

# Régression logistique

## Analyse multidimensionnelle appliquée

Léo Belzile

HEC Montréal

automne 2022

La régression logistique spécifie un modèle pour la probabilité de succès

$$p = \Pr(Y = 1 \mid \mathbf{x}) = \frac{1}{1 + \exp(-\eta)}$$

où  $\eta = \beta_0 + \dots + \beta_p x_p$ .

En substituant l'estimation  $\hat{\beta}_0, \dots, \hat{\beta}_p$ , on calcule

- le prédicteur linéaire  $\hat{\eta}_i$  et
- la probabilité de succès  $\hat{p}_i$

pour chaque ligne de la base de données.

Choisir un point de coupure  $c$ :

- si  $\hat{p} < c$ , on assigne  $\hat{Y} = 0$ .
- si  $\hat{p} \geq c$ , on assigne  $\hat{Y} = 1$ .
- Un point de coupure de  $c = 0.5$  revient à assigner l'observation à la classe (catégorie) la plus probable.
- Si  $c = 0$ , on catégorise toutes les observations en succès avec  $\hat{Y}_i = 1$  ( $i = 1, \dots, n$ ).

L'erreur quadratique pour une variable binaire est

$$(Y - \hat{Y})^2 = \begin{cases} 1, & Y \neq \hat{Y}; \\ 0, & Y = \hat{Y}. \end{cases}$$

et donc on obtient le **taux de mauvaise classification** si on calcule la moyenne.

Plus le taux de mauvaise classification est petit, meilleure est la capacité prédictive du modèle.

Utiliser les mêmes données pour l'ajustement et l'estimation de la performance n'est (toujours) pas recommandé.

Plutôt, considérer

- la validation croisée
- la division de l'échantillon

On considère un modèle pour yachat, le fait qu'une personne achète suite à l'envoi d'un catalogue

```
1 data(dbm, package = "hecmulti")
2 formule <- formula("yachat ~ x1 + x2 + x3 +
3                   x4 + x5 + x6 + x7 + x8 + x9 + x10")
4 dbm_class <- dbm |>
5   dplyr::filter(test == 0) |>
6   # pour caret, convertir 0/1 en facteurs
7   dplyr::mutate(yachat = factor(yachat))
```

On utilise la fonction `train` du paquet `caret`, avec le modèle linéaire généralisé

```
1 set.seed(202209)
2 cv_glm <-
3   caret::train(form = formule,
4                 data = dbm_class,
5                 method = "glm",
6                 family = binomial(link = "logit"),
7                 trControl = caret::trainControl(
8                   method = "cv",
9                   number = 10))
```



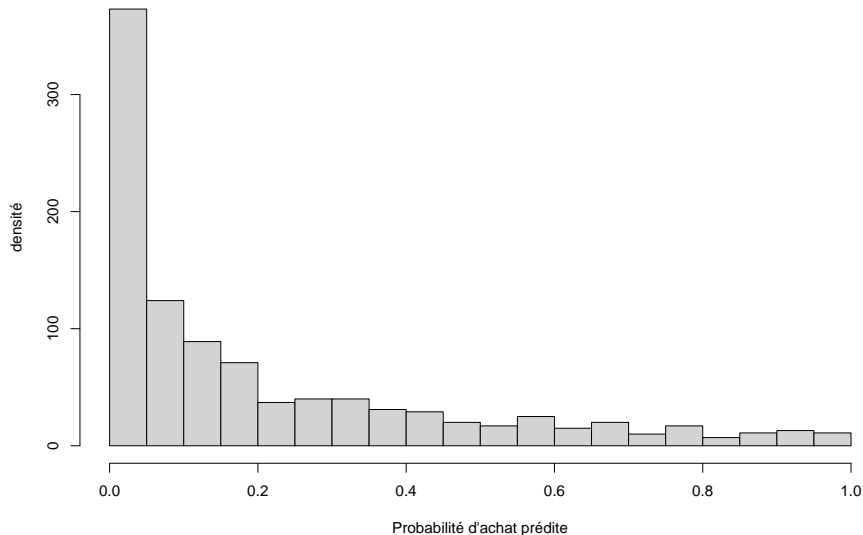


Figure 1: Répartition des probabilités de succès prédites par validation croisée.

On peut varier le point de coupure et regarder pour chaque valeur de  $c$  la classification résultante.

```
1 # predict retourne une matrice n x 2
2 # avec [P(Y=1), P(Y=0)]
3 predprob <- predict(cv_glm, type = "prob")[,1]
4 # Tableau de la performance
5 hecmulti::perfo_logistique(
6   prob = predprob,
7   resp = with(dbm, yachat[test == 0]))
```

On peut classer les observations dans un tableau pour un point de coupure donné.

Table 1: Matrice de confusion avec point de coupure 0.465.

	$Y = 1$	$Y = 0$
$\hat{Y} = 1$	109	52
$\hat{Y} = 0$	101	738