

Modelos

IZASKUN LOPEZ-SAMANIEGO

19 de noviembre de 2017

Preparación del entorno

```
library(data.table)
library(lubridate)
```

```
##
## Attaching package: 'lubridate'
```

```
## The following objects are masked from 'package:data.table':
##
##     hour, isoweek, mday, minute, month, quarter, second, wday,
##     week, yday, year
```

```
## The following object is masked from 'package:base':
##
##     date
```

```
library(caret)
```

```
## Loading required package: ggplot2
```

```
## Warning: package 'ggplot2' was built under R version 3.3.3
```

```
setwd(ruta)
source('./src/definitivos/funciones_opendata.R')
```

Cargar ficheros

a. Cargamos el fichero con la información normalizada y nos quedamos con los datos necesarios para ejecutar el modelo.

```
dt.analisis <- as.data.table(read.csv('F:/201711_dataton_opendata_madrid/d
at/PM16_dataset.csv'))

dt.analisis <- dt.analisis[,list(identif, ds,
                                intensidad, ocupacion, carga,
                                vmed, vel.med, carga.med,
                                diaSemana, diaMes, Mes, fechaTrunc,
                                prec, prec_norm, prec_zscore,
                                Dia_semana, laborable...festivo...domingo.f
estivo,
                                Tipo.de.Festivo, Festividad)]
```

```
dt.datos.prev <- as.data.table(read.csv('F:/201711_dataton_opendata_madri
d/dat/trafico_outlier_datos_previos.csv', sep = ';'))
dt.datos.prev <- dt.datos.prev[,list(identif, ds,
                                     carga.1 = carga.30/100,
                                     vmed.1 = vmed.30/100,
                                     carga.2 = carga.45/100,
                                     vmed.2 = vmed.45/100,
                                     carga.3 = carga.60/100,
                                     vmed.3 = vmed.60/100,
                                     carga.4 = carga.75/100,
                                     vmed.4 = vmed.75/100)]
```

```
dt.analisis <- merge(dt.analisis, dt.datos.prev,
                     by.x = c('identif', 'ds'),
                     by.y = c('identif', 'ds'),
                     all.x = FALSE, all.y = FALSE)
dt.analisis <- Transformacion_variables(dt.analisis)
```

Dividimos la muestra en casos de test y casos de training

```
inTrain <- sample(1:nrow(dt.analisis),
                  nrow(dt.analisis)*0.3)

train.analisis <- dt.analisis[-inTrain,]
test.analisis <- dt.analisis[inTrain,]
```

Regresión Líneal Múltivariante

```
lm.M30 <- lm(carga ~ vel.med +
              carga.med +
              carga.1 +
              vmed.1 +
              carga.2 +
              vmed.2 +
              carga.3 +
              vmed.3 +
              diaMes +
              Mes +
              prec_norm +
              #   var.carga.1 +
              #   var.carga.2 +
              var.carga.3 +
              var.vmed.1 +
              var.vmed.2 +
              var.vmed.3 +
              diaLunes +
              diaMartes +
              diaMiercoles +
              diaJueves +
              diaViernes +
              diaSabado +
              #   diaDomingo +
              n.festivo ,
              data = train.analisis)
print(lm.M30$coefficients)
```

```
##   (Intercept)      vel.med      carga.med      carga.1      vmed.1
## -1.841954e-02 -3.628076e-03  8.050415e-02  1.147173e+00 -7.845003e-02
##      carga.2      vmed.2      carga.3      vmed.3      diaMes
## -2.508606e-02  3.069201e-02 -2.004331e-01  7.657098e-02 -7.500979e-05
##      Mes      prec_norm      var.carga.3      var.vmed.1      var.vmed.2
##  4.751994e-05  2.605157e-03 -2.633825e-01 -1.320161e-03 -6.185385e-04
##      var.vmed.3      diaLunes      diaMartes      diaMiercoles      diaJueves
##  1.197587e-02  2.344749e-03  1.924341e-03  1.831088e-03  2.147832e-03
##      diaViernes      diaSabado      n.festivo
##  2.177215e-03 -3.911486e-03 -4.968884e-03
```

```
summary(lm.M30)
```

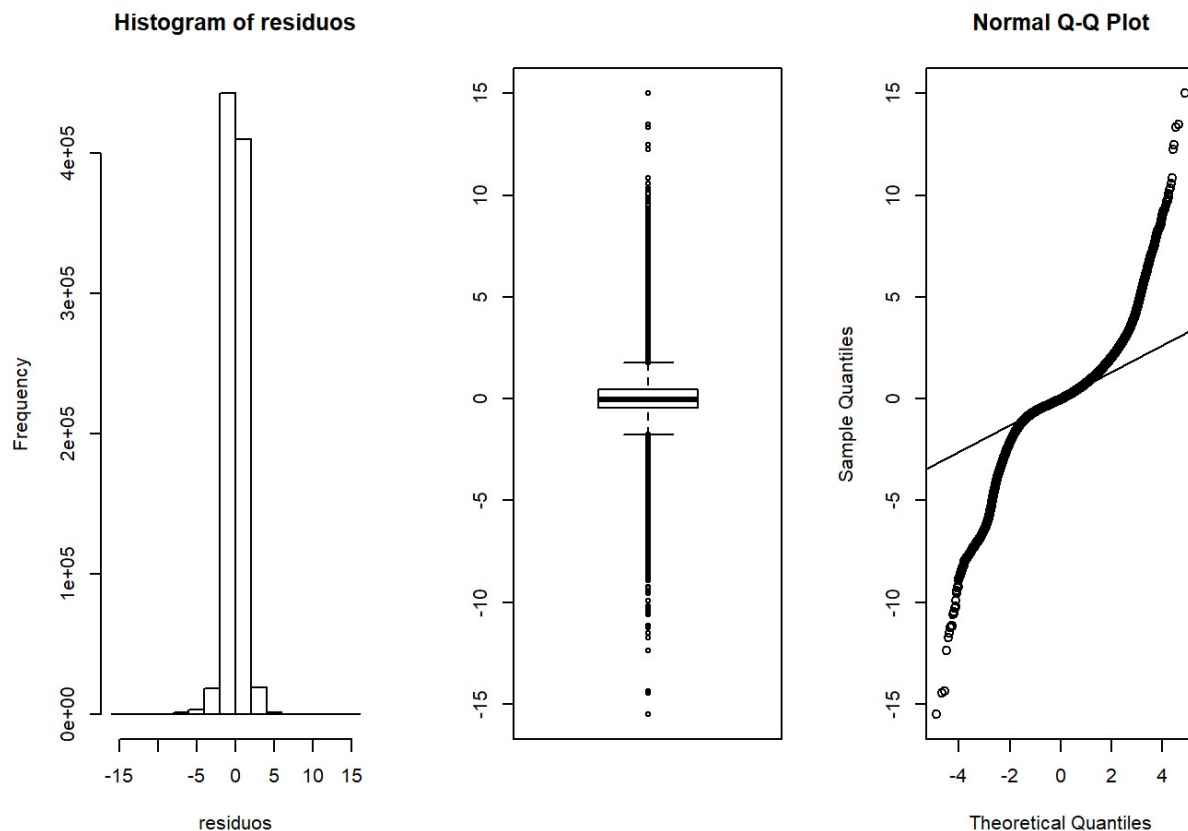
```
##
## Call:
## lm(formula = carga ~ vel.med + carga.med + carga.1 + vmed.1 +
##      carga.2 + vmed.2 + carga.3 + vmed.3 + diaMes + Mes + prec_norm +
##      var.carga.3 + var.vmed.1 + var.vmed.2 + var.vmed.3 + diaLunes +
##      diaMartes + diaMiercoles + diaJueves + diaViernes + diaSabado +
##      n.festivo, data = train.analisis)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -1.10258 -0.03027 -0.00221  0.03232  1.06615
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  -1.842e-02  6.725e-04  -27.389  < 2e-16 ***
## vel.med      -3.628e-03  1.242e-04  -29.220  < 2e-16 ***
## carga.med     8.050e-02  9.812e-04   82.043  < 2e-16 ***
## carga.1      1.147e+00  1.553e-03  738.482  < 2e-16 ***
## vmed.1       -7.845e-02  1.526e-03  -51.395  < 2e-16 ***
## carga.2      -2.509e-02  2.288e-03  -10.963  < 2e-16 ***
## vmed.2       3.069e-02  2.037e-03   15.069  < 2e-16 ***
## carga.3      -2.004e-01  1.545e-03 -129.691  < 2e-16 ***
## vmed.3       7.657e-02  1.624e-03   47.146  < 2e-16 ***
## diaMes       -7.501e-05  8.615e-06   -8.707  < 2e-16 ***
## Mes          4.752e-05  2.202e-05    2.158  0.03092 *
## prec_norm    2.605e-03  5.383e-04    4.840  1.30e-06 ***
## var.carga.3  -2.634e-01  1.540e-03 -170.974  < 2e-16 ***
## var.vmed.1   -1.320e-03  4.341e-04   -3.041  0.00236 **
## var.vmed.2   -6.185e-04  4.433e-04   -1.395  0.16289
## var.vmed.3    1.198e-02  3.282e-04   36.491  < 2e-16 ***
## diaLunes     2.345e-03  3.628e-04    6.463  1.02e-10 ***
## diaMartes    1.924e-03  3.627e-04    5.305  1.13e-07 ***
## diaMiercoles 1.831e-03  3.651e-04    5.015  5.31e-07 ***
## diaJueves    2.148e-03  3.658e-04    5.871  4.32e-09 ***
## diaViernes   2.177e-03  3.588e-04    6.069  1.29e-09 ***
## diaSabado    -3.911e-03  3.060e-04  -12.784  < 2e-16 ***
## n.festivo    -4.969e-03  1.285e-04  -38.677  < 2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.07111 on 897412 degrees of freedom
## Multiple R-squared:  0.9206, Adjusted R-squared:  0.9206
## F-statistic: 4.73e+05 on 22 and 897412 DF, p-value: < 2.2e-16
```

```
setwd(ruta)
saveRDS(lm.M30, './modelos/lmM30_30min.RData')
```

Análisis de los residuos

a. Supuesto 1: Normalidad

```
residuos<-rstandard(lm.M30) # residuos estándares del modelo ajustado (completo)
par(mfrow=c(1,3))
hist(residuos) # histograma de los residuos estandarizados
boxplot(residuos) # diagrama de cajas de los residuos estandarizados
qqnorm(residuos) # gráfico de cuantiles de los residuos estandarizados
qqline(residuos)
```

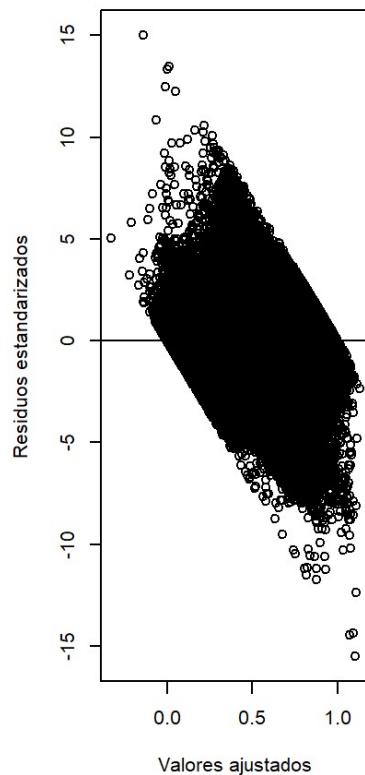


b. Supuesto 2: Varianza de los errores es constante:

- No es constante, tiene tendencia lo que indica que hay una variable desconocida que impacta en el tráfico pero no la hemos detectado.

```
par(mfrow=c(1,3))

# gráfico 2D de los valores ajustados vs. los residuos estandarizados
plot(fitted.values(lm.M30), rstandard(lm.M30),
     xlab="Valores ajustados",
     ylab="Residuos estandarizados")
# dibuja la recta en cero
abline(h=0)
```



CALCULO RMSE

a. Training

```
predict.M30 <- predict(lm.M30, interval = "prediction")
```

```
## Warning in predict.lm(lm.M30, interval = "prediction"): predictions on current data refer to _future_ responses
```

```
calculo_error(train.analysis, as.data.table(predict.M30))
```

```
##           error
## 1: 8.28673e-24
```

b. Test

```
predict.M30 <- predict(lm.M30, test.analysis, interval = "prediction")
calculo_error(test.analysis, as.data.table(predict.M30))
```

```
##           error
## 1: 0.02571554
```