RÉGRESSION LOGISTIQUE (LR) Modèle additif généralisé (GAM) ANALYSE DISCRIMINANTE LINÉAIRE (LDA)

Jérémie Cabessa Laboratoire DAVID. UVSQ

- ▶ Dans le cadre de l'apprentissage supervisé, on distingue deux types de méthodes:
 - Méthodes de régression
 La variable d'output (réponse) est quantitative.
- Méthodes de classification
 La variable d'output (réponse) est qualitative.

- ► Dans le cadre de l'apprentissage supervisé, on distingue deux types de méthodes:
- Méthodes de régression
 La variable d'output (réponse) est quantitative.
- Méthodes de classification
 La variable d'output (réponse) est qualitative

- ▶ Dans le cadre de l'apprentissage supervisé, on distingue deux types de méthodes:
- Méthodes de régression
 La variable d'output (réponse) est quantitative.
- Méthodes de classification
 La variable d'output (réponse) est qualitative.

CLASSIFICATION

Ce chapitre rappelle/présente les méthodes de classification suivantes:

- K-Nearest Neighbors (KNN)
- Logistic Regression (LR, régression logistique)
- Generalized Additive Models (GAM)

CLASSIFICATION

Ce chapitre rappelle/présente les méthodes de classification suivantes:

- K-Nearest Neighbors (KNN)
- Logistic Regression (LR, régression logistique)
- Generalized Additive Models (GAM)

CLASSIFICATION

Ce chapitre rappelle/présente les méthodes de classification suivantes:

- K-Nearest Neighbors (KNN)
- Logistic Regression (LR, régression logistique)
- Generalized Additive Models (GAM)

- ▶ Soient $X = (X_1, ..., X_p)$ des variables explicatives et Y une variable réponse qualitative *binaire*.
- ightharpoonup On code les réalisations possibles de Y par 0 ou 1.
- On aimerait modéliser la probabilité que Y=1 étant données les réalisations des variables X_1, \ldots, X_p , i.e.,:

$$\Pr(Y = 1 \mid X) = \Pr(Y = 1 \mid X_1, \dots, X_p)$$

- ▶ Soient $X = (X_1, ..., X_p)$ des variables explicatives et Y une variable réponse qualitative *binaire*.
- ightharpoonup On code les réalisations possibles de Y par 0 ou 1.
- On aimerait modéliser la probabilité que Y=1 étant données les réalisations des variables X_1, \ldots, X_p , i.e.,:

$$\Pr(Y = 1 \mid X) = \Pr(Y = 1 \mid X_1, \dots, X_p)$$

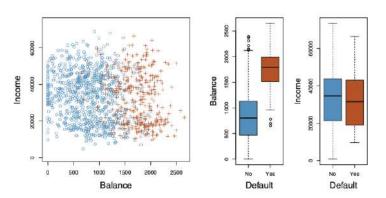
- ▶ Soient $X = (X_1, ..., X_p)$ des variables explicatives et Y une variable réponse qualitative *binaire*.
- On code les réalisations possibles de Y par 0 ou 1.
- On aimerait modéliser la probabilité que Y=1 étant données les réalisations des variables X_1,\ldots,X_p , i.e.,:

$$\Pr(Y = 1 \mid X) = \Pr(Y = 1 \mid X_1, \dots, X_p)$$

- ▶ Soient $X = (X_1, ..., X_p)$ des variables explicatives et Y une variable réponse qualitative *binaire*.
- ➤ On code les réalisations possibles de Y par 0 ou 1.
- On aimerait modéliser la probabilité que Y=1 étant données les réalisations des variables X_1,\ldots,X_p , i.e.,:

$$\Pr(Y = 1 \mid X) = \Pr(Y = 1 \mid X_1, \dots, X_p)$$

Pour diverses raisons, les régressions de types linéaires ne sont pas appropriées...



Pour diverses raisons, les régressions de types linéaires ne sont pas appropriées...

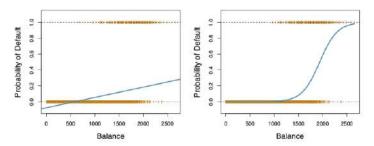


FIGURE 4.2. Classification using the Default data. Left: Estimated probability of default using linear regression. Some estimated probabilities are negative! The orange ticks indicate the 0/1 values coded for default (No or Yes). Right: Predicted probabilities of default using logistic regression. All probabilities lie between 0 and 1.

- ▶ On aimerait que $\Pr(Y=1 \mid \boldsymbol{X}) \in [0,1]$, puisque c'est une probabilité.
- ▶ On suppose alors que (la vraie probabilité) $\Pr(Y = 1 \mid X)$ est donnée par la fonction logistique suivante:

$$\Pr(Y = 1 \mid \mathbf{X}) = \frac{e^{\beta_0 + \beta_1 X_1 + \dots + \beta_p X_p}}{1 + e^{\beta_0 + \beta_1 X_1 + \dots + \beta_p X_p}}$$

Remarques

$$eta_0 + eta_1 X_1 + \dots + eta_p X_p \to +\infty$$
 implique $\Pr(Y = 1 \mid \boldsymbol{X}) \to 1$
 $eta_0 + eta_1 X_1 + \dots + eta_p X_p \to -\infty$ implique $\Pr(Y = 1 \mid \boldsymbol{X}) \to 0$

Si on connait $\Pr(Y=1\mid \boldsymbol{X})$ alors on peut immédiatement déduire $\Pr(Y=0\mid \boldsymbol{X})=1-\Pr(Y=1\mid \boldsymbol{X})$

- ▶ On aimerait que $\Pr(Y = 1 \mid \mathbf{X}) \in [0, 1]$, puisque c'est une probabilité.
- ▶ On suppose alors que (la vraie probabilité) $\Pr(Y = 1 \mid X)$ est donnée par la **fonction logistique** suivante:

$$\Pr(Y = 1 \mid \mathbf{X}) = \frac{e^{\beta_0 + \beta_1 X_1 + \dots + \beta_p X_p}}{1 + e^{\beta_0 + \beta_1 X_1 + \dots + \beta_p X_p}}$$

Remarques

$$eta_0 + eta_1 X_1 + \dots + eta_p X_p o + \infty$$
 implique $\Pr(Y = 1 \mid \boldsymbol{X}) o 1$
 $eta_0 + eta_1 X_1 + \dots + eta_p X_p o - \infty$ implique $\Pr(Y = 1 \mid \boldsymbol{X}) o 0$

Si on connait $\Pr(Y=1\mid \boldsymbol{X})$ alors on peut immédiatement déduire $\Pr(Y=0\mid \boldsymbol{X})=1-\Pr(Y=1\mid \boldsymbol{X})$

- ▶ On aimerait que $\Pr(Y = 1 \mid \mathbf{X}) \in [0, 1]$, puisque c'est une probabilité.
- ▶ On suppose alors que (la vraie probabilité) $\Pr(Y = 1 \mid X)$ est donnée par la fonction logistique suivante:

$$\Pr(Y = 1 \mid \mathbf{X}) = \frac{e^{\beta_0 + \beta_1 X_1 + \dots + \beta_p X_p}}{1 + e^{\beta_0 + \beta_1 X_1 + \dots + \beta_p X_p}}$$

Remarques:

$$\beta_0 + \beta_1 X_1 + \dots + \beta_p X_p \to +\infty \text{ implique } \Pr(Y = 1 \mid \boldsymbol{X}) \to 1$$

$$\beta_0 + \beta_1 X_1 + \dots + \beta_p X_p \to -\infty \text{ implique } \Pr(Y = 1 \mid \boldsymbol{X}) \to 0$$

Si on connait $\Pr(Y = 1 \mid \boldsymbol{X})$ alors on peut immédiatement déduire $\Pr(Y = 0 \mid \boldsymbol{X}) = 1 - \Pr(Y = 1 \mid \boldsymbol{X})$

- ▶ On aimerait que $\Pr(Y=1\mid \pmb{X})\in [0,1]$, puisque c'est une probabilité.
- ▶ On suppose alors que (la vraie probabilité) $\Pr(Y = 1 \mid \mathbf{X})$ est donnée par la fonction logistique suivante:

$$\Pr(Y = 1 \mid X) = \frac{e^{\beta_0 + \beta_1 X_1 + \dots + \beta_p X_p}}{1 + e^{\beta_0 + \beta_1 X_1 + \dots + \beta_p X_p}}$$

Remarques:

$$\beta_0 + \beta_1 X_1 + \dots + \beta_p X_p \to +\infty \text{ implique } \Pr(Y = 1 \mid \boldsymbol{X}) \to 1$$

$$\beta_0 + \beta_1 X_1 + \dots + \beta_p X_p \to -\infty \text{ implique } \Pr(Y = 1 \mid \boldsymbol{X}) \to 0$$

Si on connait $\Pr(Y = 1 \mid \boldsymbol{X})$ alors on peut immédiatement déduire $\Pr(Y = 0 \mid \boldsymbol{X}) = 1 - \Pr(Y = 1 \mid \boldsymbol{X})$

▶ Note hypothèse sur la forme de $Pr(Y = 1 \mid \boldsymbol{X})$

$$\Pr(Y = 1 \mid \mathbf{X}) = \frac{e^{\beta_0 + \beta_1 X_1 + \dots + \beta_p X_p}}{1 + e^{\beta_0 + \beta_1 X_1 + \dots + \beta_p X_p}}$$

implique que la fonction logit est de la forme linéaire suivante:

$$\log \left(\frac{\Pr(Y=1 \mid \boldsymbol{X})}{1 - \Pr(Y=1 \mid \boldsymbol{X})} \right) = \beta_0 + \beta_1 X_1 + \dots + \beta_p X_p$$

- Étant donné un training set $S_{\text{train}} = \{(\boldsymbol{x_1}, y_1), \dots, (\boldsymbol{x_N}, y_N)\}$, on aimerait estimer les paramètres $\boldsymbol{\beta} = (\beta_0, \beta_1, \dots, \beta_p)$ de la fonction logit de manière pertinente...
- On aimerait obtenir des estimateurs $\hat{\beta} = (\hat{\beta}_0, \hat{\beta}_1, \dots, \hat{\beta}_p)$ tels que, pour tout $(x_i, y_i) \in S_{\text{train}}$, on a:

$$y_i = 1 \implies \hat{\Pr}(Y = 1 \mid X = x_i) > 0.5$$

 $y_i = 0 \implies \hat{\Pr}(Y = 0 \mid X = x_i) = 1 - \hat{\Pr}(Y = 1 \mid X = x_i) \le 0.5$

où
$$\hat{\Pr}(Y=1\mid \boldsymbol{X}=\boldsymbol{x_i}):=\frac{e^{\hat{\beta}_0+\sum_{k=1}^p\hat{\beta}_kx_{ik}}}{1+e^{\hat{\beta}_0+\sum_{k=1}^p\hat{\beta}_kx_{ik}}}$$
 est l'estimateur de $\Pr(Y=1\mid \boldsymbol{X}=\boldsymbol{x_i})$ donné par les $\hat{\beta}_0,\hat{\beta}_1,\ldots,\hat{\beta}_p$.

- Étant donné un training set $S_{\text{train}} = \{(\boldsymbol{x_1}, y_1), \dots, (\boldsymbol{x_N}, y_N)\}$, on aimerait estimer les paramètres $\boldsymbol{\beta} = (\beta_0, \beta_1, \dots, \beta_p)$ de la fonction logit de manière pertinente...
- On aimerait obtenir des estimateurs $\hat{\beta}=(\hat{\beta}_0,\hat{\beta}_1,\ldots,\hat{\beta}_p)$ tels que, pour tout $(\boldsymbol{x_i},y_i)\in S_{\mathrm{train}}$, on a:

$$y_i = 1 \implies \hat{\Pr}(Y = 1 \mid \mathbf{X} = \mathbf{x_i}) > 0.5$$

 $y_i = 0 \implies \hat{\Pr}(Y = 0 \mid \mathbf{X} = \mathbf{x_i}) = 1 - \hat{\Pr}(Y = 1 \mid \mathbf{X} = \mathbf{x_i}) \le 0.5.$

où
$$\hat{\Pr}(Y=1\mid \boldsymbol{X}=\boldsymbol{x_i}):=\frac{e^{\hat{\beta}_0+\sum_{k=1}^p\hat{\beta}_kx_{ik}}}{1+e^{\hat{\beta}_0+\sum_{k=1}^p\hat{\beta}_kx_{ik}}}$$
 est l'estimateur de $\Pr(Y=1\mid \boldsymbol{X}=\boldsymbol{x_i})$ donné par les $\hat{\beta}_0,\hat{\beta}_1,\ldots,\hat{\beta}_p$.

$$\hat{\boldsymbol{\beta}} = \arg\max_{\boldsymbol{\beta}} L(\boldsymbol{\beta}) \quad \text{où}$$

$$L(\boldsymbol{\beta}) = L(\beta_0, \beta_1, \dots, \beta_p) := \prod_{\{i: y_i = 1\}} \Pr(Y = 1 \mid \boldsymbol{X} = \boldsymbol{x}_i) \prod_{\{j: y_j = 0\}} \Pr(Y = 0 \mid \boldsymbol{X} = \boldsymbol{x}_j) = \prod_{\{i: y_i = 1\}} \Pr(Y = 1 \mid \boldsymbol{X} = \boldsymbol{x}_i) \prod_{\{j: y_j = 0\}} \left(1 - \Pr(Y = 1 \mid \boldsymbol{X} = \boldsymbol{x}_j)\right) = \prod_{\{i: y_i = 1\}} \frac{e^{\beta_0 + \sum_{k=1}^p \beta_k x_{ik}}}{1 + e^{\beta_0 + \sum_{k=1}^p \beta_k x_{ik}}} \prod_{\{j: y_j = 0\}} \left(1 - \frac{e^{\beta_0 + \sum_{k=1}^p \beta_k x_{jk}}}{1 + e^{\beta_0 + \sum_{k=1}^p \beta_k x_{jk}}}\right) = \prod_{\{i: y_i = 1\}} \frac{e^{\beta^T x_i}}{1 + e^{\beta^T x_i}} \prod_{\{i: y_i = 0\}} \left(\frac{1}{1 + e^{\beta^T x_j}}\right) \quad \text{où } x_i = (1, x_i) \text{ (abus de notation)}$$

$$\begin{split} \widehat{\boldsymbol{\beta}} &= \arg\max_{\boldsymbol{\beta}} L(\boldsymbol{\beta}) \quad \text{où} \\ L(\boldsymbol{\beta}) &= L(\beta_0, \beta_1, \dots, \beta_p) := \\ \prod_{\{i: y_i = 1\}} \Pr(Y = 1 \mid \boldsymbol{X} = \boldsymbol{x_i}) \prod_{\{j: y_j = 0\}} \Pr(Y = 0 \mid \boldsymbol{X} = \boldsymbol{x_j}) = \\ \prod_{\{i: y_i = 1\}} \Pr(Y = 1 \mid \boldsymbol{X} = \boldsymbol{x_i}) \prod_{\{j: y_j = 0\}} \left(1 - \Pr(Y = 1 \mid \boldsymbol{X} = \boldsymbol{x_j})\right) = \\ \prod_{\{i: y_i = 1\}} \frac{e^{\beta_0 + \sum_{k=1}^p \beta_k x_{ik}}}{1 + e^{\beta_0 + \sum_{k=1}^p \beta_k x_{ik}}} \prod_{\{j: y_j = 0\}} \left(1 - \frac{e^{\beta_0 + \sum_{k=1}^p \beta_k x_{jk}}}{1 + e^{\beta_0 + \sum_{k=1}^p \beta_k x_{ik}}}\right) = \\ \prod_{\{i: y_i = 1\}} \frac{e^{\beta^T x_i}}{1 + e^{\beta^T x_i}} \prod_{\{j: y_j = 0\}} \left(\frac{1}{1 + e^{\beta^T x_j}}\right) \quad \text{où } \boldsymbol{x_i} = (1, \boldsymbol{x_i}) \\ \text{(abus de notation)} \end{split}$$

$$\hat{\boldsymbol{\beta}} = \arg\max_{\boldsymbol{\beta}} L(\boldsymbol{\beta})$$
 où

$$\begin{split} L(\pmb{\beta}) &= L(\beta_0, \beta_1, \dots, \beta_p) := \\ &\prod_{\{i: y_i = 1\}} \Pr(Y = 1 \mid \pmb{X} = \pmb{x}_i) \prod_{\{j: y_j = 0\}} \Pr(Y = 0 \mid \pmb{X} = \pmb{x}_j) = \\ &\prod_{\{i: y_i = 1\}} \Pr(Y = 1 \mid \pmb{X} = \pmb{x}_i) \prod_{\{j: y_j = 0\}} \left(1 - \Pr(Y = 1 \mid \pmb{X} = \pmb{x}_j)\right) = \\ &\prod_{\{i: y_i = 1\}} \frac{e^{\beta_0 + \sum_{k=1}^p \beta_k x_{ik}}}{1 + e^{\beta_0 + \sum_{k=1}^p \beta_k x_{ik}}} \prod_{\{j: y_j = 0\}} \left(1 - \frac{e^{\beta_0 + \sum_{k=1}^p \beta_k x_{jk}}}{1 + e^{\beta_0 + \sum_{k=1}^p \beta_k x_{jk}}}\right) = \\ &\prod_{\{i: y_i = 1\}} \frac{e^{\beta^T x_i}}{1 + e^{\beta^T x_i}} \prod_{\{j: y_j = 0\}} \left(\frac{1}{1 + e^{\beta^T x_j}}\right) \quad \text{où } x_i = (1, x_i) \\ \text{(abus de notation)} \end{split}$$

$$\hat{oldsymbol{eta}} = rg \max_{oldsymbol{eta}} L(oldsymbol{eta})$$
 où

$$L(\boldsymbol{\beta}) = L(\beta_0, \beta_1, \dots, \beta_p) := \prod_{\{i: y_i = 1\}} \Pr(Y = 1 \mid \boldsymbol{X} = \boldsymbol{x_i}) \prod_{\{j: y_j = 0\}} \Pr(Y = 0 \mid \boldsymbol{X} = \boldsymbol{x_j}) = \prod_{\{i: y_i = 1\}} \Pr(Y = 1 \mid \boldsymbol{X} = \boldsymbol{x_i}) \prod_{\{j: y_j = 0\}} (1 - \Pr(Y = 1 \mid \boldsymbol{X} = \boldsymbol{x_j})) = \prod_{\{i: y_i = 1\}} \frac{e^{\beta_0 + \sum_{k=1}^p \beta_k x_{ik}}}{1 + e^{\beta_0 + \sum_{k=1}^p \beta_k x_{ik}}} \prod_{\{j: y_j = 0\}} \left(1 - \frac{e^{\beta_0 + \sum_{k=1}^p \beta_k x_{jk}}}{1 + e^{\beta_0 + \sum_{k=1}^p \beta_k x_{jk}}}\right) = \prod_{\{i: y_i = 1\}} \frac{e^{\beta^T x_i}}{1 + e^{\beta^T x_i}} \prod_{\{j: y_j = 0\}} \left(\frac{1}{1 + e^{\beta^T x_j}}\right) \quad \text{où } x_i = (1, x_i) \text{ (abus de notation)}$$

$$\hat{oldsymbol{eta}} = rg\max_{oldsymbol{eta}} L(oldsymbol{eta})$$
 où

$$L(\boldsymbol{\beta}) = L(\beta_0, \beta_1, \dots, \beta_p) := \\ \prod_{\{i: y_i = 1\}} \Pr(Y = 1 \mid \boldsymbol{X} = \boldsymbol{x_i}) \prod_{\{j: y_j = 0\}} \Pr(Y = 0 \mid \boldsymbol{X} = \boldsymbol{x_j}) = \\ \prod_{\{i: y_i = 1\}} \Pr(Y = 1 \mid \boldsymbol{X} = \boldsymbol{x_i}) \prod_{\{j: y_j = 0\}} \left(1 - \Pr(Y = 1 \mid \boldsymbol{X} = \boldsymbol{x_j})\right) = \\ \prod_{\{i: y_i = 1\}} \frac{e^{\beta_0 + \sum_{k=1}^p \beta_k x_{ik}}}{1 + e^{\beta_0 + \sum_{k=1}^p \beta_k x_{ik}}} \prod_{\{j: y_j = 0\}} \left(1 - \frac{e^{\beta_0 + \sum_{k=1}^p \beta_k x_{jk}}}{1 + e^{\beta_0 + \sum_{k=1}^p \beta_k x_{jk}}}\right) = \\ \prod_{\{i: y_i = 1\}} \frac{e^{\beta^T x_i}}{1 + e^{\beta^T x_i}} \prod_{\{j: y_i = 0\}} \left(\frac{1}{1 + e^{\beta^T x_j}}\right) \quad \text{où } x_i = (1, x_i) \\ \text{(abus de notation)}$$

$$\hat{oldsymbol{eta}} = rg \max_{oldsymbol{eta}} L(oldsymbol{eta})$$
 où

$$\begin{split} L(\boldsymbol{\beta}) &= L(\beta_0, \beta_1, \dots, \beta_p) := \\ &\prod_{\{i: y_i = 1\}} \Pr(Y = 1 \mid \boldsymbol{X} = \boldsymbol{x_i}) \prod_{\{j: y_j = 0\}} \Pr(Y = 0 \mid \boldsymbol{X} = \boldsymbol{x_j}) = \\ &\prod_{\{i: y_i = 1\}} \Pr(Y = 1 \mid \boldsymbol{X} = \boldsymbol{x_i}) \prod_{\{j: y_j = 0\}} \left(1 - \Pr(Y = 1 \mid \boldsymbol{X} = \boldsymbol{x_j})\right) = \\ &\prod_{\{i: y_i = 1\}} \frac{e^{\beta_0 + \sum_{k=1}^p \beta_k x_{ik}}}{1 + e^{\beta_0 + \sum_{k=1}^p \beta_k x_{ik}}} \prod_{\{j: y_j = 0\}} \left(1 - \frac{e^{\beta_0 + \sum_{k=1}^p \beta_k x_{jk}}}{1 + e^{\beta_0 + \sum_{k=1}^p \beta_k x_{jk}}}\right) = \\ &\prod_{\{i: y_i = 1\}} \frac{e^{\beta^T x_i}}{1 + e^{\beta^T x_i}} \prod_{\{j: y_j = 0\}} \left(\frac{1}{1 + e^{\beta^T x_j}}\right) \quad \text{où } x_i = (1, x_i) \\ &\text{(abus de notation)} \end{split}$$

$$\hat{oldsymbol{eta}} = rg \max_{oldsymbol{eta}} L(oldsymbol{eta})$$
 où

$$\begin{split} L(\boldsymbol{\beta}) &= L(\beta_0, \beta_1, \dots, \beta_p) := \\ &\prod_{\{i: y_i = 1\}} \Pr(Y = 1 \mid \boldsymbol{X} = \boldsymbol{x_i}) \prod_{\{j: y_j = 0\}} \Pr(Y = 0 \mid \boldsymbol{X} = \boldsymbol{x_j}) = \\ &\prod_{\{i: y_i = 1\}} \Pr(Y = 1 \mid \boldsymbol{X} = \boldsymbol{x_i}) \prod_{\{j: y_j = 0\}} \left(1 - \Pr(Y = 1 \mid \boldsymbol{X} = \boldsymbol{x_j})\right) = \\ &\prod_{\{i: y_i = 1\}} \frac{e^{\beta_0 + \sum_{k=1}^p \beta_k x_{ik}}}{1 + e^{\beta_0 + \sum_{k=1}^p \beta_k x_{ik}}} \prod_{\{j: y_j = 0\}} \left(1 - \frac{e^{\beta_0 + \sum_{k=1}^p \beta_k x_{jk}}}{1 + e^{\beta_0 + \sum_{k=1}^p \beta_k x_{jk}}}\right) = \\ &\prod_{\{i: y_i = 1\}} \frac{e^{\boldsymbol{\beta}^T \boldsymbol{x_i}}}{1 + e^{\boldsymbol{\beta}^T \boldsymbol{x_i}}} \prod_{\{j: y_i = 0\}} \left(\frac{1}{1 + e^{\boldsymbol{\beta}^T \boldsymbol{x_j}}}\right) \quad \text{où } \boldsymbol{x_i} = (1, \boldsymbol{x_i}) \\ \text{(abus de notation)} \end{split}$$

ightharpoonup Dans la formule de $L(\beta)$, plus les probabilités

$$Pr(Y = 1 \mid \boldsymbol{X} = \boldsymbol{x_i})$$
 et $Pr(Y = 0 \mid \boldsymbol{X} = \boldsymbol{x_j})$

sont proches de 1 (i.e. sont corrects), plus la valeur de $L(\beta)$ est grande $(\to 1)$.

- ▶ D'où l'idée de chercher les estimateurs $\hat{\beta} = (\hat{\beta}_0, \hat{\beta}_1, \dots, \hat{\beta}_p)$ qui maximisent $L(\beta)$.
- ▶ En pratique, les paramètres $\hat{\beta} = (\hat{\beta}_0, \hat{\beta}_1, \dots, \hat{\beta}_p)$ qui constituent la solution de ce problème de maximisation sont calculés par une méthode de gradient itérative...

▶ Dans la formule de $L(\beta)$, plus les probabilités

$$Pr(Y = 1 \mid \boldsymbol{X} = \boldsymbol{x_i})$$
 et $Pr(Y = 0 \mid \boldsymbol{X} = \boldsymbol{x_i})$

sont proches de 1 (i.e. sont corrects), plus la valeur de $L(\beta)$ est grande $(\to 1)$.

- ▶ D'où l'idée de chercher les estimateurs $\hat{\beta} = (\hat{\beta}_0, \hat{\beta}_1, \dots, \hat{\beta}_p)$ qui maximisent $L(\beta)$.
- ▶ En pratique, les paramètres $\hat{\beta} = (\hat{\beta}_0, \hat{\beta}_1, \dots, \hat{\beta}_p)$ qui constituent la solution de ce problème de maximisation sont calculés par une méthode de gradient itérative...

▶ Dans la formule de $L(\beta)$, plus les probabilités

$$Pr(Y = 1 \mid \boldsymbol{X} = \boldsymbol{x_i})$$
 et $Pr(Y = 0 \mid \boldsymbol{X} = \boldsymbol{x_i})$

sont proches de 1 (i.e. sont corrects), plus la valeur de $L(\beta)$ est grande $(\to 1)$.

- ▶ D'où l'idée de chercher les estimateurs $\hat{\beta} = (\hat{\beta}_0, \hat{\beta}_1, \dots, \hat{\beta}_p)$ qui maximisent $L(\beta)$.
- ▶ En pratique, les paramètres $\hat{\beta} = (\hat{\beta}_0, \hat{\beta}_1, \dots, \hat{\beta}_p)$ qui constituent la solution de ce problème de maximisation sont calculés par une méthode de gradient itérative...

- Une fois les paramètres $\hat{\beta} = (\hat{\beta}_0, \hat{\beta}_1, \dots, \hat{\beta}_p)$ calculés, on peut faire des prédictions de manière très simple.
- ▶ Soit $x = (x_1, \dots, x_p)$ un point. La prédiction \hat{y} associée à x est donnée par:

$$\hat{y} := \begin{cases} 1, \text{ si } \hat{\Pr}(Y = 1 \mid \boldsymbol{X} = \boldsymbol{x}) = \frac{e^{\hat{\boldsymbol{\beta}}^T \boldsymbol{x}}}{1 + e^{\hat{\boldsymbol{\beta}}^T \boldsymbol{x}}} > 0.5 \\ 0, \text{ si } \hat{\Pr}(Y = 1 \mid \boldsymbol{X} = \boldsymbol{x}) = \frac{e^{\hat{\boldsymbol{\beta}}^T \boldsymbol{x}}}{1 + e^{\hat{\boldsymbol{\beta}}^T \boldsymbol{x}}} \le 0.5 \end{cases}$$

- Une fois les paramètres $\hat{\beta} = (\hat{\beta}_0, \hat{\beta}_1, \dots, \hat{\beta}_p)$ calculés, on peut faire des prédictions de manière très simple.
- ▶ Soit $x = (x_1, \dots, x_p)$ un point. La prédiction \hat{y} associée à x est donnée par:

$$\hat{y} := \begin{cases} 1, \text{ si } \hat{\Pr}(Y = 1 \mid \boldsymbol{X} = \boldsymbol{x}) = \frac{e^{\hat{\boldsymbol{\beta}}^T \boldsymbol{x}}}{1 + e^{\hat{\boldsymbol{\beta}}^T \boldsymbol{x}}} > 0.5\\ 0, \text{ si } \hat{\Pr}(Y = 1 \mid \boldsymbol{X} = \boldsymbol{x}) = \frac{e^{\hat{\boldsymbol{\beta}}^T \boldsymbol{x}}}{1 + e^{\hat{\boldsymbol{\beta}}^T \boldsymbol{x}}} \le 0.5 \end{cases}$$

- Le processus peut se généraliser à un contexte multi-classes quelconque. Supposons que Y possède $k \geq 2$ valeurs possibles (k classes): c_1, \ldots, c_k .
- On code Y par une variable $Y = (Y_1, \ldots, Y_k)$ telle que $Y_i = 1$ ssi $Y = c_i$ (1-hot encoding).
- On estime les paramètres $\hat{\beta} = (\beta_0, \beta_1, \dots, \beta_p)$ de manière similaire, ce qui permet l'estimation $\Pr(\bar{Y} = \bar{y} \mid X = x)$.
- La prédiction \hat{y} associée à x est alors donnée par (on note $\mathbf{1}_i$ le vecteur de dim k avec un 1 en position i et des 0 partout ailleurs): $\hat{y} = c_i$ is et seulement si

$$\Pr(ar{Y}=oldsymbol{1}_i\mid X=oldsymbol{x})>\Pr(ar{Y}=oldsymbol{1}_j\mid X=oldsymbol{x})$$
 pour tout $j
eq i$

- Le processus peut se généraliser à un contexte multi-classes quelconque. Supposons que Y possède $k \geq 2$ valeurs possibles (k classes): c_1, \ldots, c_k .
- ▶ On code Y par une variable $\bar{Y} = (Y_1, \dots, Y_k)$ telle que $Y_i = 1$ ssi $Y = c_i$ (1-hot encoding).
- On estime les paramètres $\hat{\beta} = (\beta_0, \beta_1, \dots, \beta_p)$ de manière similaire, ce qui permet l'estimation $\Pr(\bar{Y} = \bar{y} \mid X = x)$.
- La prédiction \hat{y} associée à x est alors donnée par (on note 1_i le vecteur de dim k avec un 1 en position i et des 0 partout ailleurs): $\hat{y} = c_i$ si et seulement si

$$\Pr(\bar{Y} = \mathbf{1}_i \mid X = x) > \Pr(\bar{Y} = \mathbf{1}_j \mid X = x)$$
 pour tout $j \neq i$

- Le processus peut se généraliser à un contexte multi-classes quelconque. Supposons que Y possède $k \geq 2$ valeurs possibles (k classes): c_1, \ldots, c_k .
- ▶ On code Y par une variable $\bar{Y} = (Y_1, \dots, Y_k)$ telle que $Y_i = 1$ ssi $Y = c_i$ (1-hot encoding).
- On estime les paramètres $\hat{\beta} = (\hat{\beta}_0, \hat{\beta}_1, \dots, \hat{\beta}_p)$ de manière similaire, ce qui permet l'estimation $\Pr(\bar{Y} = \bar{y} \mid X = x)$.
- La prédiction \hat{y} associée à x est alors donnée par (on note $\mathbf{1}_i$ le vecteur de dim k avec un 1 en position i et des 0 partout ailleurs): $\hat{y} = c_i$ si et seulement si

$$\Pr(\bar{Y} = \mathbf{1}_i \mid X = x) > \Pr(\bar{Y} = \mathbf{1}_j \mid X = x)$$
 pour tout $j \neq i$

- Le processus peut se généraliser à un contexte multi-classes quelconque. Supposons que Y possède $k \geq 2$ valeurs possibles (k classes): c_1, \ldots, c_k .
- ▶ On code Y par une variable $\bar{Y} = (Y_1, \dots, Y_k)$ telle que $Y_i = 1$ ssi $Y = c_i$ (1-hot encoding).
- On estime les paramètres $\hat{\boldsymbol{\beta}} = (\hat{\beta}_0, \hat{\beta}_1, \dots, \hat{\beta}_p)$ de manière similaire, ce qui permet l'estimation $\Pr(\bar{\boldsymbol{Y}} = \bar{\boldsymbol{y}} \mid \boldsymbol{X} = \boldsymbol{x})$.
- La prédiction \hat{y} associée à x est alors donnée par (on note 1_i le vecteur de dim k avec un 1 en position i et des 0 partout ailleurs): $\hat{y} = c_i$ si et seulement si

$$\hat{\Pr}(\bar{\pmb{Y}}=\pmb{1}_i\mid X=\pmb{x})>\hat{\Pr}(\bar{\pmb{Y}}=\pmb{1}_j\mid X=\pmb{x}) \text{ pour tout } j\neq i.$$

GENERALIZED ADDITIVE MODELS (GAM)

- Les Generalized Additive Models (GAM) peuvent être adaptés dans le cas de la classification.

$$\log\left(\frac{\Pr(Y=1\mid \boldsymbol{X})}{1-\Pr(Y=1\mid \boldsymbol{X})}\right) = \beta_0 + f_1(X_1) + \dots + f_p(X_p)$$

GENERALIZED ADDITIVE MODELS (GAM)

- Les Generalized Additive Models (GAM) peuvent être adaptés dans le cas de la classification.
- ▶ Dans ce cas, on fat l'hypothèse que la fonction logit est de la forme linéaire suivante:

$$\log \left(\frac{\Pr(Y=1 \mid \boldsymbol{X})}{1 - \Pr(Y=1 \mid \boldsymbol{X})} \right) = \beta_0 + f_1(X_1) + \dots + f_p(X_p)$$

- Les Generalized Additive Models (GAM) peuvent être adaptés dans le cas de la classification.
- Dans ce cas, on fat l'hypothèse que la fonction logit est de la forme linéaire suivante:

$$\log \left(\frac{\Pr(Y=1 \mid \boldsymbol{X})}{1 - \Pr(Y=1 \mid \boldsymbol{X})} \right) = \beta_0 + f_1(X_1) + \dots + f_p(X_p)$$

• où f_1, \ldots, f_p peuvent être différents types de fonctions linéaires ou non-linéaires: polynomial regressions, step functions, basis functions, regression splines, smoothing splines, local regressions... (on ne présentera pas ces méthodes en détail ici)

Les GAM possèdent les avantages suivants:

- Peuvent capturer des relations non-linéaires f_i entre les prédicteurs X_i et la réponse Y
- Donnent donc souvent de meilleures prédictions.
- Restent interprétables: la fonction f_i représente la contribution de X_i à la réponse Y si toutes les autres variables sont fixes.

Les inconvénients:

Les modèles restent *additifs*: les relations non-linaires impliquant plusieurs prédicteurs, e.g. $12X_i^2X_j^3$, ne sont donc pas capturées.

GENERALIZED ADDITIVE MODELS (GAM)

Les GAM possèdent les avantages suivants:

- \triangleright Peuvent capturer des relations non-linéaires f_i entre les prédicteurs X_i et la réponse Y
- Donnent donc souvent de meilleures prédictions.

Les inconvénients:

Les GAM possèdent les avantages suivants:

- Peuvent capturer des relations non-linéaires f_i entre les prédicteurs X_i et la réponse Y
- Donnent donc souvent de meilleures prédictions.
- Restent interprétables: la fonction f_i représente la contribution de X_i à la réponse Y si toutes les autres variables sont fixes.

Les inconvénients:

Les modèles restent *additifs*: les relations non-linaires impliquant plusieurs prédicteurs, e.g. $12X_i^2X_j^3$, ne sont donc pas capturées.

Les GAM possèdent les avantages suivants:

- Peuvent capturer des relations non-linéaires f_i entre les prédicteurs X_i et la réponse Y
- Donnent donc souvent de meilleures prédictions.
- Restent interprétables: la fonction f_i représente la contribution de X_i à la réponse Y si toutes les autres variables sont fixes.

Les inconvénients:

Les modèles restent additifs: les relations non-linaires impliquant plusieurs prédicteurs, e.g. $12X_i^2X_j^3$, ne sont donc pas capturées.

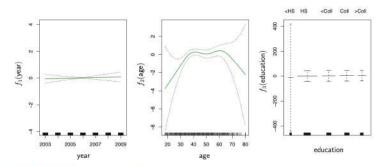


FIGURE 7.13. For the Wage data, the logistic regression GAM given in (7.19) is fit to the binary response I(wage>250). Each plot displays the fitted function and pointwise standard errors. The first function is linear in year, the second function a smoothing spline with five degrees of freedom in age, and the third a step function for education. There are very wide standard errors for the first level <HS of education.

BIBLIOGRAPHIE



James, G., Witten, D., Hastie, T., and Tibshirani, R. (2013).

An Introduction to Statistical Learning: with Applications in R, volume 103 of Springer Texts in Statistics.

Springer, New York.