problème

- Nécessité de se dégager des observations x_1, \dots, x_n pour répondre à de telles questions.
- Si on mesure la durée du trajet pendant 100 nouveaux jours, on peut en effet penser que les nouvelles observations ne seront pas exactement les mêmes que les anciennes.

Idée

Considérer que les n valeurs observées x_1, \dots, x_n sont des réalisations de variables aléatoires X_1, \dots, X_n .

Attention

X_i est une variable aléatoire et x_i est une réalisation de cette variable, c'est-à-dire un nombre !

Problème

- Nécessité de se dégager des observations x_1, \dots, x_n pour répondre à de telles questions.
- Si on mesure la durée du trajet pendant 100 nouveaux jours, on peut en effet penser que les nouvelles observations ne seront pas exactement les mêmes que les anciennes.

Idée

Considérer que les n valeurs observées x_1, \dots, x_n sont des réalisations de variables aléatoires X_1, \dots, X_n .

Attention

X_i est une variable aléatoire et x_i est une réalisation de cette variable, c'est-à-dire un nombre !

Définition

Une **variable aléatoire réelle** est une application

$$X : (\Omega, \mathcal{A}) \rightarrow (\mathbb{R}, \mathcal{B}(\mathbb{R}))$$

telle que

$$\forall B \in \mathcal{B}(\mathbb{R}), X^{-1}(B) \in \mathcal{A}.$$

- Lors de la modélisation statistique, l'espace Ω n'est généralement jamais caractérisé.
- Il contient tous les "phénomènes" pouvant expliquer les sources d'aléa (qui ne sont pas explicables...).
- En pratique, l'espace d'arrivée est généralement suffisant.

Définition

Une **variable aléatoire réelle** est une application

$$X : (\Omega, \mathcal{A}) \rightarrow (\mathbb{R}, \mathcal{B}(\mathbb{R}))$$

telle que

$$\forall B \in \mathcal{B}(\mathbb{R}), X^{-1}(B) \in \mathcal{A}.$$

- Lors de la modélisation statistique, l'espace Ω n'est généralement jamais caractérisé.
- Il contient tous les "phénomènes" pouvant expliquer les sources d'aléa (qui ne sont pas explicables...).
- En pratique, l'espace d'arrivée est généralement suffisant.

Définition

Une **variable aléatoire réelle** est une application

$$X : (\Omega, \mathcal{A}) \rightarrow (\mathbb{R}, \mathcal{B}(\mathbb{R}))$$

telle que

$$\forall B \in \mathcal{B}(\mathbb{R}), X^{-1}(B) \in \mathcal{A}.$$

- Lors de la modélisation statistique, l'espace Ω n'est généralement jamais caractérisé.
- Il contient tous les "phénomènes" pouvant expliquer les sources d'aléa (qui ne sont pas explicables...).
- En pratique, l'espace d'arrivée est généralement suffisant.

Loi de probabilité

Etant donnée \mathbf{P} une probabilité sur (Ω, \mathcal{A}) et X une variable aléatoire réelle définie sur Ω , on appelle loi de probabilité de X la mesure \mathbf{P}_X définie par

$$\mathbf{P}_X(B) = \mathbf{P}(X^{-1}(B)) = \mathbf{P}(X \in B) = \mathbf{P}(\{\omega \in \Omega : X(\omega) \in B\}) \quad \forall B \in \mathcal{B}(\mathbb{R}).$$

Une loi de probabilité est caractérisée par

- sa fonction de répartition : $F_X(x) = \mathbf{P}(X \leq x)$.
- sa densité : $f_X : \mathbb{R} \rightarrow \mathbb{R}^+$ telle que $\forall B \in \mathcal{B}(\mathbb{R})$

$$\mathbf{P}_X(B) = \int_B f_X(x) dx.$$

Loi de probabilité

Etant donnée \mathbf{P} une probabilité sur (Ω, \mathcal{A}) et X une variable aléatoire réelle définie sur Ω , on appelle loi de probabilité de X la mesure \mathbf{P}_X définie par

$$\mathbf{P}_X(B) = \mathbf{P}(X^{-1}(B)) = \mathbf{P}(X \in B) = \mathbf{P}(\{\omega \in \Omega : X(\omega) \in B\}) \quad \forall B \in \mathcal{B}(\mathbb{R}).$$

Une loi de probabilité est caractérisée par

- sa fonction de répartition : $F_X(x) = \mathbf{P}(X \leq x)$.
- sa densité : $f_X : \mathbb{R} \rightarrow \mathbb{R}^+$ telle que $\forall B \in \mathcal{B}(\mathbb{R})$

$$\mathbf{P}_X(B) = \int_B f_X(x) dx.$$

Un modèle pour l'exemple 1

- On note $x_i = 1$ si le $i^{\text{ème}}$ patient a guéri, 0 sinon.
- On peut supposer que x_i est la réalisation d'une variable aléatoire X_i de loi de bernoulli de paramètre p_0 .
- Si les individus sont choisis de manière **indépendante** et ont tous la **même probabilité de guérir** (ce qui peut revenir à dire qu'ils en sont au même stade de la pathologie), il est alors raisonnable de supposer que les variables aléatoires X_1, \dots, X_n sont indépendantes et de même loi (i.i.d.).

On dit que X_1, \dots, X_n est un **n -échantillon** de variables aléatoires indépendantes de même loi $B(p_0)$.

Un modèle pour l'exemple 1

- On note $x_i = 1$ si le $i^{\text{ème}}$ patient a guéri, 0 sinon.
- On peut supposer que x_i est la réalisation d'une variable aléatoire X_i de loi de bernoulli de paramètre p_0 .
- Si les individus sont choisis de manière **indépendante** et ont tous la **même probabilité de guérir** (ce qui peut revenir à dire qu'ils en sont au même stade de la pathologie), il est alors raisonnable de supposer que les variables aléatoires X_1, \dots, X_n sont indépendantes et de même loi (i.i.d.).

On dit que X_1, \dots, X_n est un n -échantillon de variables aléatoires indépendantes de même loi $B(p_0)$.

Un modèle pour l'exemple 1

- On note $x_i = 1$ si le $i^{\text{ème}}$ patient a guéri, 0 sinon.
- On peut supposer que x_i est la réalisation d'une variable aléatoire X_i de loi de bernoulli de paramètre p_0 .
- Si les individus sont choisis de manière **indépendante** et ont tous la **même probabilité de guérir** (ce qui peut revenir à dire qu'ils en sont au même stade de la pathologie), il est alors raisonnable de supposer que les variables aléatoires X_1, \dots, X_n sont indépendantes et de même loi (i.i.d.).

On dit que X_1, \dots, X_n est un **n -échantillon** de variables aléatoires indépendantes de même loi $B(p_0)$.

Modèle

On appelle **modèle statistique** tout triplet $(\mathcal{H}, \mathcal{A}, \mathcal{P})$ où

- \mathcal{H} est l'espace des observations (l'ensemble de tous les résultats possibles de l'expérience) ;
- \mathcal{A} est une tribu sur \mathcal{H} ;
- \mathcal{P} est une famille de probabilités définie sur $(\mathcal{H}, \mathcal{A})$.

Le problème du statisticien

- n variables aléatoires i.i.d. X_1, \dots, X_n de loi \mathbf{P} .
- Trouver une famille de lois \mathcal{P} susceptible de contenir \mathbf{P} .
- Trouver dans \mathcal{P} une loi qui soit **la plus proche** de \mathbf{P}

Modèle

On appelle **modèle statistique** tout triplet $(\mathcal{H}, \mathcal{A}, \mathcal{P})$ où

- \mathcal{H} est l'espace des observations (l'ensemble de tous les résultats possibles de l'expérience) ;
- \mathcal{A} est une tribu sur \mathcal{H} ;
- \mathcal{P} est une famille de probabilités définie sur $(\mathcal{H}, \mathcal{A})$.

Le problème du statisticien

- n variables aléatoires i.i.d. X_1, \dots, X_n de loi \mathbf{P} .
- Trouver une famille de lois \mathcal{P} susceptible de contenir \mathbf{P} .
- Trouver dans \mathcal{P} une loi qui soit **la plus proche** de \mathbf{P}

Exemples

	\mathcal{H}	\mathcal{A}	\mathcal{P}
Exemple 1	$\{0, 1\}$	$\mathcal{P}(\{0, 1\})$	$\{B(p), p \in [0, 1]\}$
Exemple 2	\mathbb{N}	$\mathcal{P}(\mathbb{N})$	$\{\mathcal{P}(\lambda), \lambda > 0\}$
Exemple 3	\mathbb{R}	$\mathcal{B}(\mathbb{R})$	$\{\mathcal{N}(\mu, \sigma^2), \mu \in \mathbb{R}, \sigma \in \mathbb{R}^+\}$

Définition

- Si $\mathcal{P} = \{\mathbf{P}_\theta, \theta \in \Theta\}$ où $\Theta \in \mathbb{R}^d$ alors on parle de **modèle paramétrique** et Θ est l'espace des paramètres.
- Si $\mathcal{P} = \{\mathbf{P}, \mathbf{P} \in \mathcal{F}\}$ où \mathcal{F} est de dimension infinie, on parle de **modèle non paramétrique**.

Exemple : modèle de densité

- $\mathcal{P} = \{N(\mu, \sigma^2), (\mu, \sigma^2) \in \mathbb{R} \times \mathbb{R}^+\}$ est un modèle paramétrique.
- $\mathcal{P} = \{\text{densités } f \text{ 2 fois dérivables}\}$ est un modèle non paramétrique.

Le problème sera d'estimer (μ, σ^2) ou f à partir de l'échantillon X_1, \dots, X_n .

Définition

- Si $\mathcal{P} = \{\mathbf{P}_\theta, \theta \in \Theta\}$ où $\Theta \in \mathbb{R}^d$ alors on parle de **modèle paramétrique** et Θ est l'espace des paramètres.
- Si $\mathcal{P} = \{\mathbf{P}, \mathbf{P} \in \mathcal{F}\}$ où \mathcal{F} est de dimension infinie, on parle de **modèle non paramétrique**.

Exemple : modèle de densité

- $\mathcal{P} = \{N(\mu, \sigma^2), (\mu, \sigma^2) \in \mathbb{R} \times \mathbb{R}^+\}$ est un modèle paramétrique.
- $\mathcal{P} = \{\text{densités } f \text{ 2 fois dérivables}\}$ est un modèle non paramétrique.

Le problème sera d'estimer (μ, σ^2) ou f à partir de l'échantillon X_1, \dots, X_n .

Définition

- Si $\mathcal{P} = \{\mathbf{P}_\theta, \theta \in \Theta\}$ où $\Theta \in \mathbb{R}^d$ alors on parle de **modèle paramétrique** et Θ est l'espace des paramètres.
- Si $\mathcal{P} = \{\mathbf{P}, \mathbf{P} \in \mathcal{F}\}$ où \mathcal{F} est de dimension infinie, on parle de **modèle non paramétrique**.

Exemple : modèle de densité

- $\mathcal{P} = \{\mathcal{N}(\mu, \sigma^2), (\mu, \sigma^2) \in \mathbb{R} \times \mathbb{R}^+\}$ est un modèle paramétrique.
- $\mathcal{P} = \{\text{densités } f \text{ 2 fois dérivables}\}$ est un modèle non paramétrique.

Le problème sera d'estimer (μ, σ^2) ou f à partir de l'échantillon X_1, \dots, X_n .

Définition

- Si $\mathcal{P} = \{\mathbf{P}_\theta, \theta \in \Theta\}$ où $\Theta \in \mathbb{R}^d$ alors on parle de **modèle paramétrique** et Θ est l'espace des paramètres.
- Si $\mathcal{P} = \{\mathbf{P}, \mathbf{P} \in \mathcal{F}\}$ où \mathcal{F} est de dimension infinie, on parle de **modèle non paramétrique**.

Exemple : modèle de densité

- $\mathcal{P} = \{\mathcal{N}(\mu, \sigma^2), (\mu, \sigma^2) \in \mathbb{R} \times \mathbb{R}^+\}$ est un modèle paramétrique.
- $\mathcal{P} = \{\text{densités } f \text{ 2 fois dérivables}\}$ est un modèle non paramétrique.

Le problème sera d'estimer (μ, σ^2) ou f à partir de l'échantillon X_1, \dots, X_n .

Modèle de régression

Modèle de régression

- On cherche à expliquer une variable Y par p variables explicatives $\mathbf{X}_1, \dots, \mathbf{X}_p$. On dispose d'un n échantillon i.i.d. $(X_i, Y_i), i = 1, \dots, n$.

Modèle linéaire (paramétrique)

- On pose

$$Y = \beta_0 + \beta_1 \mathbf{X}_1 + \dots + \beta_p \mathbf{X}_p + \varepsilon \quad \text{où} \quad \varepsilon \sim \mathcal{N}(0, \sigma^2).$$

- Le problème est d'estimer $\beta = (\beta_0, \dots, \beta_p) \in \mathbb{R}^{p+1}$ à l'aide de $(X_1, Y_1), \dots, (X_n, Y_n)$.

Un modèle non paramétrique

- On pose

$$Y = m(\mathbf{X}_1, \dots, \mathbf{X}_p) + \varepsilon$$

où $m : \mathbb{R}^p \rightarrow \mathbb{R}$ est une fonction continue.

- Le problème est d'estimer m à l'aide de $(X_1, Y_1), \dots, (X_n, Y_n)$.

Modèle de régression

- On cherche à expliquer une variable Y par p variables explicatives $\mathbf{X}_1, \dots, \mathbf{X}_p$. On dispose d'un n échantillon i.i.d. $(X_i, Y_i), i = 1, \dots, n$.

Modèle linéaire (paramétrique)

- On pose

$$Y = \beta_0 + \beta_1 \mathbf{X}_1 + \dots + \beta_p \mathbf{X}_p + \varepsilon \quad \text{où} \quad \varepsilon \sim \mathcal{N}(0, \sigma^2).$$

- Le problème est d'estimer $\beta = (\beta_0, \dots, \beta_p) \in \mathbb{R}^{p+1}$ à l'aide de $(X_1, Y_1), \dots, (X_n, Y_n)$.

Un modèle non paramétrique

- On pose

$$Y = m(\mathbf{X}_1, \dots, \mathbf{X}_p) + \varepsilon$$

où $m : \mathbb{R}^p \rightarrow \mathbb{R}$ est une fonction continue.

- Le problème est d'estimer m à l'aide de $(X_1, Y_1), \dots, (X_n, Y_n)$.

Modèle de régression

- On cherche à expliquer une variable Y par p variables explicatives $\mathbf{X}_1, \dots, \mathbf{X}_p$. On dispose d'un n échantillon i.i.d. $(X_i, Y_i), i = 1, \dots, n$.

Modèle linéaire (paramétrique)

- On pose

$$Y = \beta_0 + \beta_1 \mathbf{X}_1 + \dots + \beta_p \mathbf{X}_p + \varepsilon \quad \text{où} \quad \varepsilon \sim \mathcal{N}(0, \sigma^2).$$

- Le problème est d'estimer $\beta = (\beta_0, \dots, \beta_p) \in \mathbb{R}^{p+1}$ à l'aide de $(X_1, Y_1), \dots, (X_n, Y_n)$.

Un modèle non paramétrique

- On pose

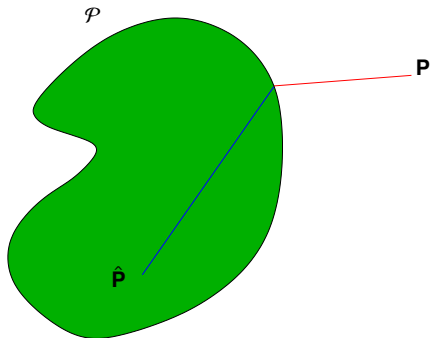
$$Y = m(\mathbf{X}_1, \dots, \mathbf{X}_p) + \varepsilon$$

où $m : \mathbb{R}^p \rightarrow \mathbb{R}$ est une fonction continue.

- Le problème est d'estimer m à l'aide de $(X_1, Y_1), \dots, (X_n, Y_n)$.

2 types d'erreur

- Poser un modèle revient à choisir une famille de loi candidates pour reconstruire la loi des données \mathbf{P} .

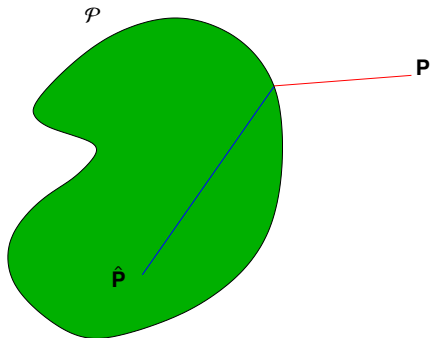


On distingue deux types d'erreurs :

- **Erreur d'estimation** : erreur commise par le choix d'une loi dans \mathcal{P} par rapport au meilleur choix.
- **Erreur d'approximation** : erreur commise par le choix de \mathcal{P} .

2 types d'erreur

- Poser un modèle revient à choisir une famille de loi candidates pour reconstruire la loi des données \mathbf{P} .



On distingue deux types d'erreurs :

- **Erreur d'estimation** : erreur commise par le choix d'une loi dans \mathcal{P} par rapport au meilleur choix.
- **Erreur d'approximation** : erreur commise par le choix de \mathcal{P} .

- 1 On récolte n observations (n valeurs) x_1, \dots, x_n qui sont le résultats de n expériences aléatoires indépendantes.
- 2 **Modélisation** : on suppose que les n valeurs sont des réalisations de n variables aléatoires indépendantes X_1, \dots, X_n et de même loi \mathbf{P}_{θ_0} .
- 3 **Estimation** : chercher dans le modèle une loi $\mathbf{P}_{\hat{\theta}}$ qui soit le plus proche possible de $\mathbf{P}_{\theta_0} \implies$ chercher un **estimateur** $\hat{\theta}$ de θ_0 .
- 4 **"Validation" de modèle** : on revient en arrière et on tente de vérifier si l'hypothèse de l'étape 2 est raisonnable (test d'adéquation, etc...)

- 1 On récolte n observations (n valeurs) x_1, \dots, x_n qui sont le résultats de n expériences aléatoires indépendantes.
- 2 **Modélisation** : on **suppose** que les n valeurs sont des réalisations de n variables aléatoires indépendantes X_1, \dots, X_n et de même loi \mathbf{P}_{θ_0} .
- 3 **Estimation** : chercher dans le modèle une loi $\mathbf{P}_{\hat{\theta}}$ qui soit le plus proche possible de $\mathbf{P}_{\theta_0} \implies$ chercher un **estimateur** $\hat{\theta}$ de θ_0 .
- 4 **"Validation" de modèle** : on revient en arrière et on tente de vérifier si l'hypothèse de l'étape 2 est raisonnable (test d'adéquation, etc...)

- 1 On récolte n observations (n valeurs) x_1, \dots, x_n qui sont le résultats de n expériences aléatoires indépendantes.
- 2 **Modélisation** : on **suppose** que les n valeurs sont des réalisations de n variables aléatoires indépendantes X_1, \dots, X_n et de même loi \mathbf{P}_{θ_0} .
- 3 **Estimation** : chercher dans le modèle une loi $\mathbf{P}_{\hat{\theta}}$ qui soit le plus proche possible de $\mathbf{P}_{\theta_0} \implies$ chercher un **estimateur** $\hat{\theta}$ de θ_0 .
- 4 **"Validation" de modèle** : on revient en arrière et on tente de vérifier si l'hypothèse de l'étape 2 est raisonnable (test d'adéquation, etc...)

- 1 On récolte n observations (n valeurs) x_1, \dots, x_n qui sont le résultats de n expériences aléatoires indépendantes.
- 2 **Modélisation** : on **suppose** que les n valeurs sont des réalisations de n variables aléatoires indépendantes X_1, \dots, X_n et de même loi \mathbf{P}_{θ_0} .
- 3 **Estimation** : chercher dans le modèle une loi $\mathbf{P}_{\hat{\theta}}$ qui soit le plus proche possible de $\mathbf{P}_{\theta_0} \implies$ chercher un **estimateur** $\hat{\theta}$ de θ_0 .
- 4 **"Validation" de modèle** : on revient en arrière et on tente de vérifier si l'hypothèse de l'étape 2 est raisonnable (test d'adéquation, etc...)

Rappels sur le modèle de régression linéaire

Le problème de régression

- On cherche à expliquer une variable Y par p variables $\mathbf{X}_1, \dots, \mathbf{X}_p$.
- Il s'agit de trouver une fonction $m : \mathbb{R}^p \rightarrow \mathbb{R}$ telle que $Y \approx m(\mathbf{X}_1, \dots, \mathbf{X}_p)$.
- Sauf cas (très) particulier, le lien n'est jamais "parfait"

$$Y = m(\mathbf{X}_1, \dots, \mathbf{X}_p) + \varepsilon.$$

Modèle de régression

- Poser un modèle de régression revient à supposer que la fonction m appartient à un certain espace \mathcal{M} .
- Le problème du statisticien sera alors de trouver la "meilleure" fonction dans \mathcal{M} à l'aide d'un n -échantillon i.i.d. $(X_1, Y_1), \dots, (X_n, Y_n)$.

Le problème de régression

- On cherche à expliquer une variable Y par p variables $\mathbf{X}_1, \dots, \mathbf{X}_p$.
- Il s'agit de trouver une fonction $m : \mathbb{R}^p \rightarrow \mathbb{R}$ telle que $Y \approx m(\mathbf{X}_1, \dots, \mathbf{X}_p)$.
- Sauf cas (très) particulier, le lien n'est jamais "parfait"

$$Y = m(\mathbf{X}_1, \dots, \mathbf{X}_p) + \varepsilon.$$

Modèle de régression

- Poser un modèle de régression revient à supposer que la fonction m appartient à un certain espace \mathcal{M} .
- Le problème du statisticien sera alors de trouver la "meilleure" fonction dans \mathcal{M} à l'aide d'un n -échantillon i.i.d. $(X_1, Y_1), \dots, (X_n, Y_n)$.

Le problème de régression

- On cherche à expliquer une variable Y par p variables $\mathbf{X}_1, \dots, \mathbf{X}_p$.
- Il s'agit de trouver une fonction $m : \mathbb{R}^p \rightarrow \mathbb{R}$ telle que $Y \approx m(\mathbf{X}_1, \dots, \mathbf{X}_p)$.
- Sauf cas (très) particulier, le lien n'est jamais "parfait"

$$Y = m(\mathbf{X}_1, \dots, \mathbf{X}_p) + \varepsilon.$$

Modèle de régression

- Poser un modèle de régression revient à supposer que la fonction m appartient à un certain espace \mathcal{M} .
- Le problème du statisticien sera alors de trouver la "meilleure" fonction dans \mathcal{M} à l'aide d'un n -échantillon i.i.d. $(X_1, Y_1), \dots, (X_n, Y_n)$.

Le problème de régression

- On cherche à expliquer une variable Y par p variables $\mathbf{X}_1, \dots, \mathbf{X}_p$.
- Il s'agit de trouver une fonction $m : \mathbb{R}^p \rightarrow \mathbb{R}$ telle que $Y \approx m(\mathbf{X}_1, \dots, \mathbf{X}_p)$.
- Sauf cas (très) particulier, le lien n'est jamais "parfait"

$$Y = m(\mathbf{X}_1, \dots, \mathbf{X}_p) + \varepsilon.$$

Modèle de régression

- Poser un modèle de régression revient à supposer que la fonction m appartient à un certain espace \mathcal{M} .
- Le problème du statisticien sera alors de trouver la "meilleure" fonction dans \mathcal{M} à l'aide d'un n -échantillon i.i.d. $(X_1, Y_1), \dots, (X_n, Y_n)$.

Modèle non paramétrique

- L'espace \mathcal{M} est de dimension infinie.
- **Exemple** : On pose $Y = m(\mathbf{X}_1, \dots, \mathbf{X}_p) + \varepsilon$ où m appartient à l'espace des fonctions continues.

Modèle paramétrique

- L'espace \mathcal{M} est de dimension finie.
- **Exemple** : on suppose que la fonction m est linéaire

$$Y = \beta_0 + \beta_1 \mathbf{X}_1 + \dots + \beta_p \mathbf{X}_p + \varepsilon.$$

Le problème est alors d'estimer $\beta = (\beta_0, \beta_1, \dots, \beta_p)$ à l'aide de $(X_1, Y_1), \dots, (X_n, Y_n)$.

- C'est le modèle de **régression linéaire**.

Modèle non paramétrique

- L'espace \mathcal{M} est de dimension infinie.
- **Exemple** : On pose $Y = m(\mathbf{X}_1, \dots, \mathbf{X}_p) + \varepsilon$ où m appartient à l'espace des fonctions continues.

Modèle paramétrique

- L'espace \mathcal{M} est de dimension finie.
- **Exemple** : on suppose que la fonction m est linéaire

$$Y = \beta_0 + \beta_1 \mathbf{X}_1 + \dots + \beta_p \mathbf{X}_p + \varepsilon.$$

Le problème est alors d'estimer $\beta = (\beta_0, \beta_1, \dots, \beta_p)$ à l'aide de $(X_1, Y_1), \dots, (X_n, Y_n)$.

- C'est le modèle de **régression linéaire**.

- On cherche à **expliquer** ou à **prédire** la concentration en ozone.
- On dispose de $n = 112$ observations de la concentration en ozone ainsi que de 12 autres variables susceptibles d'expliquer cette concentration :
 - Température relevée à différents moments de la journée.
 - Indice de nébulosité relevé à différents moments de la journée.
 - Direction du vent.
 - Pluie.

Question

Comment expliquer (modéliser) la concentration en ozone à l'aide de toutes ces variables ?

- On cherche à **expliquer** ou à **prédire** la concentration en ozone.
- On dispose de $n = 112$ observations de la concentration en ozone ainsi que de 12 autres variables susceptibles d'expliquer cette concentration :
 - Température relevée à différents moments de la journée.
 - Indice de nébulosité relevé à différents moments de la journée.
 - Direction du vent.
 - Pluie.

Question

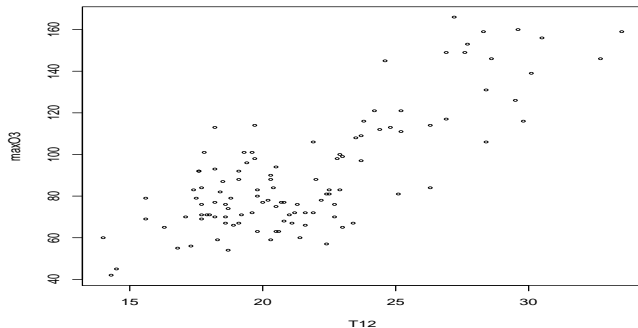
Comment expliquer (modéliser) la concentration en ozone à l'aide de toutes ces variables ?

Ozone en fonction de la température à 12h ?

MaxO3	87	82	92	114	94	80	...
T12	18.5	18.4	17.6	19.7	20.5	19.8	...

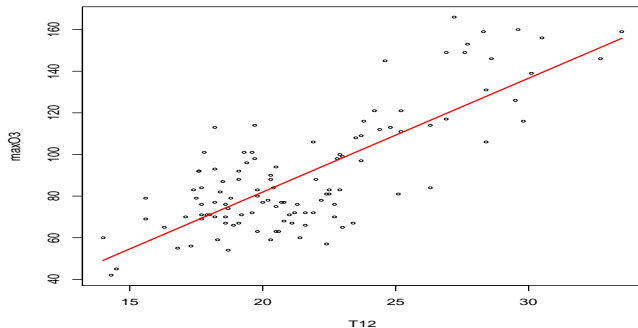
Ozone en fonction de la température à 12h ?

MaxO3	87	82	92	114	94	80	...
T12	18.5	18.4	17.6	19.7	20.5	19.8	...



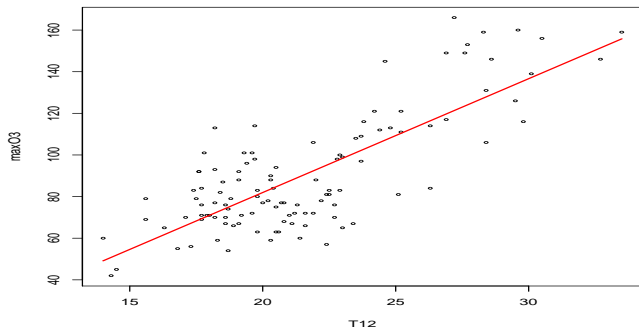
Ozone en fonction de la température à 12h ?

MaxO3	87	82	92	114	94	80	...
T12	18.5	18.4	17.6	19.7	20.5	19.8	...



Ozone en fonction de la température à 12h ?

MaxO3	87	82	92	114	94	80	...
T12	18.5	18.4	17.6	19.7	20.5	19.8	...



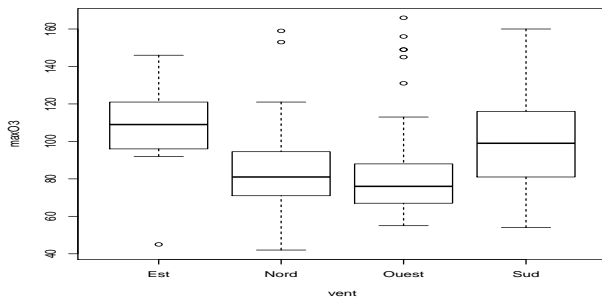
Comment ajuster le nuage de points ?

Ozone en fonction du vent ?

MaxO3	87	82	92	114	94	80	...
Vent	Nord	Nord	Est	Nord	Ouest	Ouest	...

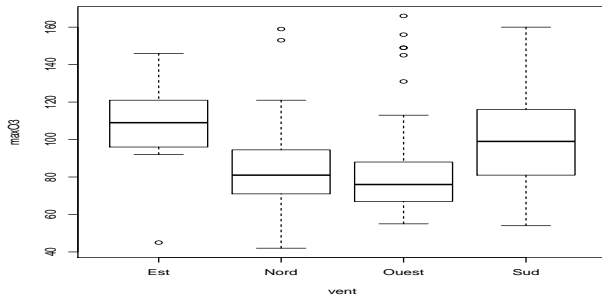
Ozone en fonction du vent ?

MaxO3	87	82	92	114	94	80	...
Vent	Nord	Nord	Est	Nord	Ouest	Ouest	...



Ozone en fonction du vent ?

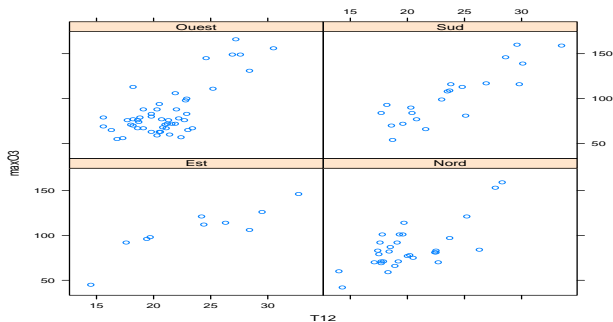
MaxO3	87	82	92	114	94	80	...
Vent	Nord	Nord	Est	Nord	Ouest	Ouest	...



$$\text{MaxO3} \approx \alpha_1 \mathbf{1}_{\text{Vent=est}} + \dots + \alpha_4 \mathbf{1}_{\text{Vent=sud}}$$

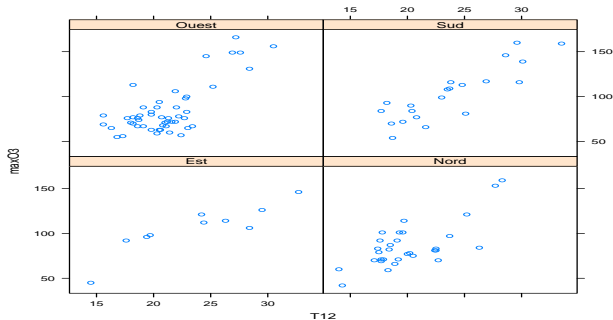
$$\alpha_j = ???$$

Ozone en fonction de la température à 12h et du vent ?



$$\max O_3 \approx \begin{cases} \beta_{01} + \beta_{11} T_{12} & \text{si vent=est} \\ \vdots & \vdots \\ \beta_{04} + \beta_{14} T_{12} & \text{si vent=ouest} \end{cases}$$

Ozone en fonction de la température à 12h et du vent ?



$$\max O_3 \approx \begin{cases} \beta_{01} + \beta_{11} T_{12} & \text{si vent=est} \\ \vdots & \vdots \\ \beta_{04} + \beta_{14} T_{12} & \text{si vent=ouest} \end{cases}$$

- Généralisation

$$\text{maxO3} \approx \beta_0 + \beta_1 V_1 + \dots + \beta_{12} V_{12}$$

Questions

- Comment calculer (ou plutôt **estimer**) les paramètres β_j ?
- Le modèle avec les 12 variables est-il "meilleur" que des modèles avec moins de variables ?
- Comment trouver le "meilleur" sous-groupe de variables ?

- Généralisation

$$\text{maxO3} \approx \beta_0 + \beta_1 V_1 + \dots + \beta_{12} V_{12}$$

Questions

- Comment calculer (ou plutôt **estimer**) les paramètres β_j ?
- Le modèle avec les 12 variables est-il "meilleur" que des modèles avec moins de variables ?
- Comment trouver le "meilleur" sous-groupe de variables ?

- Y : variable (aléatoire) à expliquer à valeurs dans \mathbb{R} .
- X_1, \dots, X_p : p variables explicatives à valeurs dans \mathbb{R} .
- n observations $(x_1, Y_1), \dots, (x_n, Y_n)$ avec $x_i = (x_{i1}, \dots, x_{ip})$.

Le modèle de régression linéaire multiple

Le modèle s'écrit :

$$Y_i = \beta_0 + \beta_1 x_{i1} + \dots + \beta_p x_{ip} + \varepsilon_i$$

où les erreurs aléatoires ε_i sont i.i.d. de loi $\mathcal{N}(0, \sigma^2)$.

- Y : variable (aléatoire) à expliquer à valeurs dans \mathbb{R} .
- X_1, \dots, X_p : p variables explicatives à valeurs dans \mathbb{R} .
- n observations $(x_1, Y_1), \dots, (x_n, Y_n)$ avec $x_i = (x_{i1}, \dots, x_{ip})$.

Le modèle de régression linéaire multiple

Le modèle s'écrit :

$$Y_i = \beta_0 + \beta_1 x_{i1} + \dots + \beta_p x_{ip} + \varepsilon_i$$

où les erreurs aléatoires ε_i sont i.i.d. de loi $\mathcal{N}(0, \sigma^2)$.

- On note

$$\mathbf{Y} = \begin{pmatrix} Y_1 \\ \vdots \\ Y_n \end{pmatrix}, \quad \mathbf{X} = \begin{pmatrix} 1 & x_{11} & \dots & x_{1p} \\ \vdots & \vdots & & \vdots \\ 1 & x_{n1} & \dots & x_{np} \end{pmatrix}, \quad \boldsymbol{\beta} = \begin{pmatrix} \beta_0 \\ \vdots \\ \beta_p \end{pmatrix}, \quad \boldsymbol{\varepsilon} = \begin{pmatrix} \varepsilon_1 \\ \vdots \\ \varepsilon_n \end{pmatrix}$$

Ecriture matricielle

Le modèle se réécrit

$$\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$$

où $\boldsymbol{\varepsilon} \sim \mathcal{N}(0, \sigma^2 \mathbf{I}_n)$.

- On note

$$\mathbf{Y} = \begin{pmatrix} Y_1 \\ \vdots \\ Y_n \end{pmatrix}, \quad \mathbf{X} = \begin{pmatrix} 1 & x_{11} & \dots & x_{1p} \\ \vdots & \vdots & & \vdots \\ 1 & x_{n1} & \dots & x_{np} \end{pmatrix}, \quad \boldsymbol{\beta} = \begin{pmatrix} \beta_0 \\ \vdots \\ \beta_p \end{pmatrix}, \quad \boldsymbol{\varepsilon} = \begin{pmatrix} \varepsilon_1 \\ \vdots \\ \varepsilon_n \end{pmatrix}$$

Ecriture matricielle

Le modèle se réécrit

$$\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$$

où $\boldsymbol{\varepsilon} \sim \mathcal{N}(0, \sigma^2 \mathbf{I}_n)$.

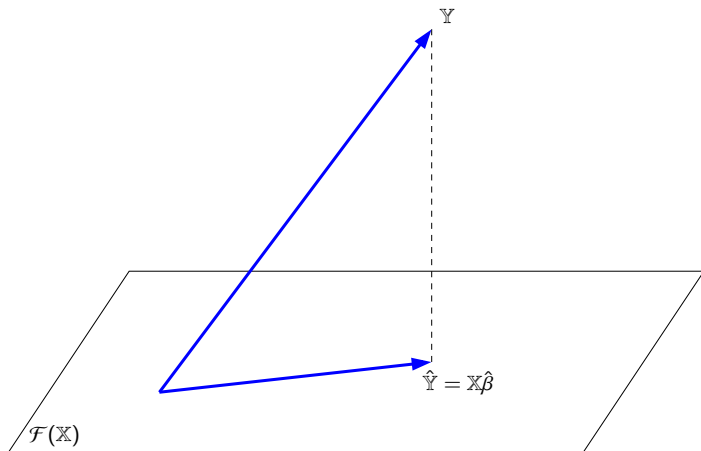
Définition

On appelle **estimateur des moindres carrés** $\hat{\beta}$ de β la statistique suivante :

$$\hat{\beta} = \underset{\beta_0, \dots, \beta_p}{\operatorname{argmin}} \sum_{i=1}^n (Y_i - \beta_0 - \beta_1 x_{i1} - \dots - \beta_p x_{ip})^2 = \underset{\beta \in \mathbb{R}^{p+1}}{\operatorname{argmin}} \|\mathbf{Y} - \mathbf{X}\beta\|^2.$$

- On note $\mathcal{F}(\mathbf{X})$ le s.e.v. de \mathbb{R}^n de dimension $p + 1$ engendré par les $p + 1$ colonnes de \mathbf{X} .
- Chercher l'estimateur des moindres carrés revient à minimiser la distance entre $\mathbf{Y} \in \mathbb{R}^n$ et $\mathcal{F}(\mathbf{X})$.

Représentation géométrique



- On déduit que $\mathbb{X}\hat{\beta}$ est le projeté orthogonal de \mathbb{Y} sur $\mathcal{F}(\mathbb{X})$:

$$\mathbb{X}\hat{\beta} = \mathbf{P}_{\mathcal{F}(\mathbb{X})}(\mathbb{Y}) = \mathbb{X}(\mathbb{X}'\mathbb{X})^{-1}\mathbb{X}'\mathbb{Y}.$$

Théorème

Si la matrice \mathbb{X} est de plein rang, l'estimateur des MC est donné par :

$$\hat{\beta} = (\mathbb{X}'\mathbb{X})^{-1}\mathbb{X}'\mathbb{Y}.$$

- On déduit que $\mathbb{X}\hat{\beta}$ est le projeté orthogonal de \mathbb{Y} sur $\mathcal{F}(\mathbb{X})$:

$$\mathbb{X}\hat{\beta} = \mathbf{P}_{\mathcal{F}(\mathbb{X})}(\mathbb{Y}) = \mathbb{X}(\mathbb{X}'\mathbb{X})^{-1}\mathbb{X}'\mathbb{Y}.$$

Théorème

Si la matrice \mathbb{X} est de plein rang, l'estimateur des MC est donné par :

$$\hat{\beta} = (\mathbb{X}'\mathbb{X})^{-1}\mathbb{X}'\mathbb{Y}.$$

- On déduit que $\mathbb{X}\hat{\beta}$ est le projeté orthogonal de \mathbb{Y} sur $\mathcal{F}(\mathbb{X})$:

$$\mathbb{X}\hat{\beta} = \mathbf{P}_{\mathcal{F}(\mathbb{X})}(\mathbb{Y}) = \mathbb{X}(\mathbb{X}'\mathbb{X})^{-1}\mathbb{X}'\mathbb{Y}.$$

Théorème

Si la matrice \mathbb{X} est de plein rang, l'estimateur des MC est donné par :

$$\hat{\beta} = (\mathbb{X}'\mathbb{X})^{-1}\mathbb{X}'\mathbb{Y}.$$

Propriété

- 1 $\hat{\beta}$ est un estimateur sans biais de β .
- 2 La matrice de variance-covariance de $\hat{\beta}$ est donnée par

$$\mathbf{V}(\hat{\beta}) = \sigma^2(\mathbf{X}'\mathbf{X})^{-1}.$$

- 3 $\hat{\beta}$ est VUMSB.

- Soit $\hat{\varepsilon} = \mathbf{Y} - \hat{\mathbf{Y}}$ le vecteur des résidus et $\widehat{\sigma^2}$ l'estimateur de σ^2 défini par

$$\widehat{\sigma^2} = \frac{\|\hat{\varepsilon}\|^2}{n - (p + 1)}.$$

Proposition

- 1 $\hat{\beta}$ est un vecteur gaussien d'espérance β et de matrice de variance-covariance $\sigma^2(\mathbf{X}'\mathbf{X})^{-1}$.
- 2 $(n - (p + 1))\frac{\widehat{\sigma^2}}{\sigma^2} \sim \chi^2_{n-(p+1)}$.
- 3 $\hat{\beta}$ et $\widehat{\sigma^2}$ sont indépendantes.

- Soit $\hat{\varepsilon} = \mathbb{Y} - \hat{\mathbb{Y}}$ le vecteur des résidus et $\widehat{\sigma^2}$ l'estimateur de σ^2 défini par

$$\widehat{\sigma^2} = \frac{\|\hat{\varepsilon}\|^2}{n - (p + 1)}.$$

Proposition

- 1 $\hat{\beta}$ est un vecteur gaussien d'espérance β et de matrice de variance-covariance $\sigma^2(\mathbb{X}'\mathbb{X})^{-1}$.
- 2 $(n - (p + 1))\frac{\widehat{\sigma^2}}{\sigma^2} \sim \chi^2_{n-(p+1)}$.
- 3 $\hat{\beta}$ et $\widehat{\sigma^2}$ sont indépendantes.

- Soit $\hat{\varepsilon} = \mathbf{Y} - \hat{\mathbf{Y}}$ le vecteur des résidus et $\widehat{\sigma^2}$ l'estimateur de σ^2 défini par

$$\widehat{\sigma^2} = \frac{\|\hat{\varepsilon}\|^2}{n - (p + 1)}.$$

Proposition

- 1 $\hat{\beta}$ est un vecteur gaussien d'espérance β et de matrice de variance-covariance $\sigma^2(\mathbf{X}'\mathbf{X})^{-1}$.
- 2 $(n - (p + 1))\frac{\widehat{\sigma^2}}{\sigma^2} \sim \chi^2_{n-(p+1)}$.
- 3 $\hat{\beta}$ et $\widehat{\sigma^2}$ sont indépendantes.

Corollaire

On note $\widehat{\sigma}_j^2 = \widehat{\sigma}^2 [\mathbb{X}'\mathbb{X}]_{jj}^{-1}$ pour $j = 0, \dots, p$. On a

$$\forall j = 0, \dots, p, \quad \frac{\hat{\beta}_j - \beta_j}{\widehat{\sigma}_j} \sim \mathcal{T}(n - (p + 1)).$$

On déduit de ce corollaire :

- des intervalles de confiance de niveau $1 - \alpha$ pour β_j .
- des procédures de test pour des hypothèses du genre $H_0 : \beta_j = 0$ contre $H_1 : \beta_j \neq 0$.

Corollaire

On note $\widehat{\sigma}_j^2 = \widehat{\sigma}^2 [\mathbb{X}'\mathbb{X}]_{jj}^{-1}$ pour $j = 0, \dots, p$. On a

$$\forall j = 0, \dots, p, \quad \frac{\hat{\beta}_j - \beta_j}{\widehat{\sigma}_j} \sim \mathcal{T}(n - (p + 1)).$$

On déduit de ce corollaire :

- des intervalles de confiance de niveau $1 - \alpha$ pour β_j .
- des procédures de test pour des hypothèses du genre $H_0 : \beta_j = 0$ contre $H_1 : \beta_j \neq 0$.

- On dispose d'une nouvelle observation $x_{n+1} = (x_{n+1,1}, \dots, x_{n+1,p})$ et on souhaite prédire la valeur $y_{n+1} = x'_{n+1}\beta$ associée à cette nouvelle observation.

- Un estimateur (naturel) de y_{n+1} est $\hat{y}_{n+1} = x'_{n+1}\hat{\beta}$.
- Un intervalle de confiance de niveau $1 - \alpha$ pour y_{n+1} est donné par

$$\left[\hat{y}_{n+1} \pm t_{n-(p+1)}(\alpha/2)\hat{\sigma} \sqrt{x'_{n+1}(\mathbb{X}'\mathbb{X})^{-1}x_{n+1} + 1} \right].$$

- On dispose d'une nouvelle observation $x_{n+1} = (x_{n+1,1}, \dots, x_{n+1,p})$ et on souhaite prédire la valeur $y_{n+1} = x'_{n+1}\beta$ associée à cette nouvelle observation.

- Un estimateur (naturel) de y_{n+1} est $\hat{y}_{n+1} = x'_{n+1}\hat{\beta}$.
- Un intervalle de confiance de niveau $1 - \alpha$ pour y_{n+1} est donné par

$$\left[\hat{y}_{n+1} \pm t_{n-(p+1)}(\alpha/2)\hat{\sigma} \sqrt{x'_{n+1}(\mathbb{X}'\mathbb{X})^{-1}x_{n+1} + 1} \right].$$

- On dispose d'une nouvelle observation $x_{n+1} = (x_{n+1,1}, \dots, x_{n+1,p})$ et on souhaite prédire la valeur $y_{n+1} = x'_{n+1}\beta$ associée à cette nouvelle observation.

- Un estimateur (naturel) de y_{n+1} est $\hat{y}_{n+1} = x'_{n+1}\hat{\beta}$.
- Un intervalle de confiance de niveau $1 - \alpha$ pour y_{n+1} est donné par

$$\left[\hat{y}_{n+1} \pm t_{n-(p+1)}(\alpha/2)\hat{\sigma} \sqrt{x'_{n+1}(\mathbb{X}'\mathbb{X})^{-1}x_{n+1} + 1} \right].$$

Exemple de l'ozone

- On considère le modèle de régression multiple :

$$\text{MaxO3} = \beta_0 + \beta_1 T_{12} + \beta_2 T_{15} + \beta_3 N_{12} + \beta_4 V_{12} + \beta_5 \text{MaxO3v} + \varepsilon.$$

```
> reg.multi <- lm(maxO3~T12+T15+Ne12+Vx12+maxO3v,data=donnees)
> summary(reg.multi)
```

```
Call:
lm(formula = maxO3 ~ T12 + T15 + Ne12 + Vx12 + maxO3v, data = donnees)
```

```
Residuals:
    Min       1Q   Median       3Q      Max
-54.216  -9.446  -0.896   8.007  41.186
```

```
Coefficients:
            Estimate Std. Error t value Pr(>|t|)
(Intercept)  3.04498   13.01591   0.234   0.8155
T12          2.47747    1.09257   2.268   0.0254 *
T15          0.63177    0.96382   0.655   0.5136
Ne12        -1.83560    0.89439  -2.052   0.0426 *
Vx12         1.33295    0.58168   2.292   0.0239 *
maxO3v       0.34215    0.05989   5.713 1.03e-07 ***
```

```
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

```
Residual standard error: 14.58 on 106 degrees of freedom
Multiple R-squared:  0.7444, Adjusted R-squared:  0.7324
F-statistic: 61.75 on 5 and 106 DF,  p-value: < 2.2e-16
```

"Validation" du modèle

- Ajuster un modèle, trouver des estimateurs est un problème relativement "simple".
- Le travail difficile est de trouver un bon modèle, ou encore le meilleur modèle (ce travail est difficile car la notion de meilleur modèle n'existe pas).
- Il est donc nécessaire de trouver des procédures automatiques de choix de modèles (méthodes pas à pas utilisant un critère de type AIC, BIC, régression lasso etc...)
- Puis de vérifier que les hypothèses effectuées (normalité, linéarité) sont raisonnables (analyse des résidus, tests d'adéquation...).

"Validation" du modèle

- Ajuster un modèle, trouver des estimateurs est un problème relativement "simple".
- Le travail difficile est de trouver un bon modèle, ou encore le meilleur modèle (ce travail est difficile car la notion de meilleur modèle n'existe pas).
- Il est donc nécessaire de trouver des procédures automatiques de choix de modèles (méthodes pas à pas utilisant un critère de type AIC, BIC, régression lasso etc...)
- Puis de vérifier que les hypothèses effectuées (normalité, linéarité) sont raisonnables (analyse des résidus, tests d'adéquation...).

"Validation" du modèle

- Ajuster un modèle, trouver des estimateurs est un problème relativement "simple".
- Le travail difficile est de trouver un bon modèle, ou encore le meilleur modèle (ce travail est difficile car la notion de meilleur modèle n'existe pas).
- Il est donc nécessaire de trouver des procédures automatiques de choix de modèles (méthodes pas à pas utilisant un critère de type AIC, BIC, régression lasso etc...)
- Puis de vérifier que les hypothèses effectuées (normalité, linéarité) sont raisonnables (analyse des résidus, tests d'adéquation...).

"Validation" du modèle

- Ajuster un modèle, trouver des estimateurs est un problème relativement "simple".
- Le travail difficile est de trouver un bon modèle, ou encore le meilleur modèle (ce travail est difficile car la notion de meilleur modèle n'existe pas).
- Il est donc nécessaire de trouver des procédures automatiques de choix de modèles (méthodes pas à pas utilisant un critère de type AIC, BIC, régression lasso etc...)
- Puis de vérifier que les hypothèses effectuées (normalité, linéarité) sont raisonnables (analyse des résidus, tests d'adéquation...).

"Validation" du modèle

- Ajuster un modèle, trouver des estimateurs est un problème relativement "simple".
- Le travail difficile est de trouver un bon modèle, ou encore le meilleur modèle (ce travail est difficile car la notion de meilleur modèle n'existe pas).
- Il est donc nécessaire de trouver des procédures automatiques de choix de modèles (méthodes pas à pas utilisant un critère de type AIC, BIC, régression lasso etc...)
- Puis de vérifier que les hypothèses effectuées (normalité, linéarité) sont raisonnables (analyse des résidus, tests d'adéquation...).

- Le modèle linéaire

$$Y_i = x_i' \beta + \varepsilon_i, \quad \varepsilon_i \text{ i.i.d de loi } \mathcal{N}(0, \sigma^2)$$

- peut se réécrire pour $i = 1, \dots, n$

$$\mathcal{L}(Y_i) = \mathcal{N}(x_i' \beta, \sigma^2).$$

Interprétation

Au point x_i la loi de Y est une gaussienne $\mathcal{N}(x_i' \beta, \sigma^2)$.

- Le modèle linéaire

$$Y_i = x_i' \beta + \varepsilon_i, \quad \varepsilon_i \text{ i.i.d de loi } \mathcal{N}(0, \sigma^2)$$

- peut se réécrire pour $i = 1, \dots, n$

$$\mathcal{L}(Y_i) = \mathcal{N}(x_i' \beta, \sigma^2).$$

Interprétation

Au point x_i la loi de Y est une gaussienne $\mathcal{N}(x_i' \beta, \sigma^2)$.

- Le modèle linéaire

$$Y_i = x_i' \beta + \varepsilon_i, \quad \varepsilon_i \text{ i.i.d de loi } \mathcal{N}(0, \sigma^2)$$

- peut se réécrire pour $i = 1, \dots, n$

$$\mathcal{L}(Y_i) = \mathcal{N}(x_i' \beta, \sigma^2).$$

Interprétation

Au point x_i la loi de Y est une gaussienne $\mathcal{N}(x_i' \beta, \sigma^2)$.

- On peut alors calculer la (log)-vraisemblance du modèle

$$\mathcal{L}(y_1, \dots, y_n; \beta) = \frac{n}{2} \log(\sigma^2) - \frac{n}{2} \log(2\pi) - \frac{1}{2\sigma^2} \|\mathbb{Y} - \mathbb{X}\beta\|^2.$$

- **Conclusion** : l'estimateur du maximum de vraisemblance $\hat{\beta}_{MV}$ coïncide avec l'estimateur des moindres carrés $\hat{\beta}$.

Remarque

- Si les variables explicatives sont **aléatoires**, ce n'est plus la loi de Y_i qui est modélisée mais celle de Y_i sachant $X_i = x_i$

$$\mathcal{L}(Y_i | X_i = x_i) = \mathcal{N}(x_i' \beta, \sigma^2).$$

- Plus généralement, lorsque les variables explicatives sont supposées **aléatoires** (économétrie), poser un modèle de régression revient à "mettre" **une famille de loi sur Y sachant $X = x$** .

- On peut alors calculer la (log)-vraisemblance du modèle

$$\mathcal{L}(y_1, \dots, y_n; \beta) = \frac{n}{2} \log(\sigma^2) - \frac{n}{2} \log(2\pi) - \frac{1}{2\sigma^2} \|\mathbb{Y} - \mathbb{X}\beta\|^2.$$

- **Conclusion** : l'estimateur du maximum de vraisemblance $\hat{\beta}_{MV}$ coïncide avec l'estimateur des moindres carrés $\hat{\beta}$.

Remarque

- Si les variables explicatives sont **aléatoires**, ce n'est plus la loi de Y_i qui est modélisée mais celle de Y_i sachant $X_i = x_i$

$$\mathcal{L}(Y_i | X_i = x_i) = \mathcal{N}(x_i' \beta, \sigma^2).$$

- Plus généralement, lorsque les variables explicatives sont supposées **aléatoires** (économétrie), poser un modèle de régression revient à "mettre" **une famille de loi sur Y sachant $X = x$** .

- On peut alors calculer la (log)-vraisemblance du modèle

$$\mathcal{L}(y_1, \dots, y_n; \beta) = \frac{n}{2} \log(\sigma^2) - \frac{n}{2} \log(2\pi) - \frac{1}{2\sigma^2} \|\mathbb{Y} - \mathbb{X}\beta\|^2.$$

- **Conclusion** : l'estimateur du maximum de vraisemblance $\hat{\beta}_{MV}$ coïncide avec l'estimateur des moindres carrés $\hat{\beta}$.

Remarque

- Si les variables explicatives sont **aléatoires**, ce n'est plus la loi de Y_i qui est modélisée mais celle de Y_i sachant $X_i = x_i$

$$\mathcal{L}(Y_i | X_i = x_i) = \mathcal{N}(x_i' \beta, \sigma^2).$$

- Plus généralement, lorsque les variables explicatives sont supposées **aléatoires** (économétrie), poser un modèle de régression revient à "mettre" **une famille de loi sur Y sachant $X = x$** .

Introduction au modèle de régression logistique

Exemples

Détection de clients à risque

- Une chaîne de magasin a mis en place une carte de crédit.
- Elle dispose d'un historique de 145 clients dont 40 ont connu des défauts de paiement.
- Elle connaît également d'autres caractéristiques de ces clients (sexe, taux d'enttement, revenus mensuels, dépenses effectuées sur certaines gammes de produit...)

Question

Comment prédire si un nouveau client connaîtra des défauts de paiement ?

- Une chaîne de magasin a mis en place une carte de crédit.
- Elle dispose d'un historique de 145 clients dont 40 ont connu des défauts de paiement.
- Elle connaît également d'autres caractéristiques de ces clients (sexe, taux d'endettement, revenus mensuels, dépenses effectuées sur certaines gammes de produit...)

Question

Comment prédire si un nouveau client connaîtra des défauts de paiement ?

- On a mesuré sur 150 iris de 3 espèces différentes (Setosa, Versicolor, Virginica) les quantités suivantes :
 - Longueur et largeur des pétales
 - Longueur et largeur des sépales

Question

Comment identifier l'espèce d'un iris à partir de ces 4 caractéristiques ?

- On a mesuré sur 150 iris de 3 espèces différentes (Setosa, Versicolor, Virginica) les quantités suivantes :
 - Longueur et largeur des pétales
 - Longueur et largeur des sépales

Question

Comment identifier l'espèce d'un iris à partir de ces 4 caractéristiques ?

- Sur 4 601 mails, on a pu identifier 1813 spams.
- On a également mesuré sur chacun de ces mails la présence ou absence de 57 mots.

Question

Peut-on construire à partir de ces données une méthode de détection automatique de spam ?

- Sur 4 601 mails, on a pu identifier 1813 spams.
- On a également mesuré sur chacun de ces mails la présence ou absence de 57 mots.

Question

Peut-on construire à partir de ces données une méthode de détection automatique de spam ?

Pathologie concernant les artères coronaires

- **Problème** : étudier la présence d'une pathologie concernant les artères coronaires en fonction de l'âge des individus.
- **Données** : on dispose d'un échantillon de taille 100 sur lequel on a mesuré les variables :
 - chd qui vaut 1 si la pathologie est présente, 0 sinon ;
 - age qui correspond à l'âge de l'individu.

```
> artere[1:5,]  
   age agrp chd  
1.  20    1   0  
2.  23    1   0  
3.  24    1   0  
4.  25    1   0  
5.  25    1   1
```

Pathologie concernant les artères coronaires

- **Problème** : étudier la présence d'une pathologie concernant les artères coronaires en fonction de l'âge des individus.
- **Données** : on dispose d'un échantillon de taille 100 sur lequel on a mesuré les variables :
 - chd qui vaut 1 si la pathologie est présente, 0 sinon ;
 - age qui correspond à l'âge de l'individu.

```
> artere[1:5,]  
  age agrp chd  
1.  20    1   0  
2.  23    1   0  
3.  24    1   0  
4.  25    1   0  
5.  25    1   1
```

Représentation du problème

- Tous ces problèmes peuvent être appréhendés dans un contexte de **régression** : on cherche à expliquer une variable Y par d'autres variables X_1, \dots, X_p :

Y	X
Défaut de paiement	caractéristiques du client
Espèce de l'iris	Longueur, largeur pétales et sépales
Spam	présence/absence de mots

- La variable à expliquer n'est plus quantitative mais **qualitative**.
- On parle de problème de **discrimination** ou **classification supervisée**.

Représentation du problème

- Tous ces problèmes peuvent être appréhendés dans un contexte de **régression** : on cherche à expliquer une variable Y par d'autres variables X_1, \dots, X_p :

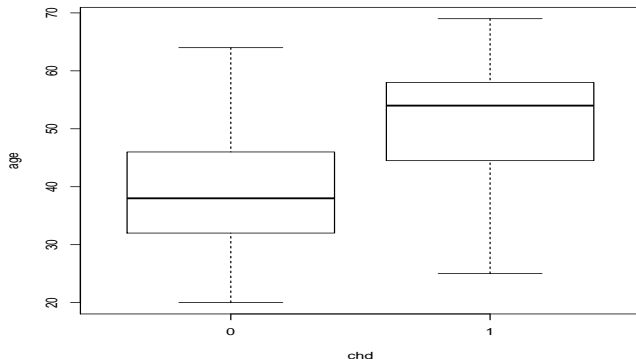
Y	X
Défaut de paiement	caractéristiques du client
Espèce de l'iris	Longueur, largeur pétales et sépales
Spam	présence/absence de mots

- La variable à expliquer n'est plus quantitative mais **qualitative**.
- On parle de problème de **discrimination** ou **classification supervisée**.

Régression logistique simple

Boxplot

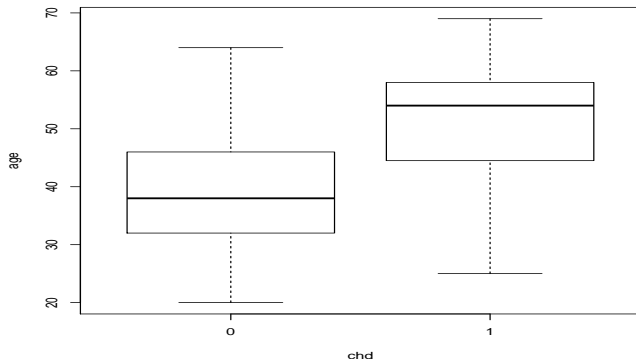
```
> plot(age~chd, data=artere)
```



Il semble que la maladie a plus de chance d'être présente chez les personnes âgées.

Boxplot

```
> plot(age~chd, data=artere)
```



Il semble que la maladie a plus de chance d'être présente chez les personnes âgées.

Question

Comment expliquer la relation entre la maladie et l'âge ?

- On désigne par
 - Y la variable aléatoire qui prend pour valeur 1 si l'individu est atteint, 0 sinon.
 - X la variable (aléatoire) qui correspond à l'âge de l'individu.

Le problème consiste ainsi à tenter de **quantifier la relation** entre Y et X à partir des données, c'est-à-dire d'un **échantillon i.i.d** $(X_1, Y_1), \dots, (X_n, Y_n)$ de taille $n = 100$.

Question

Comment expliquer la relation entre la maladie et l'âge ?

- On désigne par
 - Y la variable aléatoire qui prend pour valeur 1 si l'individu est atteint, 0 sinon.
 - X la variable (aléatoire) qui correspond à l'âge de l'individu.

Le problème consiste ainsi à tenter de **quantifier la relation** entre Y et X à partir des données, c'est-à-dire d'un **échantillon i.i.d** $(X_1, Y_1), \dots, (X_n, Y_n)$ de taille $n = 100$.

Question

Comment expliquer la relation entre la maladie et l'âge ?

- On désigne par
 - Y la variable aléatoire qui prend pour valeur 1 si l'individu est atteint, 0 sinon.
 - X la variable (aléatoire) qui correspond à l'âge de l'individu.

Le problème consiste ainsi à tenter de **quantifier la relation** entre Y et X à partir des données, c'est-à-dire d'un **échantillon i.i.d** $(X_1, Y_1), \dots, (X_n, Y_n)$ de taille $n = 100$.

- On se base sur le **modèle linéaire**.
- On suppose que les deux variables Y et X sont liées par une relation de la forme

$$Y = \beta_0 + \beta_1 X + \varepsilon \quad (1)$$

où $\beta_0 \in \mathbb{R}$ et $\beta_1 \in \mathbb{R}$ sont les **paramètres inconnus** du modèle et ε est une variable aléatoire de loi $\mathcal{N}(0, \sigma^2)$.

Problème

La variable Y est ici **qualitative**, l'écriture (1) n'a donc aucun sens.

⇒ **mauvaise idée**

- On se base sur le **modèle linéaire**.
- On suppose que les deux variables Y et X sont liées par une relation de la forme

$$Y = \beta_0 + \beta_1 X + \varepsilon \quad (1)$$

où $\beta_0 \in \mathbb{R}$ et $\beta_1 \in \mathbb{R}$ sont les **paramètres inconnus** du modèle et ε est une variable aléatoire de loi $\mathcal{N}(0, \sigma^2)$.

Problème

La variable Y est ici **qualitative**, l'écriture (1) n'a donc aucun sens.

⇒ **mauvaise idée**

- On se base sur le **modèle linéaire**.
- On suppose que les deux variables Y et X sont liées par une relation de la forme

$$Y = \beta_0 + \beta_1 X + \varepsilon \quad (1)$$

où $\beta_0 \in \mathbb{R}$ et $\beta_1 \in \mathbb{R}$ sont les **paramètres inconnus** du modèle et ε est une variable aléatoire de loi $\mathcal{N}(0, \sigma^2)$.

Problème

La variable Y est ici **qualitative**, l'écriture (1) n'a donc aucun sens.

⇒ **mauvaise idée**

- Chercher à expliquer Y par X revient à chercher de l'information sur la **loi de probabilité de Y sachant X** .
- En effet, le modèle de régression linéaire peut se réécrire en caractérisant la loi de $Y|X = x$ par la loi $\mathcal{N}(\beta_0 + \beta_1 x, \sigma^2)$.

Idée

- Étendre cette caractérisation à notre contexte (où la variable à expliquer est binaire).
- Une loi candidate naturelle pour la variable $Y|X = x$ est la loi de **Bernoulli**.

- Chercher à expliquer Y par X revient à chercher de l'information sur la **loi de probabilité de Y sachant X** .
- En effet, le modèle de régression linéaire peut se réécrire en caractérisant la loi de $Y|X = x$ par la loi $\mathcal{N}(\beta_0 + \beta_1 x, \sigma^2)$.

Idée

- Etendre cette caractérisation à notre contexte (où la variable à expliquer est binaire).
- Une loi candidate naturelle pour la variable $Y|X = x$ est la loi de **Bernoulli**.

- Chercher à expliquer Y par X revient à chercher de l'information sur la loi de probabilité de Y sachant X .
- En effet, le modèle de régression linéaire peut se réécrire en caractérisant la loi de $Y|X = x$ par la loi $\mathcal{N}(\beta_0 + \beta_1 x, \sigma^2)$.

Idée

- Étendre cette caractérisation à notre contexte (où la variable à expliquer est binaire).
- Une loi candidate naturelle pour la variable $Y|X = x$ est la loi de Bernoulli.

- Chercher à expliquer Y par X revient à chercher de l'information sur la loi de probabilité de Y sachant X .
- En effet, le modèle de régression linéaire peut se réécrire en caractérisant la loi de $Y|X = x$ par la loi $\mathcal{N}(\beta_0 + \beta_1 x, \sigma^2)$.

Idée

- Étendre cette caractérisation à notre contexte (où la variable à expliquer est binaire).
- Une loi candidate naturelle pour la variable $Y|X = x$ est la loi de Bernoulli.

Loi de Bernoulli

- On va ainsi caractériser la loi de $Y|X = x$ par la loi de Bernoulli.
- Cette loi dépend d'un paramètre

$$p(x) = \mathbf{P}(Y = 1|X = x).$$

- Sachant $X = x$, on a donc

$$Y = \begin{cases} 1 & \text{avec probabilité } p(x) \\ 0 & \text{avec probabilité } 1 - p(x) \end{cases}$$

La modélisation

Il reste maintenant à caractériser la probabilité $p(x)$.

- On va ainsi caractériser la loi de $Y|X = x$ par la loi de Bernoulli.
- Cette loi dépend d'un paramètre

$$p(x) = \mathbf{P}(Y = 1|X = x).$$

- Sachant $X = x$, on a donc

$$Y = \begin{cases} 1 & \text{avec probabilité } p(x) \\ 0 & \text{avec probabilité } 1 - p(x) \end{cases}$$

La modélisation

Il reste maintenant à caractériser la probabilité $p(x)$.

- On va ainsi caractériser la loi de $Y|X = x$ par la loi de Bernoulli.
- Cette loi dépend d'un paramètre

$$p(x) = \mathbf{P}(Y = 1|X = x).$$

- Sachant $X = x$, on a donc

$$Y = \begin{cases} 1 & \text{avec probabilité } p(x) \\ 0 & \text{avec probabilité } 1 - p(x) \end{cases}$$

La modélisation

Il reste maintenant à caractériser la probabilité $p(x)$.

- On va ainsi caractériser la loi de $Y|X = x$ par la loi de Bernoulli.
- Cette loi dépend d'un paramètre

$$p(x) = \mathbf{P}(Y = 1|X = x).$$

- Sachant $X = x$, on a donc

$$Y = \begin{cases} 1 & \text{avec probabilité } p(x) \\ 0 & \text{avec probabilité } 1 - p(x) \end{cases}$$

La modélisation

Il reste maintenant à caractériser la probabilité $p(x)$.

- Là encore, on peut se baser sur le **modèle linéaire** et proposer

$$p(x) = \beta_0 + \beta_1 x.$$

- Cette écriture n'est pas satisfaisante. En effet
 - $p(x) \in [0, 1]$ tandis que $\beta_0 + \beta_1 x \in \mathbb{R}$.
 - **Idée** : trouver une transformation φ de $p(x)$ telle que $\varphi(p(x))$ prenne ses valeurs dans \mathbb{R} .

- Là encore, on peut se baser sur le **modèle linéaire** et proposer

$$p(x) = \beta_0 + \beta_1 x.$$

- Cette écriture n'est pas satisfaisante. En effet
 - $p(x) \in [0, 1]$ tandis que $\beta_0 + \beta_1 x \in \mathbb{R}$.
 - **Idée** : trouver une transformation φ de $p(x)$ telle que $\varphi(p(x))$ prenne ses valeurs dans \mathbb{R} .

- Là encore, on peut se baser sur le **modèle linéaire** et proposer

$$p(x) = \beta_0 + \beta_1 x.$$

- Cette écriture n'est pas satisfaisante. En effet
 - $p(x) \in [0, 1]$ tandis que $\beta_0 + \beta_1 x \in \mathbb{R}$.
 - **Idée** : trouver une transformation φ de $p(x)$ telle que $\varphi(p(x))$ prenne ses valeurs dans \mathbb{R} .

- Là encore, on peut se baser sur le **modèle linéaire** et proposer

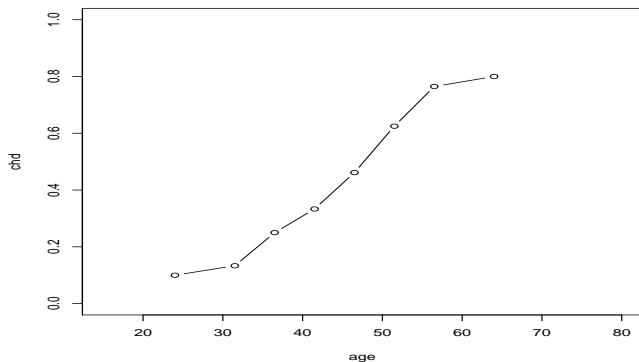
$$p(x) = \beta_0 + \beta_1 x.$$

- Cette écriture n'est pas satisfaisante. En effet
 - $p(x) \in [0, 1]$ tandis que $\beta_0 + \beta_1 x \in \mathbb{R}$.
 - **Idée** : trouver une transformation φ de $p(x)$ telle que $\varphi(p(x))$ prenne ses valeurs dans \mathbb{R} .

- On revient sur l'exemple du chd et on représente les **fréquences cumulées** d'apparition de la maladie en fonction de l'âge :

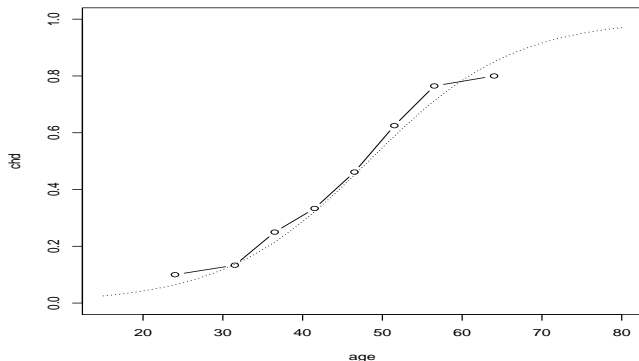
Transformation de $p(x)$

- On revient sur l'exemple du chd et on représente les **fréquences cumulées** d'apparition de la maladie en fonction de l'âge :



Transformation de $p(x)$

- On revient sur l'exemple du chd et on représente les **fréquences cumulées** d'apparition de la maladie en fonction de l'âge :



Trouver une **transformation** de $p(x)$ qui ajuste ce nuage de points.

Le modèle de régression logistique

- Il propose de **modéliser la probabilité** $p(x)$ selon

$$p(x) = \frac{\exp(\beta_0 + \beta_1 x)}{1 + \exp(\beta_0 + \beta_1 x)}.$$

- On peut réécrire

$$\text{logit } p(x) = \log\left(\frac{p(x)}{1 - p(x)}\right) = \beta_0 + \beta_1 x.$$

Le modèle de régression logistique

Le **modèle de régression logistique** consiste donc à caractériser la loi de $Y|X = x$ par une loi de **Bernoulli** de paramètre $p(x)$ tel que

$$\text{logit } p(x) = \log\left(\frac{p(x)}{1 - p(x)}\right) = \beta_0 + \beta_1 x.$$

Le modèle de régression logistique

- Il propose de **modéliser la probabilité** $p(x)$ selon

$$p(x) = \frac{\exp(\beta_0 + \beta_1 x)}{1 + \exp(\beta_0 + \beta_1 x)}.$$

- On peut réécrire

$$\text{logit } p(x) = \log\left(\frac{p(x)}{1 - p(x)}\right) = \beta_0 + \beta_1 x.$$

Le modèle de régression logistique

Le **modèle de régression logistique** consiste donc à caractériser la loi de $Y|X = x$ par une loi de **Bernoulli** de paramètre $p(x)$ tel que

$$\text{logit } p(x) = \log\left(\frac{p(x)}{1 - p(x)}\right) = \beta_0 + \beta_1 x.$$

Exemple sur R

```
> model <- glm(chd~age,data=artere,family=binomial)
> model
```

```
Call:  glm(formula = chd ~ age, family = binomial, data = artere)
```

```
Coefficients:
(Intercept)          age
   -5.3095         0.1109
```

```
Degrees of Freedom: 99 Total (i.e. Null);  98 Residual
Null Deviance:      136.7
Residual Deviance: 107.4  AIC: 111.4
```

- La fonction **glm** renvoie les estimations de β_0 et β_1 .
- On peut ainsi avoir une estimation de la **probabilité d'avoir une maladie pour un individu de 30 ans** :

$$\hat{p}(x = 30) = \frac{\exp(-5.3095 + 0.1109 * 30)}{1 + \exp(-5.3095 + 0.1109 * 30)} \approx 0.12.$$

Exemple sur R

```
> model <- glm(chd~age,data=artere,family=binomial)
> model
```

```
Call:  glm(formula = chd ~ age, family = binomial, data = artere)
```

```
Coefficients:
(Intercept)          age
   -5.3095         0.1109
```

```
Degrees of Freedom: 99 Total (i.e. Null);  98 Residual
Null Deviance:      136.7
Residual Deviance: 107.4  AIC: 111.4
```

- La fonction **glm** renvoie les estimations de β_0 et β_1 .
- On peut ainsi avoir une estimation de la probabilité d'avoir une maladie pour un individu de 30 ans :

$$\hat{p}(x = 30) = \frac{\exp(-5.3095 + 0.1109 * 30)}{1 + \exp(-5.3095 + 0.1109 * 30)} \approx 0.12.$$

Exemple sur R

```
> model <- glm(chd~age,data=artere,family=binomial)
> model
```

```
Call:  glm(formula = chd ~ age, family = binomial, data = artere)
```

```
Coefficients:
(Intercept)          age
   -5.3095         0.1109
```

```
Degrees of Freedom: 99 Total (i.e. Null);  98 Residual
Null Deviance:      136.7
Residual Deviance: 107.4  AIC: 111.4
```

- La fonction **glm** renvoie les estimations de β_0 et β_1 .
- On peut ainsi avoir une estimation de la **probabilité d'avoir une maladie pour un individu de 30 ans** :

$$\hat{p}(x = 30) = \frac{\exp(-5.3095 + 0.1109 * 30)}{1 + \exp(-5.3095 + 0.1109 * 30)} \approx 0.12.$$

Exemple sur R

```
> model <- glm(chd~age,data=artere,family=binomial)
> model
```

```
Call:  glm(formula = chd ~ age, family = binomial, data = artere)
```

```
Coefficients:
(Intercept)          age
   -5.3095         0.1109
```

```
Degrees of Freedom: 99 Total (i.e. Null);  98 Residual
Null Deviance:      136.7
Residual Deviance: 107.4  AIC: 111.4
```

- La fonction **glm** renvoie les estimations de β_0 et β_1 .
- On peut ainsi avoir une estimation de la **probabilité d'avoir une maladie pour un individu de 30 ans** :

$$\hat{p}(x = 30) = \frac{\exp(-5.3095 + 0.1109 * 30)}{1 + \exp(-5.3095 + 0.1109 * 30)} \approx 0.12.$$

Le modèle linéaire généralisé

Introduction

Le modèle logistique est un glm

- Le modèle de **régression logistique** s'ajuste sur R avec la fonction **glm**.
- Le modèle de régression logistique appartient à la famille des **modèles linéaires généralisés**.
- C'est pourquoi il faut spécifier l'argument `family=binomial` lorsque l'on veut faire une régression logistique.

Le modèle logistique est un glm

- Le modèle de **régression logistique** s'ajuste sur R avec la fonction **glm**.
- Le modèle de régression logistique appartient à la famille des **modèles linéaires généralisés**.
- C'est pourquoi il faut spécifier l'argument `family=binomial` lorsque l'on veut faire une régression logistique.

Le modèle logistique est un glm

- Le modèle de **régression logistique** s'ajuste sur R avec la fonction **glm**.
- Le modèle de régression logistique appartient à la famille des **modèles linéaires généralisés**.
- C'est pourquoi il faut spécifier l'argument **family=binomial** lorsque l'on veut faire une régression logistique.

Le modèle linéaire est un GLM

- Le modèle de **régression linéaire** s'ajuste sur R avec la fonction **lm** :

```
> Y <- rnorm(50)
> X <- runif(50)
> lm(Y~X)
```

```
Coefficients:
(Intercept)          X
    0.4245      -0.8547
```

- Mais aussi avec la fonction **glm** :

```
> glm(Y~X, family=gaussian)
```

```
Coefficients:
(Intercept)          X
    0.4245      -0.8547
```

Conclusion

Le modèle linéaire appartient également à la famille des **modèles linéaires généralisés**.

Le modèle linéaire est un GLM

- Le modèle de **régression linéaire** s'ajuste sur R avec la fonction **lm** :

```
> Y <- rnorm(50)
> X <- runif(50)
> lm(Y~X)
```

```
Coefficients:
(Intercept)          X
    0.4245      -0.8547
```

- Mais aussi avec la fonction **glm** :

```
> glm(Y~X, family=gaussian)
```

```
Coefficients:
(Intercept)          X
    0.4245      -0.8547
```

Conclusion

Le modèle linéaire appartient également à la famille des **modèles linéaires généralisés**.

Le modèle linéaire est un GLM

- Le modèle de **régression linéaire** s'ajuste sur R avec la fonction **lm** :

```
> Y <- rnorm(50)
> X <- runif(50)
> lm(Y~X)
```

```
Coefficients:
(Intercept)          X
      0.4245      -0.8547
```

- Mais aussi avec la fonction **glm** :

```
> glm(Y~X, family=gaussian)
```

```
Coefficients:
(Intercept)          X
      0.4245      -0.8547
```

Conclusion

Le modèle linéaire appartient également à la famille des **modèles linéaires généralisés**.

2 étapes identiques

- Les modèles linéaires et logistiques sont construits selon le même protocole en 2 étapes :

① Choix de la loi conditionnelle de $Y|X = x$:

- Gaussienne pour le modèle linéaire ;
- Bernoulli pour le modèle logistique.

② Choix d'une transformation g de l'espérance conditionnelle $E[Y|X = x]$:

- Logistique

$$g(E[Y|X = x]) = g(p(x)) = \text{logit } p(x) = x'\beta$$

- Linéaire

$$g(E[Y|X = x]) = x'\beta.$$

2 étapes identiques

- Les modèles linéaires et logistiques sont construits selon le même protocole en 2 étapes :

① Choix de la loi conditionnelle de $Y|X = x$:

- Gaussienne pour le modèle linéaire ;
- Bernoulli pour le modèle logistique.

② Choix d'une transformation g de l'espérance conditionnelle $E[Y|X = x]$:

- Logistique

$$g(E[Y|X = x]) = g(p(x)) = \text{logit } p(x) = x'\beta$$

- Linéaire

$$g(E[Y|X = x]) = x'\beta.$$

2 étapes identiques

- Les modèles linéaires et logistiques sont construits selon le même protocole en 2 étapes :

① Choix de la loi conditionnelle de $Y|X = x$:

- Gaussienne pour le modèle linéaire ;
- Bernoulli pour le modèle logistique.

② Choix d'une transformation g de l'espérance conditionnelle $\mathbf{E}[Y|X = x]$:

- **Logistique**

$$g(\mathbf{E}[Y|X = x]) = g(p(x)) = \text{logit } p(x) = x'\beta$$

- **Linéaire**

$$g(\mathbf{E}[Y|X = x]) = x'\beta.$$

Définitions

Rappel : Famille exponentielle

Définition

Une loi de probabilité \mathbf{P} appartient à une famille de lois de type exponentielle $\{\mathbf{P}_\theta\}_{\theta \in \mathbb{R}^p}$ si il existe une mesure dominant μ (Lebesgue ou mesure de comptage le plus souvent) telle que les lois \mathbf{P}_θ admettent pour densité par rapport à ν

$$f_\theta(y) = c(\theta)h(y) \exp\left(\sum_{j=1}^p \alpha_j(\theta) T_j(y)\right)$$

où T_1, \dots, T_p sont des fonctions réelles mesurables.

Exemple : loi de Bernoulli

La loi de Bernoulli de paramètre p admet pour densité (par rapport à la mesure de comptage)

$$f_p(y) = (1 - p) \exp(y \log(p/(1 - p))).$$

Rappel : Famille exponentielle

Définition

Une loi de probabilité \mathbf{P} appartient à une famille de lois de type exponentielle $\{\mathbf{P}_\theta\}_{\theta \in \mathbb{R}^p}$ si il existe une mesure dominante μ (Lebesgue ou mesure de comptage le plus souvent) telle que les lois \mathbf{P}_θ admettent pour densité par rapport à ν

$$f_\theta(y) = c(\theta)h(y) \exp\left(\sum_{j=1}^p \alpha_j(\theta) T_j(y)\right)$$

où T_1, \dots, T_p sont des fonctions réelles mesurables.

Exemple : loi de Bernoulli

La loi de Bernoulli de paramètre p admet pour densité (par rapport à la mesure de comptage)

$$f_p(y) = (1 - p) \exp(y \log(p/(1 - p))).$$

Rappel : Famille exponentielle

Définition

Une loi de probabilité \mathbf{P} appartient à une famille de lois de type exponentielle $\{\mathbf{P}_\theta\}_{\theta \in \mathbb{R}^p}$ si il existe une mesure dominante μ (Lebesgue ou mesure de comptage le plus souvent) telle que les lois \mathbf{P}_θ admettent pour densité par rapport à ν

$$f_\theta(y) = c(\theta)h(y) \exp\left(\sum_{j=1}^p \alpha_j(\theta) T_j(y)\right)$$

où T_1, \dots, T_p sont des fonctions réelles mesurables.

Exemple : loi de Bernoulli

La loi de Bernoulli de paramètre p admet pour densité (par rapport à la mesure de comptage)

$$f_p(y) = (1 - p) \exp(y \log(p/(1 - p))).$$

- On se place dans un contexte de **régression** : on cherche à expliquer une variable Y par p variables explicatives X_1, \dots, X_p .
- On dispose d'un n -échantillon $(x_1, Y_1), \dots, (x_n, Y_n)$ où les $x_i = (x_{i1}, \dots, x_{ip})$ sont supposées **fixes** et les Y_i sont des variables aléatoires réelles **indépendantes**.

- On se place dans un contexte de **régression** : on cherche à expliquer une variable Y par p variables explicatives X_1, \dots, X_p .
- On dispose d'un n -échantillon $(x_1, Y_1), \dots, (x_n, Y_n)$ où les $x_i = (x_{i1}, \dots, x_{ip})$ sont supposées **fixes** et les Y_i sont des variables aléatoires réelles **indépendantes**.

Modèle linéaire généralisé : GLM

Un modèle linéaire généralisé est constitué de **3 composantes** :

- 1 **Composante aléatoire** : la loi de probabilité de la réponse Y_i appartient à la famille exponentielle et est de la forme

$$f_{\alpha_i}(y_i) = \exp\left(\frac{\alpha_i y_i - b(\alpha_i)}{a(\phi)} + c(y_i, \phi)\right)$$

où a , b et c sont des fonctions spécifiées en fonction du type de la famille exponentielle.

- 2 **Composante déterministe** :

$$\eta(x_i) = \beta_0 + \beta_1 x_{i1} + \dots + \beta_p x_{ip}$$

et précise quels sont les **prédicteurs** (on peut y inclure des transformations des prédicteurs, des interactions...).

- 3 **Lien** : spécifie le **lien entre les deux composantes**, plus précisément le lien entre l'espérance de Y_i et la composante déterministe : $g(\mathbf{E}[Y_i]) = \eta(x_i) = \beta_0 + \beta_1 x_{i1} + \dots + \beta_p x_{ip}$ où g est une fonction inversible appelée **fonction de lien**.

Modèle linéaire généralisé : GLM

Un modèle linéaire généralisé est constitué de **3 composantes** :

- 1 **Composante aléatoire** : la loi de probabilité de la réponse Y_i appartient à la famille exponentielle et est de la forme

$$f_{\alpha_i}(y_i) = \exp\left(\frac{\alpha_i y_i - b(\alpha_i)}{a(\phi)} + c(y_i, \phi)\right)$$

où a , b et c sont des fonctions spécifiées en fonction du type de la famille exponentielle.

- 2 **Composante déterministe** :

$$\eta(x_i) = \beta_0 + \beta_1 x_{i1} + \dots + \beta_p x_{ip}$$

et précise quels sont les **prédicteurs** (on peut y inclure des transformations des prédicteurs, des interactions...).

- 3 **Lien** : spécifie le **lien entre les deux composantes**, plus précisément le lien entre l'espérance de Y_i et la composante déterministe : $g(\mathbf{E}[Y_i]) = \eta(x_i) = \beta_0 + \beta_1 x_{i1} + \dots + \beta_p x_{ip}$ où g est une fonction inversible appelée **fonction de lien**.

Modèle linéaire généralisé : GLM

Un modèle linéaire généralisé est constitué de **3 composantes** :

- 1 **Composante aléatoire** : la loi de probabilité de la réponse Y_i appartient à la famille exponentielle et est de la forme

$$f_{\alpha_i}(y_i) = \exp\left(\frac{\alpha_i y_i - b(\alpha_i)}{a(\phi)} + c(y_i, \phi)\right)$$

où a , b et c sont des fonctions spécifiées en fonction du type de la famille exponentielle.

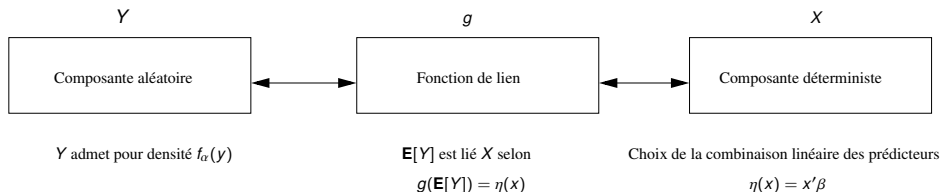
- 2 **Composante déterministe** :

$$\eta(x_i) = \beta_0 + \beta_1 x_{i1} + \dots + \beta_p x_{ip}$$

et précise quels sont les **prédicteurs** (on peut y inclure des transformations des prédicteurs, des interactions...).

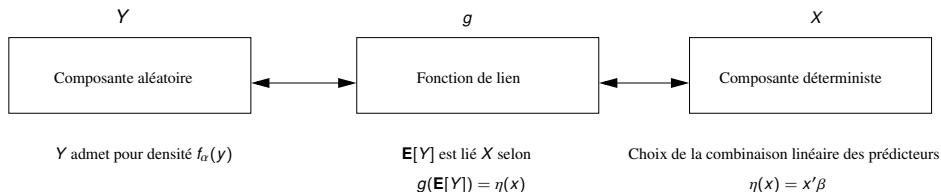
- 3 **Lien** : spécifie le **lien entre les deux composantes**, plus précisément le lien entre l'espérance de Y_i et la composante déterministe : $g(\mathbf{E}[Y_i]) = \eta(x_i) = \beta_0 + \beta_1 x_{i1} + \dots + \beta_p x_{ip}$ où g est une fonction inversible appelée **fonction de lien**.

Schéma GLM



Un modèle GLM sera caractérisé par le choix de ces trois composantes.

Schéma GLM



Un modèle GLM sera caractérisé par le choix de ces trois composantes.

- Le problème du **choix de la combinaison linéaire des variables explicatives** est similaire à tous ce qui a été vu dans le modèle linéaire :
 - Utilisation d'indicatrices pour des **variables explicatives qualitatives** (sans oublier les **contraintes d'identifiabilité**).
 - Possibilité de prendre en compte des **effets quadratique**, ou autre transformation des variables explicatives.
 - Possibilité de prendre en compte des **interactions**.
 - Méthode de **sélection de variables** (stepwise, lasso...)

Dans la suite, on notera $\eta(x) = \beta_0 + \beta_1 x_1 + \dots + \beta_p x_p$ la combinaison linéaire choisie.

- Le problème du **choix de la combinaison linéaire des variables explicatives** est similaire à tous ce qui a été vu dans le modèle linéaire :
 - Utilisation d'indicateurs pour des **variables explicatives qualitatives** (sans oublier les **contraintes d'identifiabilité**).
 - Possibilité de prendre en compte des **effets quadratique**, ou autre transformation des variables explicatives.
 - Possibilité de prendre en compte des **interactions**.
 - Méthode de **sélection de variables** (stepwise, lasso...)

Dans la suite, on notera $\eta(x) = \beta_0 + \beta_1 x_1 + \dots + \beta_p x_p$ la combinaison linéaire choisie.

- Le problème du **choix de la combinaison linéaire des variables explicatives** est similaire à tous ce qui a été vu dans le modèle linéaire :
 - Utilisation d'indicatrices pour des **variables explicatives qualitatives** (sans oublier les **contraintes d'identifiabilité**).
 - Possibilité de prendre en compte des **effets quadratique**, ou autre transformation des variables explicatives.
 - Possibilité de prendre en compte des **interactions**.
 - Méthode de **sélection de variables** (stepwise, lasso...)

Dans la suite, on notera $\eta(x) = \beta_0 + \beta_1 x_1 + \dots + \beta_p x_p$ la combinaison linéaire choisie.

- Le problème du **choix de la combinaison linéaire des variables explicatives** est similaire à tous ce qui a été vu dans le modèle linéaire :
 - Utilisation d'indicatrices pour des **variables explicatives qualitatives** (sans oublier les **contraintes d'identifiabilité**).
 - Possibilité de prendre en compte des **effets quadratique**, ou autre transformation des variables explicatives.
 - Possibilité de prendre en compte des **interactions**.
 - Méthode de **sélection de variables** (stepwise, lasso...)

Dans la suite, on notera $\eta(x) = \beta_0 + \beta_1 x_1 + \dots + \beta_p x_p$ la combinaison linéaire choisie.

- Le problème du **choix de la combinaison linéaire des variables explicatives** est similaire à tous ce qui a été vu dans le modèle linéaire :
 - Utilisation d'indicatrices pour des **variables explicatives qualitatives** (sans oublier les **contraintes d'identifiabilité**).
 - Possibilité de prendre en compte des **effets quadratique**, ou autre transformation des variables explicatives.
 - Possibilité de prendre en compte des **interactions**.
 - Méthode de **sélection de variables** (stepwise, lasso...)

Dans la suite, on notera $\eta(x) = \beta_0 + \beta_1 x_1 + \dots + \beta_p x_p$ la combinaison linéaire choisie.

Composante aléatoire et fonction de lien du modèle logistique

Propriété

Le modèle de régression logistique est un GLM.

En effet :

- La **loi exponentielle** est la loi de Bernoulli de paramètre $p_i = \mathbf{P}(Y_i = 1)$:

$$f_{\alpha_i}(y_i) = \exp[y_i x_i' \beta - \log(1 + \exp(x_i' \beta))].$$

On a donc $\alpha_i = x_i' \beta$ et $b(\alpha_i) = \log(1 + \alpha_i)$.

- La **fonction de lien** est

$$g(u) = \text{logit}(u) = \log \frac{u}{1-u}.$$

Composante aléatoire et fonction de lien du modèle logistique

Propriété

Le modèle de régression logistique est un GLM.

En effet :

- La **loi exponentielle** est la loi de Bernoulli de paramètre $p_i = \mathbf{P}(Y_i = 1)$:

$$f_{\alpha_i}(y_i) = \exp[y_i x_i' \beta - \log(1 + \exp(x_i' \beta))].$$

On a donc $\alpha_i = x_i' \beta$ et $b(\alpha_i) = \log(1 + \alpha_i)$.

- La **fonction de lien** est

$$g(u) = \text{logit}(u) = \log \frac{u}{1 - u}.$$

Composante aléatoire et fonction de lien du modèle linéaire

Propriété

Le modèle linéaire gaussien est un GLM.

En effet :

- La **loi exponentielle** est la loi gaussienne de paramètres μ_i et σ^2 :

$$f_{\alpha_i}(y_i) = \exp \left\{ \frac{y_i x_i' \beta - 0.5(x_i' \beta)^2}{\sigma^2} - \left(\frac{y_i^2}{2\sigma^2} - \frac{\log(2\pi\sigma^2)}{2} \right) \right\}.$$

- La **fonction de lien** est l'identité.

Composante aléatoire et fonction de lien du modèle linéaire

Propriété

Le modèle linéaire gaussien est un GLM.

En effet :

- La **loi exponentielle** est la loi gaussienne de paramètres μ_i et σ^2 :

$$f_{\alpha_i}(y_i) = \exp \left\{ \frac{y_i x_i' \beta - 0.5(x_i' \beta)^2}{\sigma^2} - \left(\frac{y_i^2}{2\sigma^2} - \frac{\log(2\pi\sigma^2)}{2} \right) \right\}.$$

- La **fonction de lien** est l'identité.

- Outre le choix classique de la composante déterministe (choix de la combinaison linéaire des variables explicatives), la modélisation GLM s'effectue à travers 2 choix :
 - ① Choix de la loi de Y_j dans la famille exponentielle GLM décrite plus haut.
 - ② Choix de la fonction de lien (inversible).

	logistique	linéaire
Loi expo	Bernoulli	Gaussienne
fdl	$g(u) = \text{logit}(u)$	$g(u) = u$

- Outre le choix classique de la composante déterministe (choix de la combinaison linéaire des variables explicatives), la modélisation GLM s'effectue à travers 2 choix :
 - ① Choix de la loi de Y_i dans la famille exponentielle GLM décrite plus haut.
 - ② Choix de la fonction de lien (inversible).

	logistique	linéaire
Loi expo	Bernoulli	Gaussienne
fdl	$g(u) = \text{logit}(u)$	$g(u) = u$

- Outre le choix classique de la composante déterministe (choix de la combinaison linéaire des variables explicatives), la modélisation GLM s'effectue à travers 2 choix :
 - ① Choix de la loi de Y_i dans la famille exponentielle GLM décrite plus haut.
 - ② Choix de la fonction de lien (inversible).

	logistique	linéaire
Loi expo	Bernoulli	Gaussienne
fdl	$g(u) = \text{logit}(u)$	$g(u) = u$

- Outre le choix classique de la composante déterministe (choix de la combinaison linéaire des variables explicatives), la modélisation GLM s'effectue à travers 2 choix :
 - ① Choix de la loi de Y_i dans la famille exponentielle GLM décrite plus haut.
 - ② Choix de la fonction de lien (inversible).

	logistique	linéaire
Loi expo	Bernoulli	Gaussienne
fdl	$g(u) = \text{logit}(u)$	$g(u) = u$

- 1 **Loi exponentielle**. Ce choix est généralement guidé par la nature de la variable à expliquer (Binaire : Bernoulli, Comptage : Poisson, continue : normale ou gamma).
- 2 **Fonction de lien**. Ce choix est plus délicat. La fonction de lien dite "canonique" $g(u) = (b')^{-1}(u)$ est souvent privilégiée (notamment pour des raisons d'écriture de modèles et de simplicité d'écriture)

Propriété

Les fonctions de lien des modèles logistique et linéaire sont canoniques.

- 1 **Loi exponentielle.** Ce choix est généralement guidé par la nature de la variable à expliquer (Binaire : Bernoulli, Comptage : Poisson, continue : normale ou gamma).
- 2 **Fonction de lien.** Ce choix est plus délicat. La fonction de lien dite "canonique" $g(u) = (b')^{-1}(u)$ est souvent privilégiée (notamment pour des raisons d'écriture de modèles et de simplicité d'écriture)

Propriété

Les fonctions de lien des modèles logistique et linéaire sont canoniques.

- 1 **Loi exponentielle**. Ce choix est généralement guidé par la nature de la variable à expliquer (Binaire : Bernoulli, Comptage : Poisson, continue : normale ou gamma).
- 2 **Fonction de lien**. Ce choix est plus délicat. La fonction de lien dite "canonique" $g(u) = (b')^{-1}(u)$ est souvent privilégiée (notamment pour des raisons d'écriture de modèles et de simplicité d'écriture)

Propriété

Les fonctions de lien des modèles logistique et linéaire sont canoniques.

Nom du lien	Fonction de lien
identité	$g(u) = u$
log	$g(u) = \log(u)$
cloglog	$g(u) = \log(-\log(1 - u))$
logit	$g(u) = \log(u/(1 - u))$
probit	$g(u) = \Phi^{-1}(u)$
réciproque	$g(u) = -1/u$
puissance	$g(u) = u^\gamma, \gamma \neq 0$

- Il faut bien entendu spécifier à la fonction **glm** les 3 composantes d'un modèle **glm** :

glm(formula=...,family=...(link=...))

- 1 **formula** : spécifie la composante déterministe $Y = X_1 + X_2$,
 $Y = X_1 + X_2 + X_1 : X_2$ (prendre en compte l'interaction entre X_1 et X_2 .)
- 2 **family** : spécifie composante aléatoire (**gaussian** pour le modèle linéaire gaussien, **binomial** lorsque la variable à expliquer est binaire...)
- 3 **link** : spécifie la fonction de lien (**logit** pour logistique, **probit** pour probit...)

- Il faut bien entendu spécifier à la fonction **glm** les 3 composantes d'un modèle **glm** :

glm(formula=...,family=...(link=...))

- 1 **formula** : spécifie la composante déterministe $Y = X_1 + X_2$,
 $Y = X_1 + X_2 + X_1 : X_2$ (prendre en compte l'interaction entre X_1 et X_2 .)
- 2 **family** : spécifie composante aléatoire (**gaussian** pour le modèle linéaire gaussien, **binomial** lorsque la variable à expliquer est binaire...)
- 3 **link** : spécifie la fonction de lien (**logit** pour logistique, **probit** pour probit...)

- Il faut bien entendu spécifier à la fonction **glm** les 3 composantes d'un modèle **glm** :

glm(formula=...,family=...(link=...))

- 1 **formula** : spécifie la composante déterministe $Y = X_1 + X_2$,
 $Y = X_1 + X_2 + X_1 : X_2$ (prendre en compte l'interaction entre X_1 et X_2 .)
- 2 **family** : spécifie composante aléatoire (**gaussian** pour le modèle linéaire gaussien, **binomial** lorsque la variable à expliquer est binaire...)
- 3 **link** : spécifie la fonction de lien (**logit** pour logistique, **probit** pour probit...)

- Il faut bien entendu spécifier à la fonction **glm** les 3 composantes d'un modèle **glm** :

glm(formula=...,family=...(link=...))

- 1 **formula** : spécifie la composante déterministe $Y = X_1 + X_2$,
 $Y = X_1 + X_2 + X_1 : X_2$ (prendre en compte l'interaction entre X_1 et X_2 .)
- 2 **family** : spécifie composante aléatoire (**gaussian** pour le modèle linéaire gaussien, **binomial** lorsque la variable à expliquer est binaire...)
- 3 **link** : spécifie la fonction de lien (**logit** pour logistique, **probit** pour probit...)

Exemple

- On cherche à expliquer une variable binaire Y par deux variables continues X_1 et X_2 :

```
> Y <- rbinom(50,1,0.6)
> X1 <- runif(50)
> X2 <- rnorm(50)
```

- On ajuste les modèles

```
> glm(Y~X1+X2, family=binomial)
```

Coefficients:

(Intercept)	X1	X2
-0.2849	1.8610	-0.0804

```
> glm(Y~X1+X2+X1:X2, family=binomial)
```

Coefficients:

(Intercept)	X1	X2	X1:X2
-0.3395	2.1175	-0.4568	1.0346

```
> glm(Y~X1+X2, family=binomial(link = "probit"))
```

Coefficients:

(Intercept)	X1	X2
-0.17038	1.11986	-0.04864

Modèle de Poisson

Le problème

- On cherche à quantifier l'influence d'un traitement sur l'évolution du nombre de polypes au colon. On dispose des données suivantes :

```
number  treat  age
1      63 placebo 20
2       2   drug 16
3      28 placebo 18
4      17   drug 22
5      61 placebo 13
...
```

où

- `number` : nombre de polypes après 12 mois de traitement.
- `treat` : `drug` si le traitement a été administré, `placebo` sinon.
- `age` : age de l'individu.

Le problème est d'expliquer la variable `number` par les deux autres variables à l'aide d'un GLM.

Le problème

- On cherche à quantifier l'influence d'un traitement sur l'évolution du nombre de polypes au colon. On dispose des données suivantes :

```
number  treat  age
1      63 placebo 20
2       2   drug 16
3      28 placebo 18
4      17   drug 22
5      61 placebo 13
...
```

où

- `number` : nombre de polypes après 12 mois de traitement.
- `treat` : `drug` si le traitement a été administré, `placebo` sinon.
- `age` : age de l'individu.

Le problème est d'expliquer la variable `number` par les deux autres variables à l'aide d'un GLM.

On note

- Y_i la variable aléatoire représentant le nombre de polypes du i ème patient après les 12 mois de traitement.
- x_{i1} la variable `treat` pour le i ème individu et x_{i2} l'âge du i ème individu.

GLM

- 1 La variable Y_i étant une variable de **comptage**, on choisit comme densité de Y_i la densité (par rapport à la mesure de comptage) de la loi de **Poisson** de paramètre λ_i :

$$f_{\alpha_i}(y_i) = \exp(-\lambda_i) \frac{\lambda_i^{y_i}}{y_i!} = \exp[y_i \log(\lambda_i) - \exp(\log(\lambda_i)) - \log(y_i!)].$$

- 2 La **fonction de lien canonique** est donc donnée par :

$$g(u) = \log(u).$$

On note

- Y_i la variable aléatoire représentant le nombre de polypes du i ème patient après les 12 mois de traitement.
- x_{i1} la variable `treat` pour le i ème individu et x_{i2} l'âge du i ème individu.

GLM

- 1 La variable Y_i étant une variable de **comptage**, on choisit comme densité de Y_i la densité (par rapport à la mesure de comptage) de la loi de **Poisson** de paramètre λ_i :

$$f_{\alpha_i}(y_i) = \exp(-\lambda_i) \frac{\lambda_i^{y_i}}{y_i!} = \exp[y_i \log(\lambda_i) - \exp(\log(\lambda_i)) - \log(y_i!)].$$

- 2 La **fonction de lien canonique** est donc donnée par :

$$g(u) = \log(u).$$

Définition

Le **modèle de Poisson** modélise la loi de Y_i par une loi de Poisson de paramètre $\lambda_i = \lambda(x_i)$ telle que

$$\log(\lambda(x_i)) = x_i' \beta.$$

- L'ajustement sur R s'effectue toujours à l'aide de la fonction **glm** :

```
> glm(number~treat+age, data=polyps, family=poisson)
```

Coefficients:

(Intercept)	treatdrug	age
4.52902	-1.35908	-0.03883

- **Prédiction** : pour un individu de 23 ans, ayant reçu le traitement on pourra estimer le nombre de polypes à 12 mois par

$$\exp(4.52902 - 1.35908 - 0.03883 * 23) = 9.745932.$$

Définition

Le **modèle de Poisson** modélise la loi de Y_i par une loi de Poisson de paramètre $\lambda_i = \lambda(x_i)$ telle que

$$\log(\lambda(x_i)) = x_i' \beta.$$

- L'ajustement sur R s'effectue toujours à l'aide de la fonction **glm** :

```
> glm(number~treat+age, data=polyps, family=poisson)
```

Coefficients:

(Intercept)	treatdrug	age
4.52902	-1.35908	-0.03883

- **Prédiction** : pour un individu de 23 ans, ayant reçu le traitement on pourra estimer le nombre de polypes à 12 mois par

$$\exp(4.52902 - 1.35908 - 0.03883 * 23) = 9.745932.$$

Définition

Le **modèle de Poisson** modélise la loi de Y_i par une loi de Poisson de paramètre $\lambda_i = \lambda(x_i)$ telle que

$$\log(\lambda(x_i)) = x_i' \beta.$$

- L'ajustement sur R s'effectue toujours à l'aide de la fonction **glm** :

```
> glm(number~treat+age, data=polyps, family=poisson)
```

Coefficients:

(Intercept)	treatdrug	age
4.52902	-1.35908	-0.03883

- **Prédiction** : pour un individu de 23 ans, ayant reçu le traitement on pourra estimer le nombre de polypes à 12 mois par

$$\exp(4.52902 - 1.35908 - 0.03883 * 23) = 9.745932.$$

Définition

Le **modèle de Poisson** modélise la loi de Y_i par une loi de Poisson de paramètre $\lambda_i = \lambda(x_i)$ telle que

$$\log(\lambda(x_i)) = x_i' \beta.$$

- L'ajustement sur R s'effectue toujours à l'aide de la fonction **glm** :

```
> glm(number~treat+age, data=polyps, family=poisson)
```

Coefficients:

(Intercept)	treatdrug	age
4.52902	-1.35908	-0.03883

- **Prédiction** : pour un individu de 23 ans, ayant reçu le traitement on pourra estimer le nombre de polypes à 12 mois par

$$\exp(4.52902 - 1.35908 - 0.03883 * 23) = 9.745932.$$

