Let us get back on the Titanic dataset,

```
1 loc_fichier = "http://freakonometrics.free.fr/titanic.RData"
2 download.file(loc_fichier, "titanic.RData")
3 load("titanic.RData")
4 base = base[!is.na(base$Age),]
```

On consider two variables, the age x (the continuous one) and the survivor indicator y (the qualitative one)

```
1 X = base$Age
2 Y = base$Survived
```

It looks like the age might be a valid explanatory variable in the logistic regression,

```
1  summary(glm(Survived~Age,data=base,family=binomial))
2
3  Coefficients:
4              Estimate Std. Error z value Pr(&gt;|z|)
5  (Intercept) -0.05672    0.17358  -0.327    0.7438
6  Age         -0.01096    0.00533  -2.057    0.0397 *
7  ---
8  Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
9
10 (Dispersion parameter for binomial family taken to be 1)
11
12     Null deviance: 964.52  on 713  degrees of freedom
13 Residual deviance: 960.23  on 712  degrees of freedom
14 AIC: 964.23
```

The significance test here has a p-value just below 4%. Actually, one can relate it with the value of the deviance (the null deviance and the residual deviance). Recall that $D=2\big(\log\mathcal{L}(\boldsymbol{y})-\log\mathcal{L}(\widehat{\boldsymbol{\mu}})\big)$ while $D_0=2\big(\log\mathcal{L}(\boldsymbol{y})-\log\mathcal{L}(\overline{y})\big)$ Under the assumption that x is worthless, $D_0-D$ tends to a $\chi^2$ distribution with 1 degree of freedom. And we can compute the p-value dof that likelihood ratio test,

```
1 1-pchisq(964.52-960.23,1)
2 [1] 0.03833717
```

(which is consistent with a Gaussian test). But if we consider a nonlinear transformation

```
1  summary(glm(Survived~bs(Age),data=base,family=binomial))
2
3  Coefficients:
4              Estimate Std. Error z value Pr(&gt;|z|)
5  (Intercept)   0.8648     0.3460   2.500 0.012433 *
6  bs(Age)1     -3.6772     1.0458  -3.516 0.000438 ***
7  bs(Age)2      1.7430     1.1068   1.575 0.115299
8  bs(Age)3     -3.9251     1.4544  -2.699 0.006961 **
9  ---
10 Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
11
```

```
12 (Dispersion parameter for binomial family taken to be 1)
13
14     Null deviance: 964.52  on 713  degrees of freedom
15 Residual deviance: 948.69  on 710  degrees of freedom
```

which seems to be "more significant"

```
1 1-pchisq(964.52-948.69,3)
2 [1] 0.001228712
```

So it looks like the variable x is interesting here.

To visualize the non-null correlation, one can consider the condition distribution of x given y=1, and compare it with the condition distribution of x given y=0,

```
1 ks.test(X[Y==0],X[Y==1])
2
3          Two-sample Kolmogorov-Smirnov test
4
5 data:  X[Y == 0] and X[Y == 1]
6 D = 0.088777, p-value = 0.1324
7 alternative hypothesis: two-sided
```
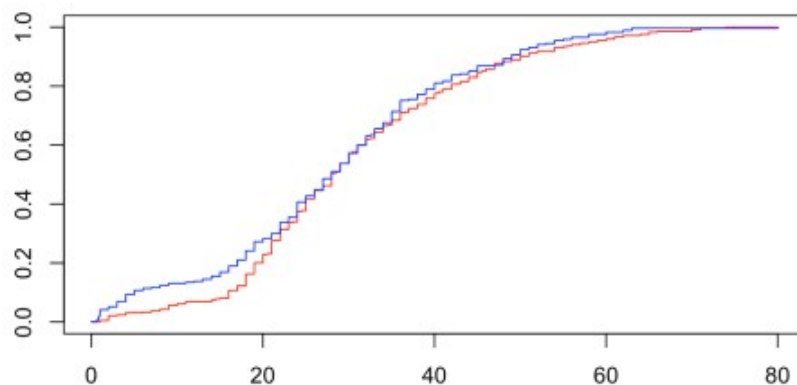
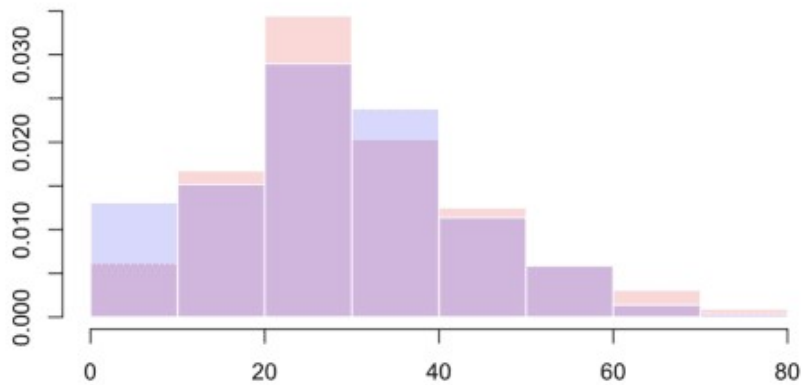i.e. with a p-value above 10%, the two distributions are not significatly different.

```
1 F0 = function(x) mean(X[Y==0]&lt;=x)
2 F1 = function(x) mean(X[Y==1]&lt;=x)
3 vx = seq(0,80,by=.1)
4 vy0 = Vectorize(F0)(vx)
5 vy1 = Vectorize(F1)(vx)
6 plot(vx,vy0,col="red",type="s")
7 lines(vx,vy1,col="blue",type="s")
```



(we can also look at the density, but it looks like that there is not much to see)

An alternative is discretize variable x and to use Pearson's independence test,

```
1  k=5
2  LV = quantile(X,(0:k)/k)
3  LV[1] = 0
4  Xc = cut(X,LV)
5  table(Xc,Y)
6               Y
7  Xc             0  1
8    (0,19]     85 79
9    (19,25]    92 45
10   (25,31.8] 77 50
11   (31.8,41] 81 63
12   (41,80]    89 53
13 chisq.test(table(Xc,Y))
14
15         Pearson's Chi-squared test
16
17 data:  table(Xc, Y)
18 X-squared = 8.6155, df = 4, p-value = 0.07146
```
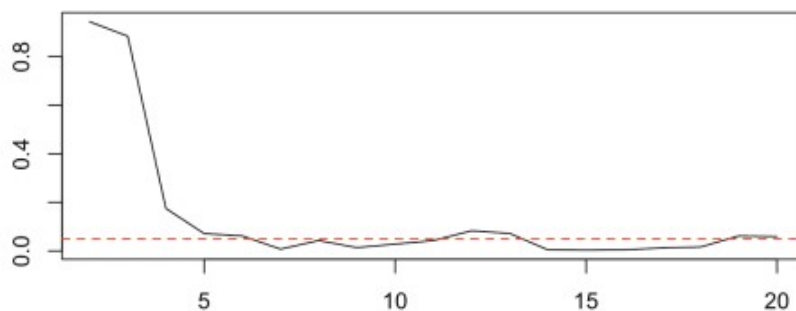
The p-value is here 7%, with five categories for the age. And actually, we can compare the p-value

```
1 pvalue = function(k=5){
2 LV = quantile(X,(0:k)/k)
3 LV[1] = 0
4 Xc = cut(X,LV)
5 chisq.test(table(Xc,Y))$p.value}
6 vk = 2:20
7 vp = Vectorize(pvalue)(vk)
8 plot(vk,vp,type="l")
9 abline(h=.05,col="red",lty=2)
```

which gives a p-value close to 5%, as soon as we have enough categories….