

Data Science for Actuaries (ACT6100)

Arthur Charpentier

Aggrégation # 5.2 (Apprentissage séquentiel)

automne 2020

 <https://github.com/freakonometrics/ACT6100/>

Bagging et Boosting

- ▶ Le **Bagging** vise à construire B modèles à partir de B bases de données obtenues par **bootstrap**. La **construction** de ces B modèles se fait de façon **parallèle**, c'est-à-dire de façon totalement indépendante.
- ▶ Le **Boosting** vise à construire B modèles de façon **séquentielle**: le b -ième modèle dépend du $(b - 1)$ -ième modèle qui dépend lui-même du $(b - 2)$ -ième modèle, etc. Chacun nouveau modèle est construit spécifiquement afin d'améliorer les prédictions faites par le modèle précédent.
- ▶ En français, on pourrait parler (mais en pratique, on ne le fait jamais) d'amplification ou de stimulation.

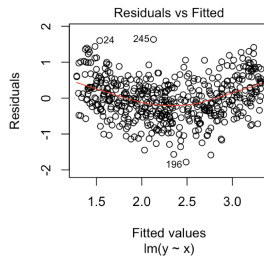
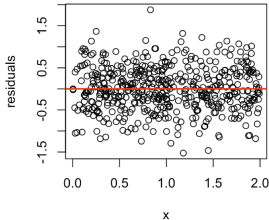
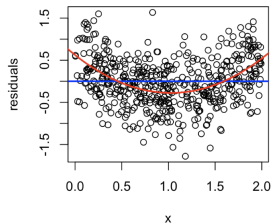
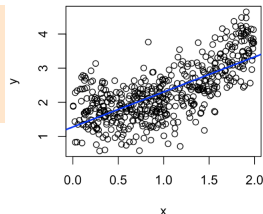
Apprendre de ses erreurs...

```
1 > df = dfr = data.frame(x=x, y=y)
2 > reg1 = lm(y~x, data=df)
3 > dfr$y = residuals(reg1)
4 > reg2 = lm(y~poly(x,2), data=dfr)
```

$$1. y_i = \beta_0 + \beta_1 x_i + r_{1,i}$$

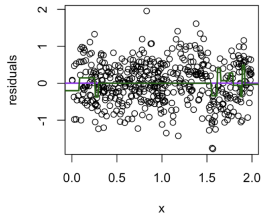
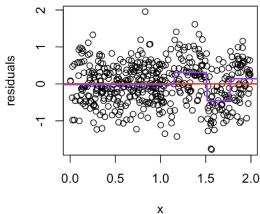
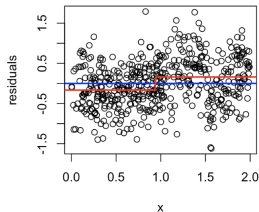
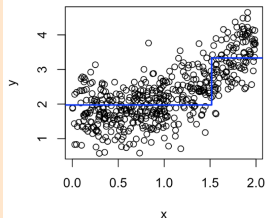
$$2. r_{1,i} = \gamma_0 + \gamma_1 x_i + \gamma_2 x_i^2 + r_{2,i}$$

(etc ?)



Apprendre de ses erreurs...

```
1 > df = dfr = data.frame(x=x, y=y)
2 > reg1 = rpart(y~x, data=df)
3 > dfr$y = residuals(reg1)
4 > reg2 = rpart(y~x, data=dfr)
5 > dfr$y = residuals(reg2)
6 > reg3 = rpart(y~x, data=dfr)
7 > dfr$y = residuals(reg3)
8 > reg4 = rpart(y~x, data=dfr)
```



Boosting

Algorithm 1: Boosting (version 1)

- 1 initialization : k (number of trees), γ , $f_0(\mathbf{x}) = \bar{y}$;
 - 2 **for** $t = 1, 2, \dots, k$ **do**
 - 3 compute $r_{i,t} \leftarrow y_i - f_{t-1}(\mathbf{x}_i)$;
 - 4 fit a model $r_{i,t} \sim h(\mathbf{x}_i)$ for some tree h ;
 - 5 update $f_t(\cdot) = f_{t-1}(\cdot) + \gamma h(\cdot)$
-

Algorithm 2: Boosting (version 2)

- 1 initialization : k (number of trees), $f_0(\mathbf{x}) = \underset{\gamma}{\operatorname{argmin}} \sum_{i=1}^n \ell(y_i, \gamma)$;
 - 2 **for** $t = 1, 2, \dots, k$ **do**
 - 3 compute $r_{i,t} \leftarrow \left. \frac{\partial \ell(y_i, \hat{y})}{\partial \hat{y}} \right|_{\hat{y}=f_{t-1}(\mathbf{x}_i)}$;
 - 4 fit a model $r_{i,t} \sim h(\mathbf{x}_i)$ for some tree h ;
 - 5 update $f_t(\cdot) = f_{t-1}(\cdot) + \gamma h(\cdot)$
-

Algorithme AdaBoost

ou voir la mise à jour à l'aide de poids (si la classification n'est pas bonne)

On s'intéresse à un problème de classification binaire (0 ou 1) à partir d'une base de données $\mathcal{D} = \{Y_i; \mathbf{X}_i\}_{i=1,\dots,n}$.

1. Toutes les observations reçoivent un poids identique:
 $p_1(\mathbf{X}_i) = 1/n, i = 1, \dots, n.$

Algorithme AdaBoost (suite)

2. Pour l'étape $t = 1, \dots, T$,
- 2a. Piger au hasard (en utilisant les poids $p_t(\mathbf{X}_i)$) une base de données d'entrainement $\mathcal{D}_t = (\mathcal{D}, p_t)$.
 - 2b. Apprendre une règle de classification r_t sur cette base de données.
 - 2c. Calculer l'erreur apparente

$$\epsilon_t = p_t(r_t(\mathbf{X}_i) \neq Y_i) = \sum_{i: r_t(\mathbf{X}_i) \neq Y_i} p_t(i)$$

pour \mathcal{D}_t et évaluer $\alpha_t = 0.5 \ln((1 - \epsilon_t)/\epsilon_t)$.

- 2d. Mettre à jour les poids:

$$p_{t+1}(\mathbf{X}_i) = \begin{cases} \frac{p_t(\mathbf{X}_i)}{Z_t} e^{-\alpha_t}, & r_t(\mathbf{X}_i) = Y_i \\ \frac{p_t(\mathbf{X}_i)}{Z_t} e^{+\alpha_t}, & r_t(\mathbf{X}_i) \neq Y_i, \end{cases}$$

où Z_t est une constante de normalisation assurant que $\sum_{i=1}^n p_{t+1}(\mathbf{X}_i) = 1$.

Algorithme AdaBoost (suite)

3. La classification finale faite par le modèle est alors

$$R(\mathbf{x}_i) = \begin{cases} 1, & \sum_{t=1}^T \alpha_t r_t(\mathbf{x}_i) > 0 \\ -1, & \sum_{t=1}^T \alpha_t r_t(\mathbf{x}_i) < 0. \end{cases}$$

Algorithm 3: Boosting (Adaboost)

```
1 initialization :  $k, \omega_i \leftarrow 1/n$ ;  
2 for  $t = 1, 2, \dots, k$  do  
3   error rate  $e_t \leftarrow \frac{\sum_{i=1}^n \omega_i \mathbf{1}(y_i \neq f_{t-1}(\mathbf{x}_i))}{\sum_{i=1}^n \omega_i}$ ;  
4   set  $\alpha_t \leftarrow \log(1 - e_t) - \log(e_t)$ ;  
5   update  $\omega_i \leftarrow \omega_i e^{\alpha_t \mathbf{1}(y_i \neq f_{t-1}(\mathbf{x}_i))}$ ;  
6   update  $f_t(\cdot) \leftarrow f_{t-1}(\cdot) + \alpha_t h(\cdot)$ 
```

Remarques

On apprend tant que $\alpha_t \geq 0$, soit $e_t < 1/5$.

Modèle classique pour h : CART (avec des poids)

► Approche théorique

Étant donné une famille de modèles \mathcal{H} , on cherche à résoudre

$$h^* = \operatorname{argmin}_{h \in \mathcal{H}} \left\{ \mathbb{E}[\ell(Y, h(\mathbf{X}))] \right\}$$

ou sa version empirique

$$\hat{h}^* = \operatorname{argmin}_{h \in \mathcal{H}} \left\{ \frac{1}{n} \sum_{i=1}^n \ell(y_i, h(\mathbf{x}_i)) \right\}$$

avec classiquement $\ell(y, \hat{y}) = \mathbf{1}(y \neq \hat{y})$.

On peut tenter de **convexifier la fonction de perte**,

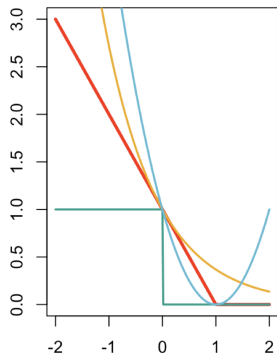
$$\bar{\ell}(y, \hat{y}) = \exp(-y\hat{y}).$$

Remarques

Supposons $y \in \{-1, +1\}$.

Les points sont mal classés si $y \cdot f(\mathbf{x}) < 0$,
on peut tracer la fonction de perte en fonction
de $y\hat{y}$

- ▶ misclassification: $\mathbf{1}(y\hat{y} < 0)$ —————
- ▶ hinge: $(1 - y\hat{y})_+$ —————
- ▶ exponential: $\exp(-y\hat{y})$ —————



L'algorithme Adaboost vérifie $f_t = f_{t-1} + 2\beta^*h^*$, où (β^*, h^*) est
solution de

$$(\beta^*, h^*) = \operatorname{argmin}_{(\beta, h)} \left\{ \sum_{i=1}^n \bar{\ell}(y_i, f_{t-1}(\mathbf{x}_i) + \beta h(\mathbf{x}_i)) \right\}$$

Remarques

$$(\beta^*, h^*) = \operatorname{argmin}_{(\beta, h)} \left\{ \sum_{i=1}^n \bar{\ell}(y_i, f_{t-1}(\mathbf{x}_i) + \beta h(\mathbf{x}_i)) \right\}$$

$$(\beta^*, h^*) = \operatorname{argmin}_{(\beta, h)} \left\{ \sum_{i=1}^n \underbrace{\exp(-y_i f_{t-1}(\mathbf{x}_i))}_{\omega_i} \exp(-y_i \beta h(\mathbf{x}_i)) \right\}$$

on distingue alors les i pour lesquels $y_i = h(\mathbf{x}_i)$ et $y_i \neq h(\mathbf{x}_i)$, et

$$h^* = \operatorname{argmin}_{h \in \mathcal{H}} \left\{ \frac{1}{n} \sum_{i=1}^n \omega_i \ell(y_i, h(\mathbf{x}_i)) \right\}$$

Remarques

- ▶ L'erreur apparente est calculée en utilisant la distribution avec laquelle le modèle est entraîné.
- ▶ Chacune des nouvelles règles r_t est pondérée par un poids α_t qui indique l'importance que cette règle aura dans la décision finale R .
- ▶ Si T , le nombre de modèles, est *trop* grand, alors il y a un risque que l'algorithme se concentre trop vers la fin sur les cas difficiles (potentiellement du bruit) et accorde à ces derniers trop d'importance dans la décision finale R .

► L'erreur apparente

$$\epsilon_t = p_t(r_t(\mathbf{X}_i) \neq Y_i) = \sum_{i:r_t(\mathbf{X}_i) \neq Y_i} p_t(i)$$

peut être réécrite comme

$$\epsilon_t = 0.5 - \gamma_t,$$

où γ_t est l'amélioration apportée à la prédiction par la règle r_t en comparaison avec le hasard (0.5).

Bornes (suite)

- ▶ On peut démontrer que l'erreur apparente finale (pour T) a comme borne supérieure

$$\epsilon_T \leq \exp \left(-2 \sum_t \gamma_t^2 \right) \leq \exp \left(-2T\gamma^2 \right),$$

où $\gamma = \min(\gamma_1, \dots, \gamma_T)$.

- ▶ Ainsi, si le **weak learner** fait légèrement mieux que le hasard, $\gamma_t > 0$ et l'erreur apparente finale diminue exponentiellement avec le nombre de modèles (T).

Bornes (suite)

- ▶ On peut également démontrer que l'erreur de généralisation, c'est-à-dire l'erreur calculée sur une base de données qui n'a pas été utilisée pour l'entraînement du modèle, a comme borne supérieure

$$\text{Err}_G \leq \epsilon_T + \sqrt{\frac{(T)(d_{\mathcal{H}})}{n}},$$

où $d_{\mathcal{H}}$ est un terme lié à la taille de l'espace d'hypothèses.

- ▶ Cette borne indique que si T devient grand, alors le **Boosting** devrait tendre à surapprendre.
- ▶ En pratique, on observe rarement ce phénomène.

Fonction objectif

1. De façon générale, on cherche une règle \hat{r} en résolvant un problème d'optimisation empirique

$$\hat{r} = \operatorname{argmin}_{r \in \mathcal{R}} \frac{1}{n} \sum_{i=1}^n \ell(Y_i, r(\mathbf{X}_i)),$$

où ℓ est une fonction de perte quelconque.

2. Si l'optimisation est faite sur l'ensemble des fonctions r possibles et imaginables, le problème est trop complexe pour être résolu (même numériquement!).

Fonction objectif

- ▶ Une solution possible est de convexifier la fonction objectif afin de rendre le problème plus simple d'un point de vue numérique.
- ▶ L'algorithme **AdaBoost** répond à ce principe de minimisation pour un risque convexifié donné par

$$\ell(Y_i, r(\mathbf{X}_i)) = \exp(-Y_i r(\mathbf{X}_i)).$$

- ▶ On va également réaliser l'optimisation de façon séquentielle plutôt que globale (**greedy**).

Fonction objectif

On peut démontrer que la règle de classification à l'étape t de l'algorithme **AdaBoost** peut s'écrire comme étant

$$r_t(\mathbf{X}_i) = r_{t-1}(\mathbf{X}_i) + \alpha_t R_t(\mathbf{X}_i),$$

où

$$(\alpha_t, R_t) = \underset{(\alpha, r) \in \mathbb{R} \times \mathcal{R}}{\operatorname{argmin}} \sum_{i=1}^n e^{-Y_i(r_{t-1}(\mathbf{X}_i) + \alpha r(\mathbf{X}_i))}$$

et où \mathcal{R} est l'espace des choix possibles pour la règle r .

Boosting par descente de gradient fonctionnelle

1. Initialiser

$$r_0(\mathbf{X}_i) = \operatorname{argmin}_{c \in \mathbb{R}} \frac{1}{n} \sum_{i=1}^n \ell(Y_i, c).$$

2. Pour les étapes $t = 1, \dots, T$,

2a. Calculer

$$U_i = - \left[\frac{\partial L(Y_i, r(\mathbf{X}_i))}{\partial r(\mathbf{X}_i)} \right]_{\{r(\mathbf{X}_i) = r_{t-1}(\mathbf{X}_i)\}}, \quad i = 1, \dots, n.$$

2b. Ajuster le **weak learner** sélectionné sur l'échantillon composé des éléments $(\mathbf{X}_1, U_1), (\mathbf{X}_2, U_2), \dots, (\mathbf{X}_n, U_n)$ afin d'obtenir la pseudo-règle $h_t(\mathbf{X}_i)$.

2c. Effectuer la mise à jour $r_t(\mathbf{X}_i) = r_{t-1}(\mathbf{X}_i) + \alpha h_t(\mathbf{X}_i)$.

3. La prédiction finale est donnée par $r_T(\mathbf{X}_i)$.

Remarques

- ▶ Le choix du paramètre α (taux d'apprentissage) est peu important. On recommande généralement de prendre une petite valeur (0.01 ou 0.001).
- ▶ Si on pose $\alpha = 1$ et $\ell(y, \hat{y}) = \exp(-y\hat{y})$, on retrouve (presque) l'algorithme **AdaBoost**.

Gradient Boosting with R, ℓ_2 loss

```
1 > loss_L2 = function(y, yhat) (mean(1/2*(y-yhat)^2))
2 > gradient_L2 = function(y, yhat) {(y-yhat)}
```

then use

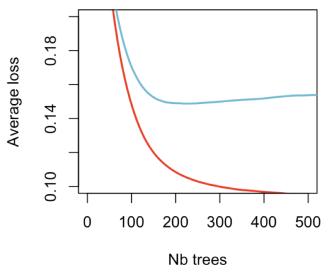
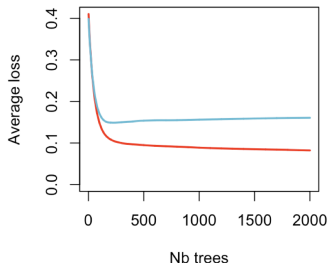
```
1 > grad_boost_train_valid <- function(formula,
    data_train, data_valid, nu = 0.01, stop=1000, =
    loss_L2) {
2   data = as.data.frame(data_train)
3   formula = terms.formula(formula)
4   noms_X = names(data_train)[-which(names(data_train)
    ==as.character(formula)[2])]
5   noms_Y = names(data_train)[which(names(data_train)==
    as.character(formula)[2])]
6   y = data_train[, noms_Y]
7   fit_train = fit_valid = mean(y)
8   data = as.data.frame(data_train[,noms_X])
9   names(data) = noms_X
10  data$u = grad.fun(y = y, yhat = fit_train)
11  loss_train = loss_valid = c()
```

Gradient Boosting with R, ℓ_2 loss

```
1  for (i in 1:stop) {
2    model_weak = rpart(u ~ ., data=data, maxdepth=3,
3                      cp=1e-9)
4    fit_train = as.numeric(fit_train + nu * predict(
5      model_weak, newdata=data_train))
6    fit_valid = fit_valid + nu * predict(model_weak,
7      newdata=data_valid)
8    data$u = as.numeric(unlist(grad.fun(y = as.numeric(
9      unlist(data_train[noms_Y])), yhat = fit_train)))
10   loss_train <- c(loss_train, loss.fun(y = as.
11     numeric(unlist(data_train[noms_Y])), yhat =
12       fit_train))
13   loss_valid <- c(loss_valid, loss.fun(y = as.
14     numeric(unlist(data_valid[noms_Y])), yhat =
15       fit_valid))
16 }
17 loss_matrix = data.frame(N = 1:stop, TRAIN =
18   loss_train, VALID = loss_valid)
19 return(list(loss = loss_matrix))
20 }
```

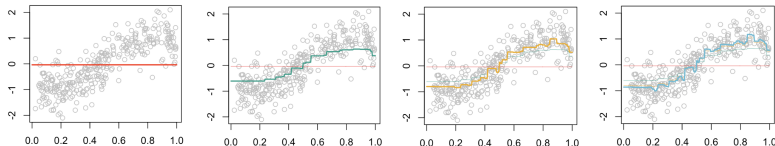
Gradient Boosting with R, ℓ_2 loss

```
1 > x=sort(runif(n))
2 > y=sin(-3*pi/4+x*pi*3/2)+rnorm(n)/2
3 DF = data.frame(x=x,y=y)
4 > set.seed(1)
5 > idx = sample(1:n,size = 2*n/3)
6 > library(rpart)
7 R = grad_boost_train_valid(y~x, data_train=DF[idx,],
  data_valid=DF[-idx,], nu = .01, stop=2000, grad.fun
  = gradient_L2, loss.fun = loss_L2)
```



Gradient Boosting with R, ℓ_2 loss

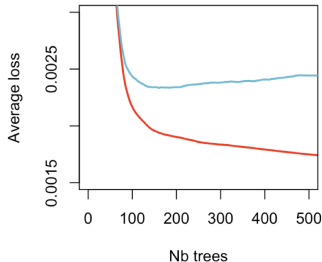
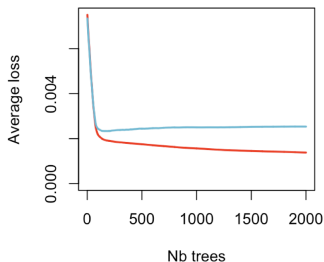
```
1 > x=sort(runif(n))
2 > y=sin(-3*pi/4+x*pi*3/2)+rnorm(n)/2
3 DF = data.frame(x=x,y=y)
4 > set.seed(1)
5 > idx = sample(1:n,size = 2*n/3)
6 > library(rpart)
7 R = grad_boost_train_valid(y~x, data_train=DF[idx,],
  data_valid=DF[-idx,], nu = .01, stop=2000,grad.fun
  = gradient_L2, loss.fun = loss_L2)
```



Sequential learning, with 0, 100, 200, 500 trees

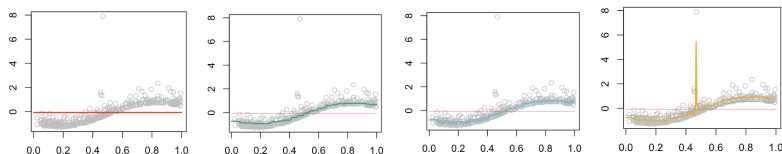
Gradient Boosting with R, ℓ_1 (or Huber) loss

```
1 > loss_L_huber = function(y, yhat, delta=.01){
2   L=1/2*(y-yhat)^2*(abs(y-yhat)<=delta) + delta*(abs(y
   -yhat)-delta/2)*(abs(y-yhat)>delta)
3   return(mean(L))}
4 > gradient_L_huber = function(y, yhat, delta=.01){
5   L=(y-yhat)*(abs(y-yhat)<=delta) + delta*sign(y-yhat)
   *(abs(y-yhat)>delta)
6   return(as.numeric(L))}
7 > y = sin(-3*pi/4+x*pi*3/2)
8 > e = rlnorm(n)
9 > DF = data.frame(x=x,y=y+(e-mean(e))/sd(e)/2)
```



Gradient Boosting with R , ℓ_1 (or Huber) loss

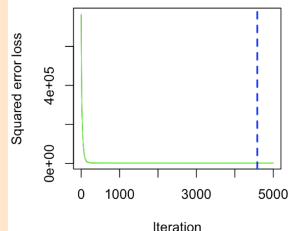
```
1 > library(rpart)
2 R = grad_boost_train_valid(y~x, data_train=DF[idx,],
  data_valid=DF[-idx,], nu = .01, stop=2000, grad.fun
  = gradient_L_huber, loss.fun = loss_L_huber)
```



Sequential learning, with 0, 100, 1000 trees,
(and 1000 trees with some ℓ_2 loss on the right)

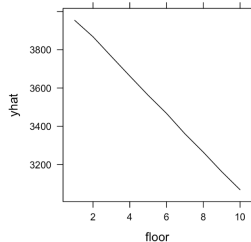
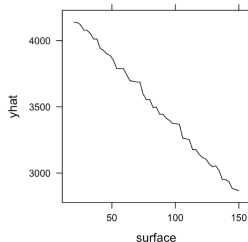
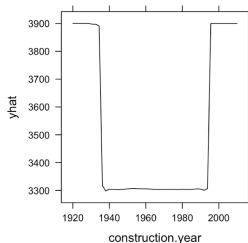
Gradient Boosting with R, ℓ_2 loss

```
1 > library(gbm)
2 > library(DALEX)
3 > G = gbm(
4   formula = m2.price ~ .,
5   data = apartments,
6   distribution = "gaussian",
7   n.trees = 5000,
8   shrinkage = 0.1,
9   cv.folds = 5
10 )
11 > gbm.perf(G, method = "cv")
12 > min_MSE <- which.min(G$cv.error)
13 > sqrt(G$cv.error[min_MSE])
14 [1] 50.67178
```



Gradient Boosting with R, ℓ_2 loss

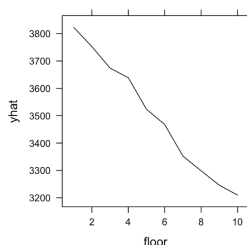
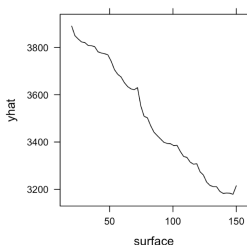
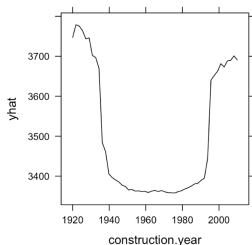
```
1 > library(pdp)
2 > pdp::partial(G, pred.var = "floor", plot = TRUE, n.
  trees=2000)
```



Gradient Boosting with R, ℓ_2 loss

one can compare with a random forest

```
1 > library(randomForest)
2 > RF = randomForest(m2.price ~ ., data = apartments)
3 > pdp::partial(G, pred.var = "floor", plot = TRUE, n.trees=2000)
```



GLM & boosting

Consider a logistic regression, with data (y, \mathbf{x}) and

$$\mathbb{P}(Y = y | \mathbf{X} = \mathbf{x}) = p(\mathbf{x})^y [1 - p(\mathbf{x})]^{1-y},$$

with

$$p(\mathbf{x}) = \frac{e^{\mathbf{x}^\top \beta}}{1 + e^{\mathbf{x}^\top \beta}} = \frac{1}{1 + e^{-\mathbf{x}^\top \beta}}$$

Here, given $\psi : \mathbb{R}^p \rightarrow \mathbb{R}$, consider

$$p(\mathbf{x}) = \frac{e^{\psi(\mathbf{x})}}{e^{-\psi(\mathbf{x})} + e^{\psi(\mathbf{x})}} = \frac{1}{1 + e^{-2\psi(\mathbf{x})}}$$

i.e.

$$\psi(\mathbf{x}) = \frac{1}{2} \log \left(\frac{p(\mathbf{x})}{1 - p(\mathbf{x})} \right)$$

We want to find ψ^* that solves $\psi^* = \operatorname{argmin} \{ -\log \mathcal{L} \}$, i.e

$$\psi^* = \operatorname{argmin} \left\{ \frac{1}{n} \sum_{i=1}^n \log [1 + \exp[-2\tilde{y}_i \psi(\mathbf{x}_i)]] \right\}$$

where $\tilde{y}_i = 2y_i - 1 \in \{-1, +1\}$.

Poisson Regression & boosting

Recall that $Y \sim \mathcal{P}(\lambda)$ if $\mathbb{P}[Y = y] = \frac{e^{-\lambda} \lambda^y}{y!}$

In the Poisson regression, with a log-link function

$$\mathbb{P}[Y = y | \mathbf{X} = \mathbf{x}] = \frac{e^{-\lambda_{\mathbf{x}}} \lambda_{\mathbf{x}}^y}{y!}, \quad \text{with } \lambda_{\mathbf{x}} = e^{\mathbf{x}^\top \boldsymbol{\beta}}$$

The log-likelihood for a sample (y_i, \mathbf{x}_i) , $i = 1, \dots, n$ is

$$\log \mathcal{L} = \sum_{i=1}^n -\exp[\mathbf{x}_i^\top \boldsymbol{\beta}] + y_i \mathbf{x}_i^\top \boldsymbol{\beta} - \log y_i!$$

As previously, assume that $\lambda = \exp[\psi(\mathbf{x})]$ so that

$$\log \mathcal{L} = \sum_{i=1}^n -\exp[\psi(\mathbf{x}_i)] + y_i \psi(\mathbf{x}_i) - \log y_i! = -\sum_{i=1}^n \ell(y_i, \lambda_i)$$

Poisson Regression & boosting

then

$$-\frac{\partial \ell(y_i, \lambda_i)}{\partial \psi(\mathbf{x}_i)} = -\exp[\psi(\mathbf{x}_i)] + y_i$$

Let $r_i = y_i - \exp[\psi(\mathbf{x}_i)]$, we want to find the step direction that best fits $-\mathbf{r}$, in the least-square sense,

$$h^* = \operatorname{argmin} \left\{ \sum_{i=1}^n [-r_i - h(\mathbf{x}_i)]^2 \right\}$$

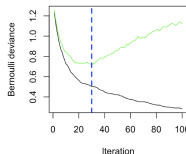
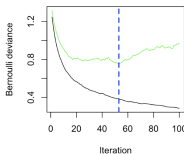
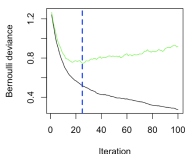
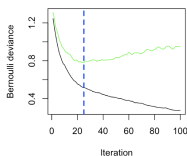
$$h^* = \operatorname{argmin} \left\{ \sum_{i=1}^n [\exp[\psi(\mathbf{x}_i)] - y_i - h(\mathbf{x}_i)]^2 \right\}$$

e.g. with a regression tree

Gradient Boosting with R, logistic loss

one can compare with a random forest

```
1 > G <- gbm(formula = (PRONO=="SURVIE") ~ ., data =  
  myocarde, distribution = "bernoulli", n.trees =  
    100, shrinkage = 0.1, cv.folds = 5)  
2 > gbm.perf(G, method = "cv")
```



```
1 > G  
2 gbm(formula = (PRONO == "SURVIE") ~ ., distribution =  
  "bernoulli")  
3 A gradient boosted model with bernoulli loss function.  
4 100 iterations were performed.  
5 The best cross-validation iteration was 30.  
6 There were 7 predictors of which 7 had non-zero  
  influence.
```

Gradient Boosting with R, logistic loss

```
1 > summary(G)
2 > pdp::partial(G, pred.var = "INSYS", plot = TRUE, n.trees=25)
```

