

Actuariat de l'Assurance Non-Vie # 9

A. Charpentier (Université de Rennes 1)

ENSAE 2017/2018



credit: Arnold Odermatt

Fourre-Tout sur la Tarification

- modèle collectif vs. modèle individuel
- cas de la grande dimension
- choix de variables
- choix de modèles

Modèle *individuel* ou modèle *collectif* ? La loi Tweedie

Consider a Tweedie distribution, with variance function power $p \in (1, 2)$, mean μ and scale parameter ϕ , then it is a compound Poisson model,

- $N \sim \mathcal{P}(\lambda)$ with $\lambda = \frac{\phi\mu^{2-p}}{2-p}$
- $Y_i \sim \mathcal{G}(\alpha, \beta)$ with $\alpha = -\frac{p-2}{p-1}$ and $\beta = \frac{\phi\mu^{1-p}}{p-1}$

Conversely, consider a compound Poisson model $N \sim \mathcal{P}(\lambda)$ and $Y_i \sim \mathcal{G}(\alpha, \beta)$, then

- variance function power is $p = \frac{\alpha+2}{\alpha+1}$
- mean is $\mu = \frac{\lambda\alpha}{\beta}$
- scale parameter is $\phi = \frac{[\lambda\alpha]^{\frac{\alpha+2}{\alpha+1}-1} \beta^{2-\frac{\alpha+2}{\alpha+1}}}{\alpha+1}$

Modèle *individuel* ou modèle *collectif* ? La régression Tweedie

In the context of regression

$$N_i \sim \mathcal{P}(\lambda_i) \text{ with } \lambda_i = \exp[\mathbf{X}_i^\top \boldsymbol{\beta}_\lambda]$$

$$Y_{j,i} \sim \mathcal{G}(\mu_i, \phi) \text{ with } \mu_i = \exp[\mathbf{X}_i^\top \boldsymbol{\beta}_\mu]$$

Then $S_i = Y_{1,i} + \dots + Y_{N,i}$ has a Tweedie distribution

- variance function power is $p = \frac{\phi + 2}{\phi + 1}$
- mean is $\lambda_i \mu_i$
- scale parameter is $\frac{\lambda_i^{\frac{1}{\phi+1}-1}}{\mu_i^{\frac{\phi}{\phi+1}}} \left(\frac{\phi}{1 + \phi} \right)$

There are $1 + 2\dim(\mathbf{X})$ degrees of freedom.

Modèle *individuel* ou modèle *collectif* ? La régression Tweedie

Remark Note that the scale parameter should not depend on i .

A Tweedie regression is

- variance function power is $p \in (1, 2)$
- mean is $\mu_i = \exp[\mathbf{X}_i^\top \boldsymbol{\beta}_{\text{Tweedie}}]$
- scale parameter is ϕ

There are $2 + \dim(\mathbf{X})$ degrees of freedom.

Double Modèle Fréquence - Coût Individuel

Considérons les bases suivantes, en RC, pour la fréquence

```
1 > freq = merge(contrat, nombre_RC)
```

pour les coûts individuels

```
1 > sinistre_RC = sinistre[(sinistre$garantie=="1RC") & (sinistre$cout > 0),]
```

```
2 > sinistre_RC = merge(sinistre_RC, contrat)
```

et pour les coûts agrégés par police

```
1 > agg_RC = aggregate(sinistre_RC$cout, by=list(sinistre_RC$nocontrat), FUN='sum')
```

```
2 > names(agg_RC)=c('nocontrat', 'cout_RC')
```

```
3 > global_RC = merge(contrat, agg_RC, all.x=TRUE)
```

```
4 > global_RC$cout_RC[is.na(global_D0$cout_RC)]=0
```

Double Modèle Fréquence - Coût Individuel

```
1 > library(splines)
2 > reg_f = glm(nb_RC~zone+bs(ageconducteur)+carburant, offset=log(
    exposition), data=freq, family=poisson)
3 > reg_c = glm(cout~zone+bs(ageconducteur)+carburant, data=sinistre_RC
    , family=Gamma(link="log"))
```

Simple Modèle Coût par Police

```
1 > library(tweedie)
2 > library(statmod)
3 > reg_a = glm(cout_RC~zone+bs(ageconducteur)+carburant, offset=log(
    exposition), data=global_RC, family=tweedie(var.power=1.5, link.
    power=0))
```

Comparaison des primes

```

1 > freq2 = freq
2 > freq2$exposition = 1
3 > P_f = predict(reg_f,newdata=freq2,type="response")
4 > P_c = predict(reg_c,newdata=freq2,type="response")
5 prime1 = P_f*P_c

1 > k = 1.5
2 > reg_a = glm(cout_D0~zone+bs(ageconducteur)+carburant, offset=log(
    exposition),data=global_D0,family=tweedie(var.power=k, link.power
    =0))
3 > prime2 = predict(reg_a,newdata=freq2,type="response")

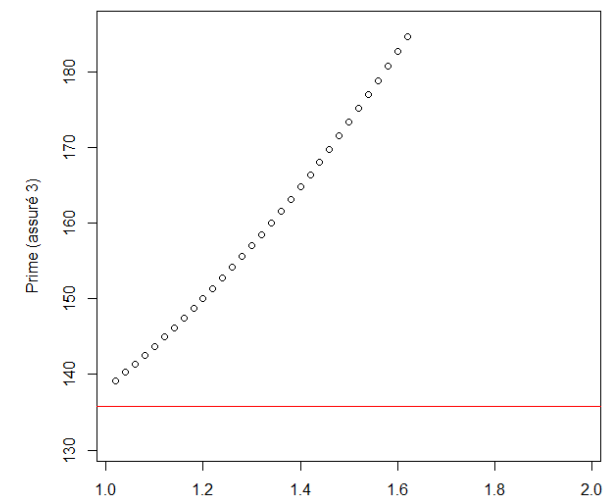
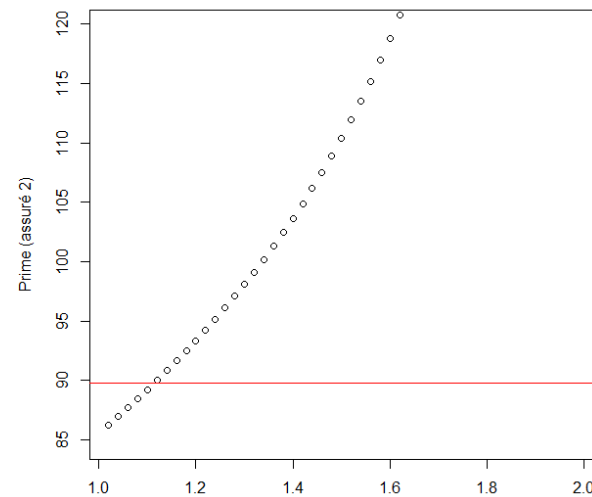
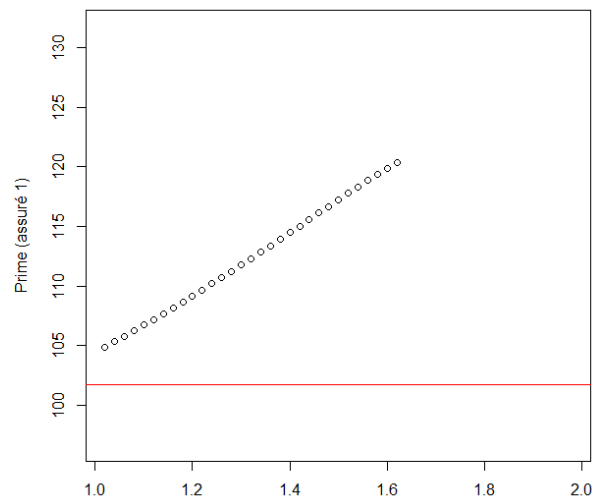
1 > arrows(1:100,prime1[1:100],1:100,prime2[1:100],length=.1)

```


Impact du degré Tweedie sur les Primes Pures

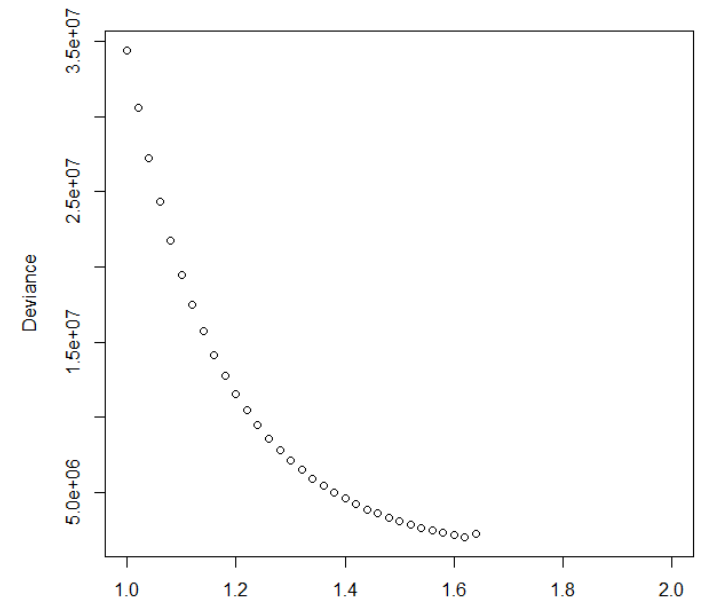
Impact du degré Tweedie sur les Primes Pures

Comparaison des primes pures, assurés no1, no2 et no 3 (DO)



'Optimisation' du Paramètre Tweedie

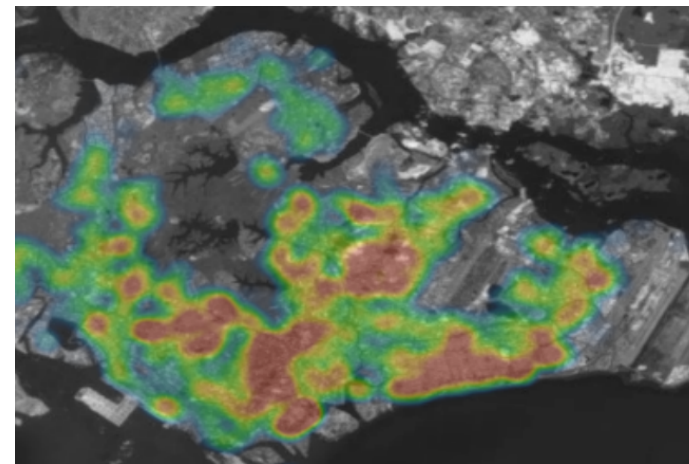
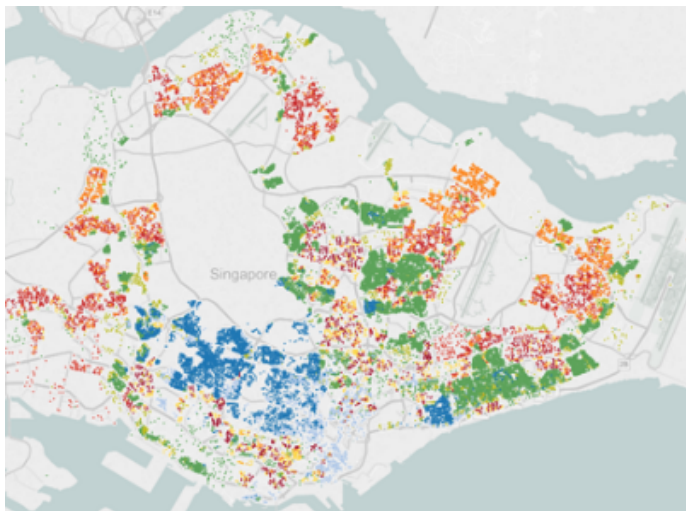
```
1 > dev = function(k){  
2 + reg = glm(cout_RC~zone+bs(ageconducteur)+  
    carburant, data=global_RC, family=  
    tweedie(var.power=k, link.power=0),  
    offset=log(exposition))  
3 + reg$deviance  
4 + }
```



Tarification et données massives (*Big Data*)

Problèmes classiques avec des données massives

- beaucoup de variables explicatives, k grand, $\mathbf{X}^T \mathbf{X}$ peut-être non inversible
- gros volumes de données, e.g. données télématiques
- données non quantitatives, e.g. texte, localisation, etc.



La fascination pour les estimateurs sans biais

En statistique mathématique, on aime les estimateurs sans biais car ils ont plusieurs propriétés intéressantes. Mais ne peut-on pas considérer des estimateurs biaisés, potentiellement meilleurs ?

Consider a sample, i.i.d., $\{y_1, \dots, y_n\}$ with distribution $\mathcal{N}(\mu, \sigma^2)$. Define $\hat{\theta} = \alpha \bar{Y}$. What is the optimal α^* to get the **best estimator of μ** ?

- bias: $\text{bias}(\hat{\theta}) = \mathbb{E}(\hat{\theta}) - \mu = (\alpha - 1)\mu$
- variance: $\text{Var}(\hat{\theta}) = \frac{\alpha^2 \sigma^2}{n}$
- mse: $\text{mse}(\hat{\theta}) = (\alpha - 1)^2 \mu^2 + \frac{\alpha^2 \sigma^2}{n}$

The optimal value is $\alpha^* = \frac{\mu^2}{\mu^2 + \frac{\sigma^2}{n}} < 1$.

Linear Model

Consider some linear model $y_i = \mathbf{x}_i^\top \boldsymbol{\beta} + \varepsilon_i$ for all $i = 1, \dots, n$.

Assume that ε_i are i.i.d. with $\mathbb{E}(\varepsilon) = 0$ (and finite variance). Write

$$\underbrace{\begin{pmatrix} y_1 \\ \vdots \\ y_n \end{pmatrix}}_{\mathbf{y}, n \times 1} = \underbrace{\begin{pmatrix} 1 & x_{1,1} & \cdots & x_{1,k} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & x_{n,1} & \cdots & x_{n,k} \end{pmatrix}}_{\mathbf{X}, n \times (k+1)} \underbrace{\begin{pmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_k \end{pmatrix}}_{\boldsymbol{\beta}, (k+1) \times 1} + \underbrace{\begin{pmatrix} \varepsilon_1 \\ \vdots \\ \varepsilon_n \end{pmatrix}}_{\boldsymbol{\varepsilon}, n \times 1}.$$

Assuming $\boldsymbol{\varepsilon} \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbb{I})$, the maximum likelihood estimator of $\boldsymbol{\beta}$ is

$$\hat{\boldsymbol{\beta}} = \operatorname{argmin}\{\|\mathbf{y} - \mathbf{X}^\top \boldsymbol{\beta}\|_{\ell_2}\} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}$$

... under the assumption that $\mathbf{X}^\top \mathbf{X}$ is a full-rank matrix.

What if $\mathbf{X}_i^\top \mathbf{X}$ cannot be inverted? Then $\hat{\boldsymbol{\beta}} = [\mathbf{X}^\top \mathbf{X}]^{-1} \mathbf{X}^\top \mathbf{y}$ does not exist, but $\hat{\boldsymbol{\beta}}_\lambda = [\mathbf{X}^\top \mathbf{X} + \lambda \mathbb{I}]^{-1} \mathbf{X}^\top \mathbf{y}$ always exist if $\lambda > 0$.

Ridge Regression

The estimator $\hat{\beta} = [\mathbf{X}^\top \mathbf{X} + \lambda \mathbb{I}]^{-1} \mathbf{X}^\top \mathbf{y}$ is the **Ridge** estimate obtained as solution of

$$\hat{\beta} = \underset{\beta}{\operatorname{argmin}} \left\{ \sum_{i=1}^n [y_i - \beta_0 - \mathbf{x}_i^\top \beta]^2 + \underbrace{\lambda \|\beta\|_{\ell_2}}_{\mathbf{1}^\top \beta^2} \right\}$$

for some tuning parameter λ . One can also write

$$\hat{\beta} = \underset{\beta; \|\beta\|_{\ell_2} \leq s}{\operatorname{argmin}} \{ \|\mathbf{Y} - \mathbf{X}^\top \beta\|_{\ell_2} \}$$

Remark Note that we solve $\hat{\beta} = \underset{\beta}{\operatorname{argmin}} \{ \text{objective}(\beta) \}$ where

$$\text{objective}(\beta) = \underbrace{\mathcal{L}(\beta)}_{\text{training loss}} + \underbrace{\mathcal{R}(\beta)}_{\text{regularization}}$$

Going further on sparsity issues

In several applications, k can be (very) large, but a lot of features are just noise: $\beta_j = 0$ for many j 's. Let s denote the number of relevant features, with $s \ll k$, cf [Hastie, Tibshirani & Wainwright \(2015\)](#),

$$s = \text{card}\{\mathcal{S}\} \text{ where } \mathcal{S} = \{j; \beta_j \neq 0\}$$

The model is now $y = \mathbf{X}_{\mathcal{S}}^{\top} \boldsymbol{\beta}_{\mathcal{S}} + \varepsilon$, where $\mathbf{X}_{\mathcal{S}}^{\top} \mathbf{X}_{\mathcal{S}}$ is a full rank matrix.

Going further on sparsity issues

Define $\|\mathbf{a}\|_{\ell_0} = \sum \mathbf{1}(|a_i| > 0)$. Ici $\dim(\beta) = s$.

We wish we could solve

$$\hat{\beta} = \underset{\beta; \|\beta\|_{\ell_0} \leq s}{\operatorname{argmin}} \{ \|\mathbf{Y} - \mathbf{X}^\top \beta\|_{\ell_2} \}$$

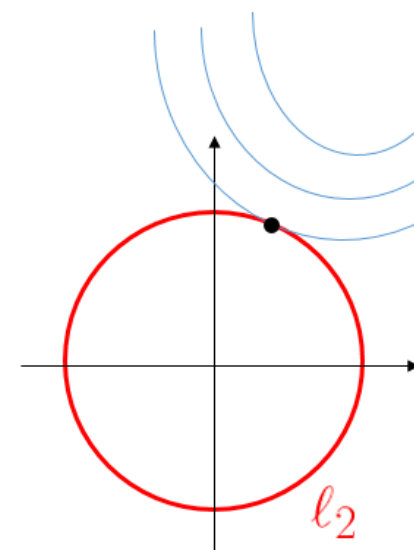
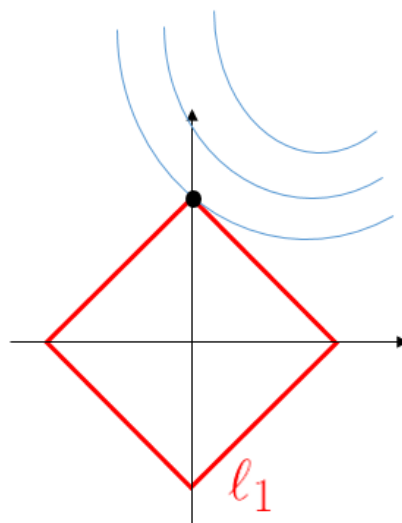
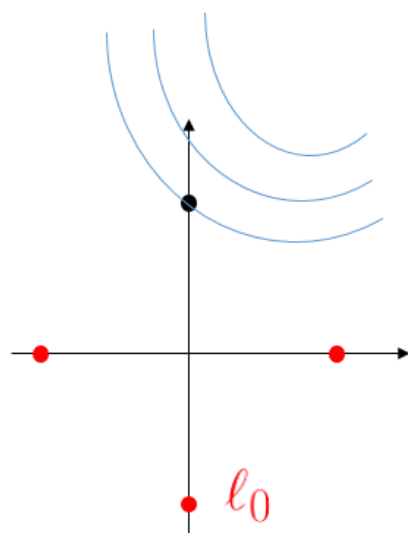
Problem: it is usually not possible to describe all possible constraints, since $\binom{s}{k}$ coefficients should be chosen here (with k (very) large).

Idea: solve the dual problem

$$\hat{\beta} = \underset{\beta; \|\mathbf{Y} - \mathbf{X}^\top \beta\|_{\ell_2} \leq h}{\operatorname{argmin}} \{ \|\beta\|_{\ell_0} \}$$

where we might convexify the ℓ_0 norm, $\|\cdot\|_{\ell_0}$.

Regularization ℓ_0 , ℓ_1 et ℓ_2



Going further on sparsity issues

On $[-1, +1]^k$, the convex hull of $\|\beta\|_{\ell_0}$ is $\|\beta\|_{\ell_1}$

On $[-a, +a]^k$, the convex hull of $\|\beta\|_{\ell_0}$ is $a^{-1}\|\beta\|_{\ell_1}$

Hence,

$$\hat{\beta} = \underset{\beta; \|\beta\|_{\ell_1} \leq \tilde{s}}{\operatorname{argmin}} \{ \|\mathbf{Y} - \mathbf{X}^T \beta\|_{\ell_2} \}$$

is equivalent (Kuhn-Tucker theorem) to the Lagrangian optimization problem

$$\hat{\beta} = \operatorname{argmin} \{ \|\mathbf{Y} - \mathbf{X}^T \beta\|_{\ell_2} + \lambda \|\beta\|_{\ell_1} \}$$

LASSO *Least Absolute Shrinkage and Selection Operator*

$$\hat{\beta} \in \operatorname{argmin}\{\|\mathbf{Y} - \mathbf{X}^T \beta\|_{\ell_2} + \lambda \|\beta\|_{\ell_1}\}$$

is a convex problem (several algorithms^{*}), but not strictly convex (no unicity of the minimum). Nevertheless, predictions $\hat{\mathbf{y}} = \mathbf{x}^T \hat{\beta}$ are unique

^{*} MM, minimize majorization, coordinate descent [Hunter \(2003\)](#).

Optimal LASSO Penalty

Use cross validation, e.g. K -fold,

$$\hat{\beta}_{(-k)}(\lambda) = \operatorname{argmin} \left\{ \sum_{i \notin \mathcal{I}_k} [y_i - \mathbf{x}_i^\top \beta]^2 + \lambda \|\beta\| \right\}$$

then compute the sum of the squared errors,

$$Q_k(\lambda) = \sum_{i \in \mathcal{I}_k} [y_i - \mathbf{x}_i^\top \hat{\beta}_{(-k)}(\lambda)]^2$$

and finally solve

$$\lambda^* = \operatorname{argmin} \left\{ \overline{Q}(\lambda) = \frac{1}{K} \sum_k Q_k(\lambda) \right\}$$

Note that this might overfit, so [Hastie, Tibshiriani & Friedman \(2009\)](#) suggest the largest λ such that

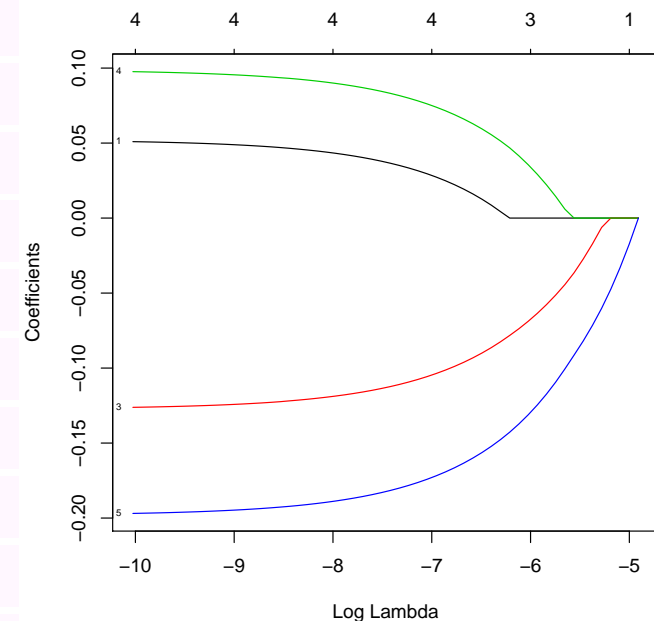
$$\overline{Q}(\lambda) \leq \overline{Q}(\lambda^*) + \operatorname{se}[\lambda^*] \quad \text{with} \quad \operatorname{se}[\lambda]^2 = \frac{1}{K^2} \sum_{k=1}^K [Q_k(\lambda) - \overline{Q}(\lambda)]^2$$

```

1 > freq = merge(contrat, nombre_RC)
2 > freq = merge(freq, nombre_D0)
3 > freq[,10]=as.factor(freq[,10])
4 > mx=cbind(freq[,c(4,5,6)],freq[,9]=="D",
             freq[,3]%in%c("A","B","C"))
5 > colnames(mx)=c(names(freq)[c(4,5,6)],"
                  diesel","zone")
6 > for(i in 1:ncol(mx)) mx[,i]=(mx[,i]-mean(
             mx[,i]))/sd(mx[,i])
7 > names(mx)
8 [1]    puissance agevehicule ageconducteur
          diesel          zone
9 > library(glmnet)
10 > fit = glmnet(x=as.matrix(mx), y=freq[,11],
                offset=log(freq[,2]), family = "poisson")
11 > plot(fit, xvar="lambda", label=TRUE)

```

LASSO, Fréquence RC



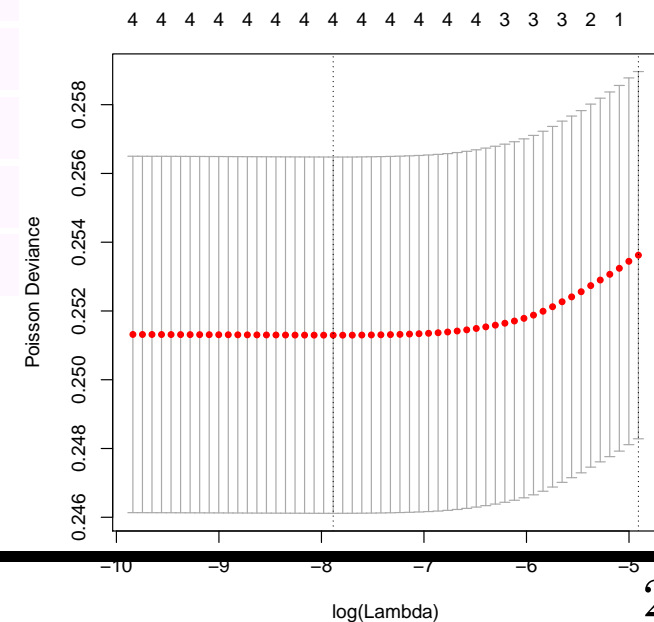
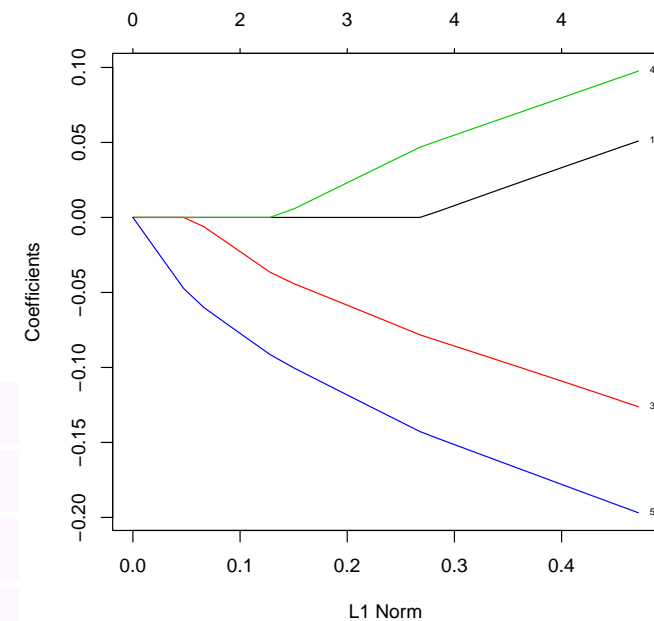
LASSO, Fréquence RC

```

1 > plot(fit,label=TRUE)
2 > cvfit = cv.glmnet(x=as.matrix(mx), y=freq
    [,11], offset=log(freq[,2]),family = "
    poisson")
3 > plot(cvfit)
4 > cvfit$lambda.min
5 [1] 0.0002845703
6 > log(cvfit$lambda.min)
7 [1] -8.16453

```

- Cross validation curve + error bars

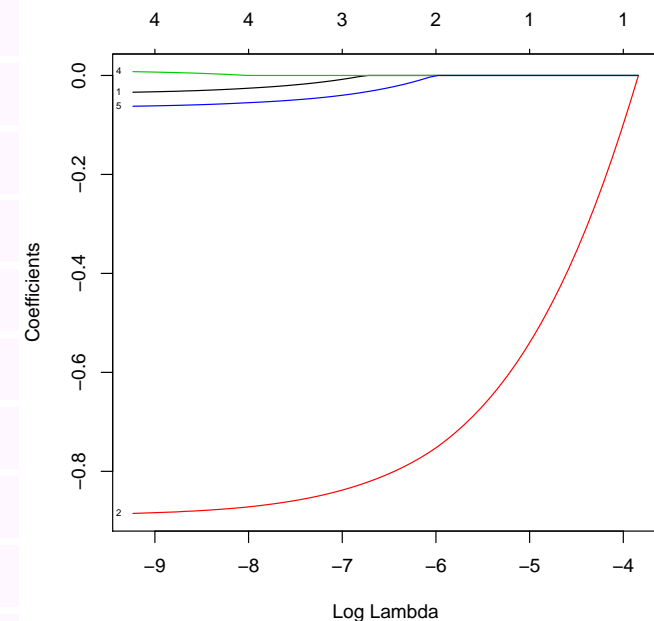


```

1 > freq = merge(contrat, nombre_RC)
2 > freq = merge(freq, nombre_DO)
3 > freq[,10]=as.factor(freq[,10])
4 > mx=cbind(freq[,c(4,5,6)],freq[,9]=="D",
             freq[,3]%in%c("A","B","C"))
5 > colnames(mx)=c(names(freq)[c(4,5,6)],"
                  diesel","zone")
6 > for(i in 1:ncol(mx)) mx[,i]=(mx[,i]-mean(
             mx[,i]))/sd(mx[,i])
7 > names(mx)
8 [1]    puissance agevehicule ageconducteur
          diesel          zone
9 > library(glmnet)
10 > fit = glmnet(x=as.matrix(mx), y=freq[,12],
                offset=log(freq[,2]), family = "poisson")
11 > plot(fit, xvar="lambda", label=TRUE)

```

LASSO, Fréquence DO



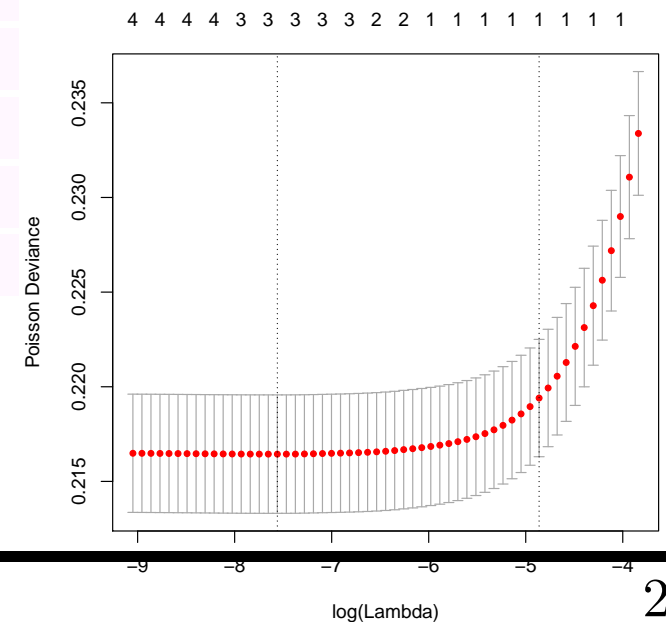
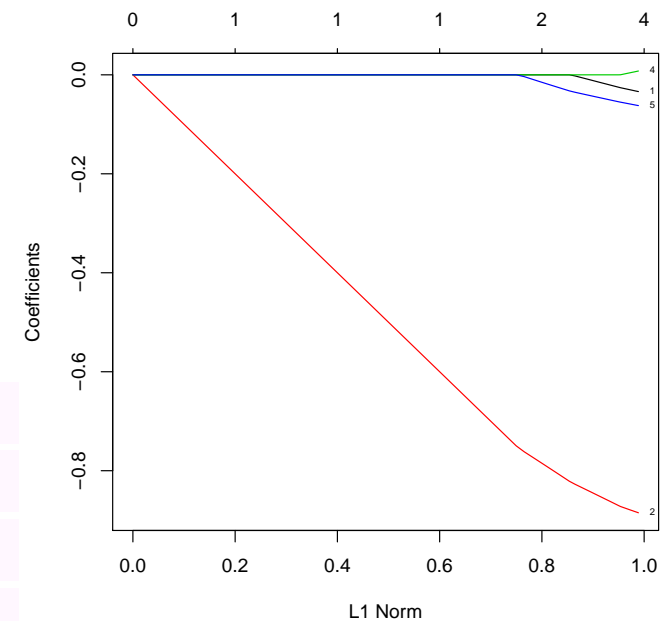
LASSO, Fréquence DO

```

1 > plot(fit,label=TRUE)
2 > cvfit = cv.glmnet(x=as.matrix(mx), y=freq
    [,12], offset=log(freq[,2]),family = "
    poisson")
3 > plot(cvfit)
4 > cvfit$lambda.min
5 [1] 0.0004744917
6 > log(cvfit$lambda.min)
7 [1] -7.653266

```

- Cross validation curve + error bars



Model Selection and Gini/Lorentz (on incomes)

Consider an ordered sample $\{y_1, \dots, y_n\}$, then Lorenz curve is

$$\{F_i, L_i\} \text{ with } F_i = \frac{i}{n} \text{ and } L_i = \frac{\sum_{j=1}^i y_j}{\sum_{j=1}^n y_j}$$

The theoretical curve, given a distribution F , is

$$u \mapsto L(u) = \frac{\int_{-\infty}^{F^{-1}(u)} t dF(t)}{\int_{-\infty}^{+\infty} t dF(t)}$$

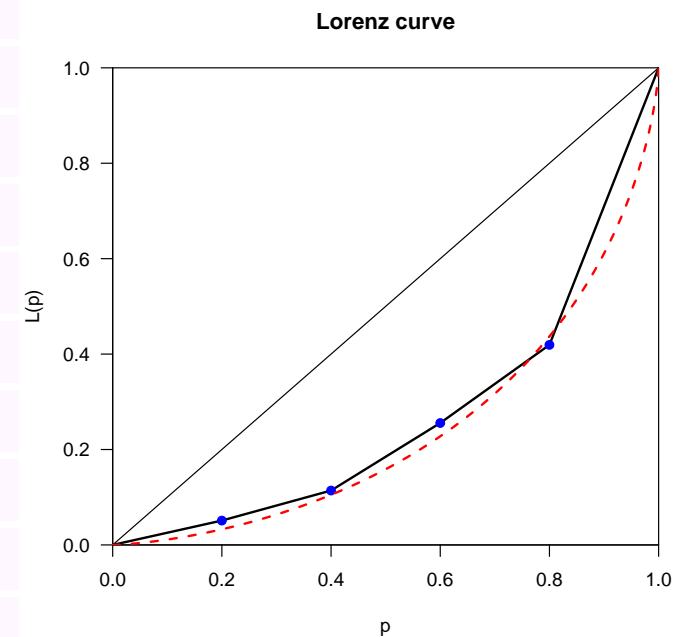
see Gastwirth (1972, econpapers.repec.org)

Model Selection and Gini/Lorentz (on incomes)

```

1 > library(ineq)
2 > set.seed(1)
3 > (x<-sort(rlnorm(5,0,1)))
4 [1] 0.4336018 0.5344838 1.2015872 1.3902836
    4.9297132
5 > Lc.sim <- Lc(x)
6 > plot(Lc.sim)
7 > points((1:4)/5,(cumsum(x)/sum(x))[1:4],pch
    =19,col="blue")
8 > lines(Lc.lognorm, parameter=1,lty=2)

```



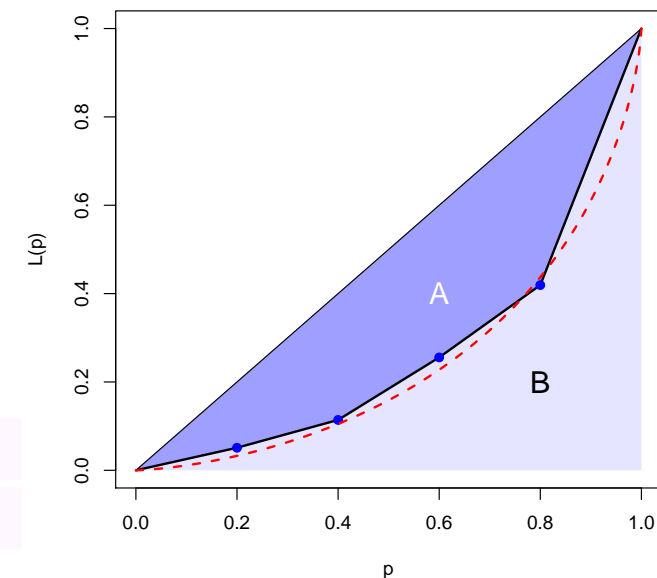
Model Selection and Gini/Lorentz (on incomes)

Gini index is the ratio of the areas $\frac{A}{A+B}$. Thus,

$$G = \frac{2}{n(n-1)\bar{x}} \sum_{i=1}^n i \cdot x_{i:n} - \frac{n+1}{n-1}$$

$$= \frac{1}{\mathbb{E}(Y)} \int_0^\infty F(y)(1 - F(y))dy$$

```
1 > Gini(x)
2 [1] 0.4640003
```



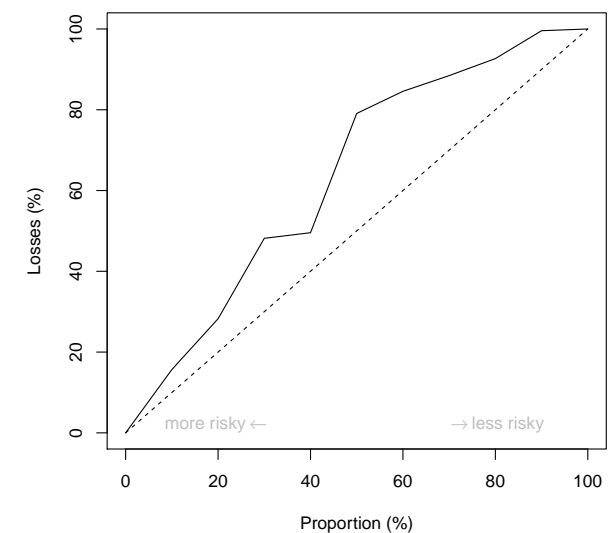
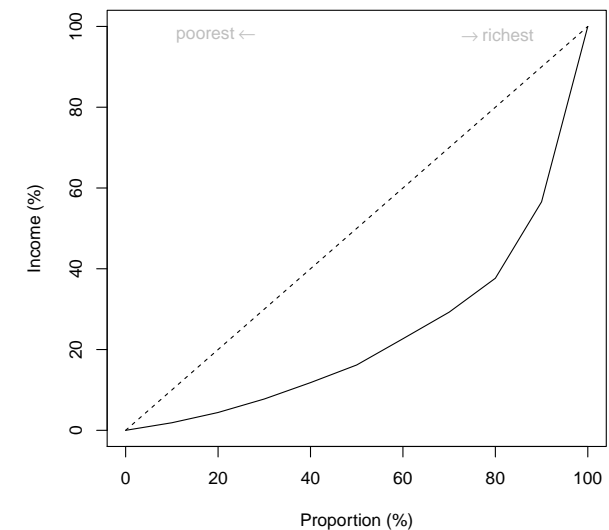
Comparing Models

Consider an ordered sample $\{y_1, \dots, y_n\}$ of incomes, with $y_1 \leq y_2 \leq \dots \leq y_n$, then Lorenz curve is

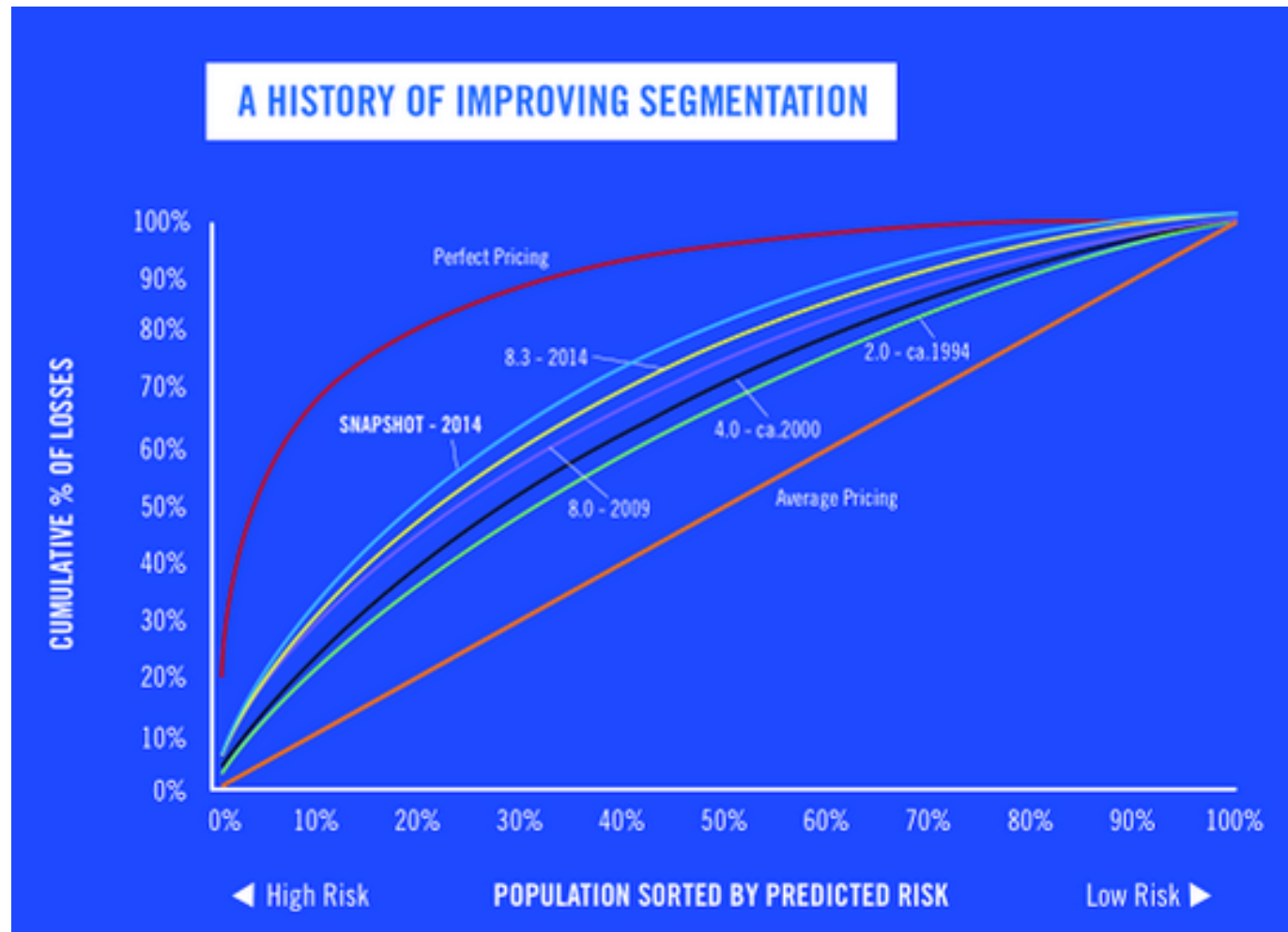
$$\{F_i, L_i\} \text{ with } F_i = \frac{i}{n} \text{ and } L_i = \frac{\sum_{j=1}^i y_j}{\sum_{j=1}^n y_j}$$

We have observed losses y_i and premiums $\hat{\pi}(x_i)$. Consider an **ordered sample by the model**, see **Frees, Meyers & Cummins (2014)**, $\hat{\pi}(x_1) \geq \hat{\pi}(x_2) \geq \dots \geq \hat{\pi}(x_n)$, then plot

$$\{F_i, L_i\} \text{ with } F_i = \frac{i}{n} \text{ and } L_i = \frac{\sum_{j=1}^i y_j}{\sum_{j=1}^n y_j}$$



Choix et comparaison de modèle en tarification



See Frees *et al.* (2010) or Tevet (2013).