# Median

Consider a sample $\{y_1,\cdots,y_n\}$. To compute the median, solve
$$\min_\mu \left\lbrace\sum_{i=1}^n|y_i-\mu|\right\rbrace$$
which can be solved using linear programming techniques. More precisely, this problem is equivalent to
$$\min_{\mu,\mathbf{a},\mathbf{b}}\left\lbrace\sum_{i=1}^n a_i+b_i\right\rbrace$$
with $a_i,b_i\geq 0$ and $y_i-\mu=a_i-b_i, \forall i=1,\cdots,n$. Heuristically, the idea is to write $y_i=\mu+\varepsilon_i$, and then define $a_i$'s and $b_i$'s so that $\varepsilon_i=a_i-b_i$ and $|\varepsilon_i|=a_i+b_i$, i.e. $a_i=(\varepsilon_i)_+=\max\lbrace0,\varepsilon_i\rbrace=|\varepsilon|\cdot\boldsymbol{1}_{\varepsilon_i>0}$ and $b_i=(-\varepsilon_i)_+=\max\lbrace0,-\varepsilon_i\rbrace=|\varepsilon|\cdot\boldsymbol{1}_{\varepsilon_i<0}$ denote respectively the positive and the negative parts.</p>
<p>Unfortunately (that was the error in my previous post), the expression of linear programs is
$$\min_{\mathbf{z}}\left\lbrace\boldsymbol{c}^\top\mathbf{z}\right\rbrace\text{ s.t. }\boldsymbol{A}\mathbf{z}=\boldsymbol{b},\mathbf{z}\geq\boldsymbol{0}$$
In the equation above, with the $a_i$'s and $b_i$'s, we're not far away. Except that we have $\mu\in\mathbb{R}$, while it should be positive. So similarly, set $\mu=\mu^+-\mu^-$ where $\mu^+=(\mu)_+$ and $\mu^-=(-\mu)_+$.

Thus, let
$$\mathbf{z}=\big(\mu^+;\mu^-;\boldsymbol{a},\boldsymbol{b}\big)^\top\in\mathbb{R}_+^{2n+2}$$
and then write the constraint as $\boldsymbol{A}\mathbf{z}=\boldsymbol{b}$ with $\boldsymbol{b}=\boldsymbol{y}$ and $\boldsymbol{A}=\big[\boldsymbol{1}_n;-\boldsymbol{1}_n;\mathbb{I}_n;-\mathbb{I}_n\big]$ And for the objective function $\boldsymbol{c}=\big(\boldsymbol{0},\boldsymbol{1}_n,-\boldsymbol{1}_n\big)^\top\in\mathbb{R}_+^{2n+2}$

To illustrate, consider a sample from a lognormal distribution,

```
1 n = 101
2 set.seed(1)
3 y = rlnorm(n)
4 median(y)
5 [1] 1.077415
```

For the optimization problem, use the matrix form, with 3n constraints, and 2n+1 parameters,

```
1 library(lpSolve)
2 X = rep(1,n)
3 A = cbind(X, -X, diag(n), -diag(n))
4 b = y
5 c = c(rep(0,2), rep(1,n),rep(1,n))
6 equal_type = rep("=", n)
7 r = lp("min", c,A,equal_type,b)
8 head(r$solution,1)
9 [1] 1.077415
```

It looks like it's working well…

# Quantile

Of course, we can adapt our previous code for quantiles

```
1 tau = .3
2 quantile(y,tau)
3       30%
4 0.6741586
```

The linear program is now$\min_{q^+,q^-,\mathbf{a},\mathbf{b}}\left\lbrace\sum_{i=1}^n\tau a_i+(1-\tau)b_i\right\rbrace$with $a_i,b_i,q^+,q^-\geq 0$ and $y_i=q^+-q^-+a_i-b_i, \forall i=1,\cdots,n$. The R code is now

```
1 c = c(rep(0,2), tau*rep(1,n),(1-tau)*rep(1,n))
2 r = lp("min", c,A,equal_type,b)
3 head(r$solution,1)
4 [1] 0.6741586
```

So far so good…

# Quantile Regression

Consider the following dataset, with rents of flat, in a major German city, as function of the surface, the year of construction, etc.

```
1 base=read.table("http://freakonometrics.free.fr/rent98_00.txt",header=TRUE)
```

The linear program for the quantile regression is now$\min_{\boldsymbol{\beta}^+,\boldsymbol{\beta}^-,\mathbf{a},\mathbf{b}}\left\lbrace\sum_{i=1}^n\tau a_i+(1-\tau)b_i\right\rbrace$with $a_i,b_i\geq 0$ and $y_i=\boldsymbol{x}^\top[\boldsymbol{\beta}^+-\boldsymbol{\beta}^-]+a_i-b_i\forall i=1,\cdots,n$ and $\beta_j^+,\beta_j^-\geq 0 \forall j=0,\cdots,k$. So use here

```
1  require(lpSolve)
2  tau = .3
3  n=nrow(base)
4  X = cbind( 1, base$area)
5  y = base$rent_euro
6  K = ncol(X)
7  N = nrow(X)
8  A = cbind(X,-X,diag(N),-diag(N))
9  c = c(rep(0,2*ncol(X)),tau*rep(1,N),(1-tau)*rep(1,N))
10 b = base$rent_euro
11 const_type = rep("=",N)
12 r = lp("min",c,A,const_type,b)
13 beta = r$sol[1:K] -  r$sol[(1:K+K)]
14 beta
15 [1] 148.946864   3.289674
```

Of course, we can use R function to fit that model

```
1 library(quantreg)
2 rq(rent_euro~area, tau=tau, data=base)
3 Coefficients:
4 (Intercept)         area
5  148.946864    3.289674
```
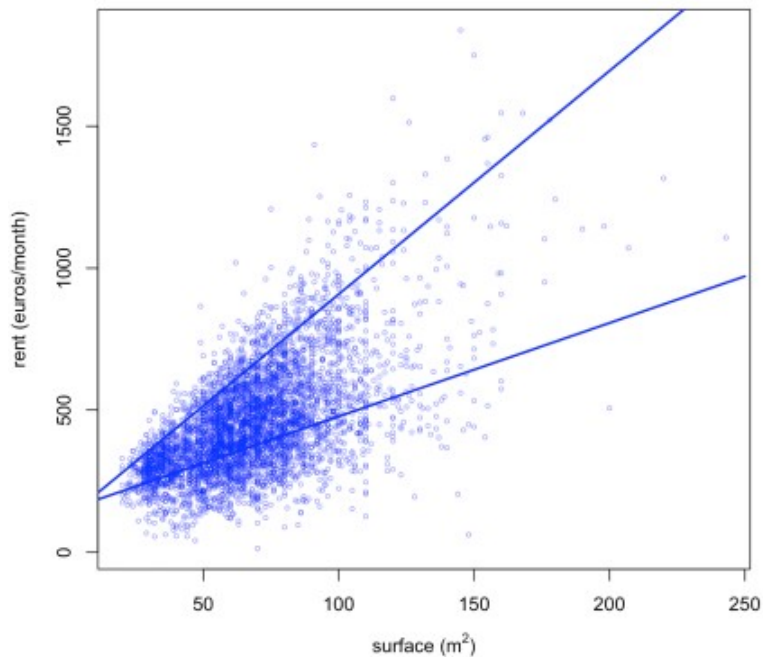
Here again, it seems to work quite well. We can use a different probability level, of course, and get a plot

```
1 plot(base$area,base$rent_euro,xlab=expression(paste("surface (",m^2,")")),
2      ylab="rent (euros/month)",col=rgb(0,0,1,.4),cex=.5)
3 sf=0:250
4 yr=r$solution[2*n+1]+r$solution[2*n+2]*sf
5 lines(sf,yr,lwd=2,col="blue")
```

```
6  tau = .9
7  r = lp("min",c,A,const_type,b)
8  tail(r$solution,2)
9  [1] 121.815505    7.865536
10 yr=r$solution[2*n+1]+r$solution[2*n+2]*sf
11 lines(sf,yr,lwd=2,col="blue")
```



And we can adapt the later to multiple regressions, of course,

```
1  X = cbind(1,base$area,base$yearc)
2  K = ncol(X)
3  N = nrow(X)
4  A = cbind(X,-X,diag(N),-diag(N))
5  c = c(rep(0,2*ncol(X)),tau*rep(1,N),(1-tau)*rep(1,N))
6  b = base$rent_euro
7  const_type = rep("=",N)
8  r = lp("min",c,A,const_type,b)
9  beta = r$sol[1:K] -  r$sol[(1:K+K)]
10 beta
11 [1] -5542.503252    3.978135    2.887234
```

to be compared with

```
1 library(quantreg)
2 rq(rent_euro~ area + yearc, tau=tau, data=base)
3
4 Coefficients:
5  (Intercept)          area          yearc
6 -5542.503252      3.978135      2.887234
7
8 Degrees of freedom: 4571 total; 4568 residual
```