

Modèles Linéaires Appliqués

Arthur Charpentier

Automne 2Q20

GLM #23 (example)

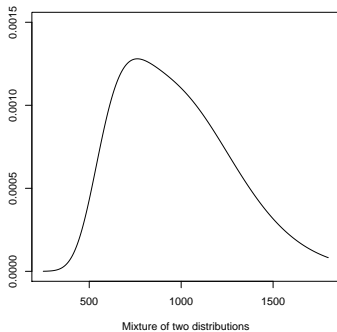
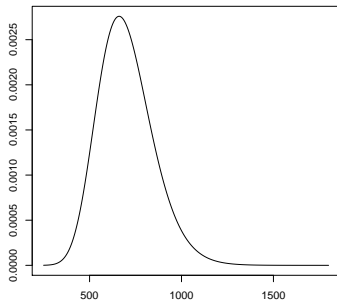
Fréquence & coût en assurance

```
1 > sinistre=read.table("http://freakonometrics.free.fr/sinistreACT2040.txt",
  header=TRUE, sep=";")
2 > contrat=read.table("http://freakonometrics.free.fr/contractACT2040.txt", header
  =TRUE, sep=";")
3 > contrat=contrat[,1:10]
4 > names(contrat)[10]="region"
5 > sinistre_D0=sinistre[(sinistre$garantie=="2D0")&(sinistre$cout>0),]
6 > sinistre_RC=sinistre[(sinistre$garantie=="1RC")&(sinistre$cout>0),]
7 > base_D0=merge(sinistre_D0, contrat)
8 > dim(base_D0)
9 [1] 1735    13
10 > base_RC=merge(sinistre_RC, contrat)
11 > dim(base_RC)
12 [1] 1924    13
```

Idée de base: si $S = \sum_{i=1}^N Y_i$,

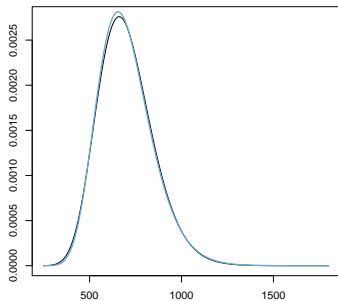
$$\pi(\mathbf{x}) = \mathbb{E}[S | \mathbf{X} = \mathbf{x}] = \mathbb{E}[N | \mathbf{X} = \mathbf{x}] \cdot \mathbb{E}[Y | \mathbf{X} = \mathbf{x}]$$

Coût : log-normale ou Gamma ?

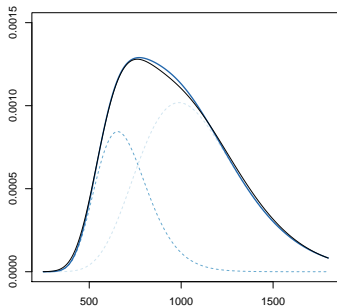


Coût : log-normale ou Gamma?

Loi Gamma ? Mélange de deux lois Gamma ?



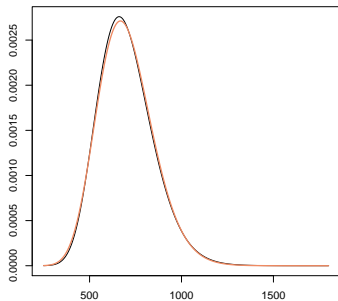
Lognormal distribution



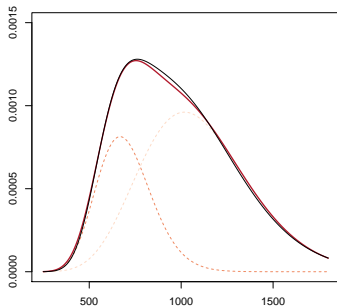
Mixture of lognormal distributions

Coût : log-normale ou Gamma?

Loi log-normale ? Mélange de deux lois log-normales ?



Gamma distribution



Mixture of Gamma distributions

Coût : Gamma

Pour la régression **Gamma** (et un lien **log** i.e.

$\mathbb{E}(Y|\mathbf{X} = \mathbf{x}) = \exp[\mathbf{x}^\top \boldsymbol{\beta}]$), on a

```
1 > regg=glm(cout~agevehicule+carburant+zone,data=base_RC,
2 +         family=Gamma(link="log"))
3 > summary(regg)
4 (Intercept)    7.72660      0.09300   83.079   < 2e-16 ***
5 agevehicule   -0.04674      0.00855   -5.466  5.27e-08 ***
6 carburantE     -0.14693      0.06329   -2.321  0.02038 *
7 zoneB         -0.14876      0.12690   -1.172  0.24124
8 zoneC         -0.04275      0.09924   -0.431  0.66668
9 zoneD         -0.11026      0.10416   -1.058  0.28998
10 zoneE         -0.12129      0.10478   -1.158  0.24719
11 zoneF         -0.47684      0.18142   -2.628  0.00865 **
12
13 (Dispersion parameter for Gamma family taken to be 1.686782)
```

Coût : Inverse-Gaussienne

Pour la régression **inverse-Gaussienne**, (et un lien **log** i.e.
 $\mathbb{E}(Y|\mathbf{X}) = \exp[\mathbf{X}'\beta]$),

```
1 > regig=glm(cout~agevehicule+carburant+zone,data=base_D0,
2 + family=inverse.gaussian(link="log"),start=coefficients(regg))
3 > summary(regig)
4 Coefficients:
5 (Intercept)  7.731661    0.093390   82.789   < 2e-16 ***
6 agevehicule -0.046699    0.007016   -6.656  3.76e-11 ***
7 carburantE   -0.153028    0.061479   -2.489   0.01290 *
8 zoneB       -0.138902    0.123192   -1.128   0.25968
9 zoneC       -0.054040    0.098951   -0.546   0.58505
10 zoneD       -0.102103    0.102734   -0.994   0.32043
11 zoneE       -0.127266    0.103662   -1.228   0.21973
12 zoneF       -0.492622    0.155715   -3.164   0.00159 **
13
14 (Dispersion parameter for inverse.gaussian family taken to be 0.001024064)
```

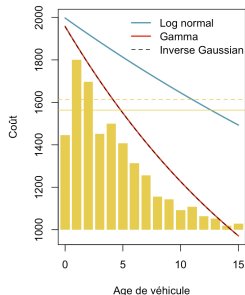
Coût : Log-normale

Pour la régression **log-normale** i.e. $\mathbb{E}(\log Y|\mathbf{X}) = \mathbf{X}'\beta$, on a

```
1 > regln=lm(log(cout)~agevehicule+carburant+zone,data=base_D0)
2 > summary(regln)
3 Coefficients:
4             Estimate Std. Error t value Pr(>|t|)
5 (Intercept)  6.776664   0.094371  71.809  <2e-16 ***
6 agevehicule -0.019397   0.008676  -2.236   0.0255 *
7 carburantE   -0.045508   0.064224  -0.709   0.4787
8 zoneB       -0.022196   0.128763  -0.172   0.8632
9 zoneC        0.056457   0.100695   0.561   0.5751
10 zoneD       -0.008894   0.105694  -0.084   0.9330
11 zoneE        0.017727   0.106321   0.167   0.8676
12 zoneF       -0.363002   0.184087  -1.972   0.0488 *
13
14 Residual standard error: 1.318 on 1727 degrees of freedom
15 > sigma=summary(regln)$sigma
```

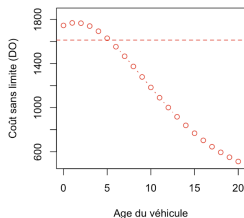
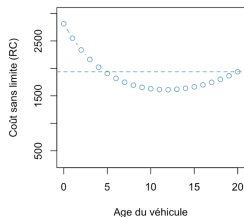

Coût : Comparison

```
1 > nd = data.frame(agevehicule=seq(0,15,by=.25),  
  carburant="E",zone="A")  
2 > Tb = table(base_D0$agevehicule)  
3 > ypln = exp(predict(regln,newdata=nd)+.5*sigma^2)  
4 > ypg = predict(regg,newdata=nd,type="response")  
5 > ypig = predict(regig,newdata=nd,type="response")  
6 > plot(nd$agevehicule,ypln,col="blue")  
7 > lines(nd$agevehicule,ypg,lwd=2,col=colr[2])  
8 > lines(nd$agevehicule,ypig,lty=2)  
9 > abline(h=mean(base_D0$cout),lty=2)  
10 abline(h=mean(base_D0$cout[(base_D0$carburant=="E")  
  &(base_D0$zone=="A")]))
```



Coût : RC vs DO

```
1 > library(bs)
2 > reg=glm(cout~bs(agevehicule),data=base_RC,family=
  Gamma(link="log"))
3 > age = 0:15
4 > yp=predict(reg,newdata=data.frame(agevehicule=age)
  )
5 > plot(age,yp,type="b")
6 > reg=glm(cout~bs(agevehicule),data=base_DO,family=
  Gamma(link="log"))
7 > age = 0:15
8 > yp=predict(reg,newdata=data.frame(agevehicule=age)
  )
9 > plot(age,yp,type="b")
```



Coût : coûts importants

On a ici quelques *gros* sinistres. L'idée est de noter que

$$\mathbb{E}(Y) = \sum_i \mathbb{E}(Y|\Theta = \theta_i) \cdot \mathbb{P}(\Theta = \theta_i)$$

Supposons que Θ prenne deux valeurs, correspondant au cas $\{Y \leq s\}$ et $\{Y > s\}$. Alors

$$\mathbb{E}(Y) = \mathbb{E}(Y|Y \leq s) \cdot \mathbb{P}(Y \leq s) + \mathbb{E}(Y|Y > s) \cdot \mathbb{P}(Y > s)$$

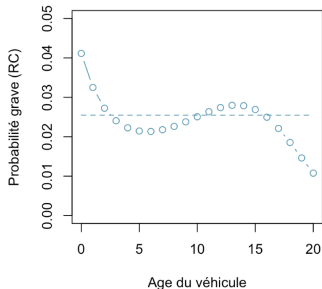
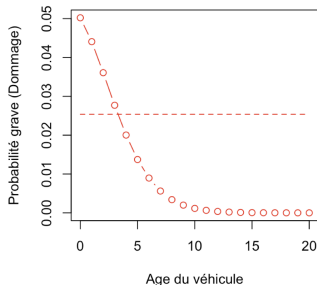
ou, en calculant l'espérance sous $\mathbb{P}_{\mathbf{X}}$ et plus \mathbb{P} ,

$$\mathbb{E}(Y|\mathbf{X}) = \underbrace{\mathbb{E}(Y|\mathbf{X}, Y \leq s)}_A \cdot \underbrace{\mathbb{P}(Y \leq s|\mathbf{X})}_B + \underbrace{\mathbb{E}(Y|Y > s, \mathbf{X})}_C \cdot \underbrace{\mathbb{P}(Y > s|\mathbf{X})}_{B^*}$$

Probabilité d'un grave

Pour le terme B , il s'agit d'une régression *standard* d'une variable de Bernoulli,

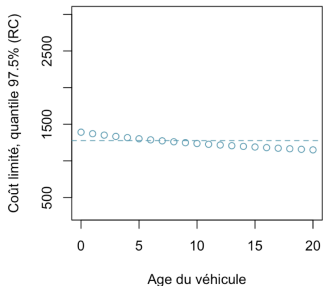
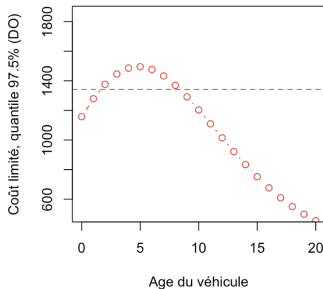
```
1 > sRC = quantile(base_RC$cout, .975)
2 > sRC
3 97.5%
4 8677.38
5 > sD0 = quantile(base_D0$cout, .975)
6 > sD0
7 97.5%
8 8203.183
9 > base_RC$anormal=(base_RC$cout >= sRC)
10 > library(splines)
11 > age=seq(0,20)
12 > regB=glm(normal~bs(agevehicule), data=base_RC,
13           family=binomial)
14 > ypB=predict(regB, newdata=data.frame(agevehicule=
15           age), type="response")
16 > plot(age, ypB, type="b")
```



Coût d'un non-grave

Pour le terme A , il s'agit d'une régression *standard* Gamma,

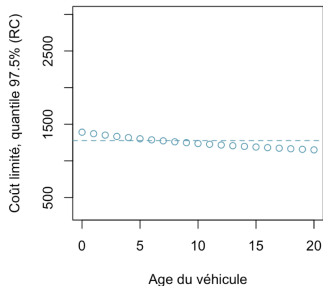
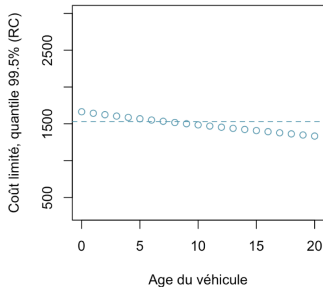
```
1 > library(splines)
2 > indice = which(base_RC$cout <= sRC)
3 > mean(base_RC$cout[indice])
4 [1] 1269.22
5 > regA = glm(cout~bs(agevehicule),data=base_RC,
6             subset=indice,family=Gamma(link="log"))
7 > ypA = predict(regA,newdata=data.frame(agevehicule=
8             age),type="response")
8 > plot(age,ypA,type="b")
```



Coût d'un non-grave

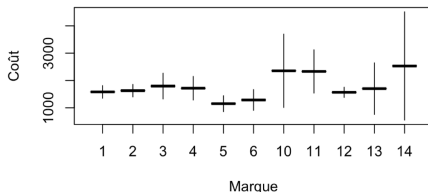
On peut prendre un seuil plus élevé

```
1 > sRC = quantile(base_RC$cout, .995)
2 > sRC
3   99.5%
4 23978.43
5 > library(splines)
6 > indice = which(base_RC$cout <= sRC)
7 > mean(base_RC$cout[indice])
8 [1] 1529.892
9 > regA = glm(cout~bs(agevehicule), data=base_RC,
10             subset=indice, family=Gamma(link="log"))
11 > ypA = predict(regA, newdata=data.frame(agevehicule=
12             age), type="response")
12 > plot(age, ypA, type="b")
```



Modéliser les coûts individuels

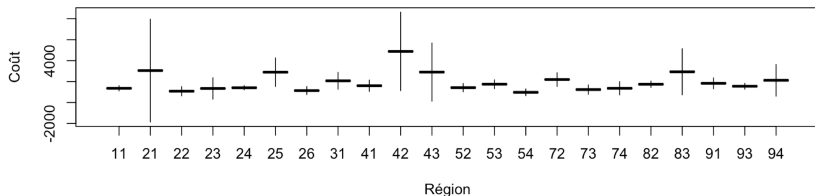
```
1 k = which(names(base_DO) %in% c("garantie","no","nocontrat","exposition"))
2 base_DO = base_DO[,-k]
3 base_DO$marque = as.factor(base_DO$marque)
4 A= aggregate(base_DO$cout, by=list(base_DO$marque),function(x) c(mean(x)-2*sd(x)
  /sqrt(length(x)),mean(x),mean(x)+2*sd(x)/sqrt(length(x))))
5 plot(A$Group.1,A$x[,2],ylim=range(A$x))
6 segments(1:nrow(A),A$x[,1],1:nrow(A),A$x[,3])
```



```
1 levels(base_DO$marque) = c("A","A","A","A","B","B","C","C","A","A","C")
2 levels(base_DO$zone) = c("A","A","A","A","A","A","F")
```

Modéliser les coûts individuels

```
1 base_DO$region = as.factor(base_DO$region)
2 A= aggregate(base_DO$cout, by=list(base_DO$region),function(x) c(mean(x)-2*sd(x)
  /sqrt(length(x)),mean(x),mean(x)+2*sd(x)/sqrt(length(x))))
3 plot(A$Group.1,A$x[,2],ylim=range(A$x))
4 segments(1:nrow(A),A$x[,1],1:nrow(A),A$x[,3])
```



```
1 levels(base_DO$region) = c("B","E","A","B", "B","E","A","D", "B","F","E","B","C",
  ",","A","D","A", "B","C","E","C", "B","D")
```


Modéliser les coûts individuels

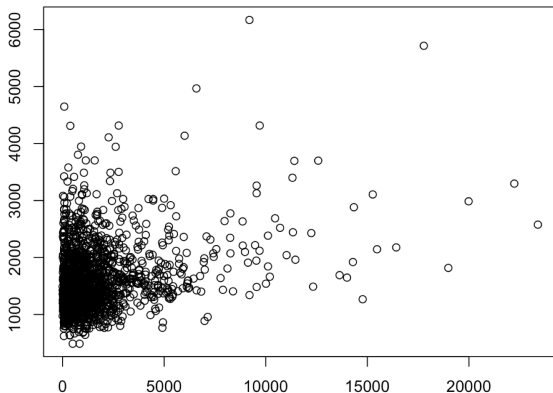
```
1 > regg=glm(cout~.,data=base_D0, family=Gamma(link="log"))
2 > summary(regg)
3
4 Coefficients:
5             Estimate Std. Error t value Pr(>|t|)
6 (Intercept)   7.159940   0.229208  31.238 < 2e-16 ***
7 zoneF        -0.331000   0.156797  -2.111 0.034915 *
8 puissance     0.035415   0.014598   2.426 0.015369 *
9 agevehicule  -0.033725   0.007937  -4.249 2.26e-05 ***
10 ageconducteur -0.006137   0.002309  -2.658 0.007933 **
11 bonus         0.006451   0.002141   3.013 0.002626 **
12 marqueB      -0.258137   0.103477  -2.495 0.012702 *
13 marqueC       0.309273   0.141523   2.185 0.029000 *
14 carburantE    -0.096018   0.059037  -1.626 0.104046
15 regionE       0.565601   0.165092   3.426 0.000627 ***
16 regionA      -0.287633   0.111578  -2.578 0.010024 *
17 regionD       0.296621   0.101872   2.912 0.003641 **
18 regionF       0.860987   0.603632   1.426 0.153951
19 regionC       0.126424   0.071222   1.775 0.076064 .
20 ---
21 Signif. codes:
22 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
23
24 (Dispersion parameter for Gamma family taken to be 1.442703)
25
26 Null deviance: 2464.2 on 1734 degrees of freedom
27 Residual deviance: 2259.6 on 1721 degrees of freedom
28 AIC: 28937
29
30 Number of Fisher Scoring iterations: 7
```

Modéliser les coûts individuels

```
1 > regln=lm(log(cout)~.,data=base_D0)
2 > summary(regln)
3
4 Coefficients:
5             Estimate Std. Error t value Pr(>|t|)
6 (Intercept)   6.396580   0.248328  25.759 < 2e-16 ***
7 zoneF        -0.369535   0.169876  -2.175  0.02974 *
8 puissance     0.009819   0.015816   0.621  0.53478
9 agevehicule  -0.015695   0.008599  -1.825  0.06813 .
10 ageconducteur -0.002956   0.002502  -1.182  0.23748
11 bonus        0.006678   0.002320   2.878  0.00405 **
12 marqueB      -0.226641   0.112109  -2.022  0.04337 *
13 marqueC       0.334085   0.153328   2.179  0.02948 *
14 carburantE    -0.022320   0.063961  -0.349  0.72716
15 regionE       0.408499   0.178863   2.284  0.02250 *
16 regionA      -0.232556   0.120885  -1.924  0.05455 .
17 regionD       0.190320   0.110369   1.724  0.08482 .
18 regionF       1.314463   0.653985   2.010  0.04459 *
19 regionC       0.157108   0.077163   2.036  0.04190 *
20 ---
21 Signif. codes:
22 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
23
24 Residual standard error: 1.301 on 1721 degrees of freedom
25 Multiple R-squared:  0.03465, Adjusted R-squared:  0.02736
26 F-statistic: 4.751 on 13 and 1721 DF,  p-value: 3.657e-08
```

Modéliser les coûts individuels

```
1 > plot(base_D0$cout, predict(regg, type="response"))
```



pour aller plus loin ? [ACT6100](#)