

Modelos velocidad

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.15/100,
                                     vmed.1 = vmed.15/100,
                                     carga.2 = carga.30/100,
                                     vmed.2 = vmed.30/100,
                                     carga.3 = carga.45/100,
                                     vmed.3 = vmed.45/100,
                                     carga.4 = carga.60/100,
                                     vmed.4 = vmed.60/100)]
```

```
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(vmed ~ 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
## -3.3985128567  0.0712716305 -0.0043842684 -0.3837662472  3.9787764459
##      carga.2      vmed.2      vmed.3      diaMes      Mes
##  0.3882943012  0.9681071491 -0.1382238980 -0.0001306308 -0.0015828375
##    prec_norm    var.vmed.1    var.vmed.2    var.vmed.3    diaLunes
## -0.0151097661  0.0167749579  0.0286037227 -0.0504845113 -0.0258240922
##    diaMartes    diaMiercoles    diaJueves    diaViernes    n.festivo
## -0.0302338145 -0.0294944149 -0.0327173784 -0.0292258918  0.0087526251
```

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

```
##
## Call:
## lm(formula = vmed ~ vel.med + carga.med + carga.1 + vmed.1 +
##      carga.2 + vmed.2 + vmed.3 + diaMes + Mes + prec_norm + var.vmed.1 +
##      var.vmed.2 + var.vmed.3 + diaLunes + diaMartes + diaMiercoles +
##      diaJueves + diaViernes + n.festivo, data = train.analisis)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -4.8821 -0.0824  0.0179  0.1167  7.5982
##
## Coefficients:
##              Estimate Std. Error  t value Pr(>|t|)
## (Intercept) -3.399e+00  3.140e-03 -1082.188 < 2e-16 ***
## vel.med      7.127e-02  6.093e-04  116.979 < 2e-16 ***
## carga.med    -4.384e-03  4.807e-03   -0.912   0.362
## carga.1     -3.838e-01  7.489e-03  -51.241 < 2e-16 ***
## vmed.1       3.979e+00  7.536e-03  527.994 < 2e-16 ***
## carga.2      3.883e-01  7.492e-03   51.829 < 2e-16 ***
## vmed.2       9.681e-01  1.012e-02   95.700 < 2e-16 ***
## vmed.3      -1.382e-01  8.089e-03  -17.088 < 2e-16 ***
## diaMes      -1.306e-04  4.228e-05   -3.090   0.002 **
## Mes         -1.583e-03  1.081e-04  -14.647 < 2e-16 ***
## prec_norm   -1.511e-02  2.635e-03   -5.735 9.74e-09 ***
## var.vmed.1    1.677e-02  2.175e-03    7.711 1.25e-14 ***
## var.vmed.2    2.860e-02  2.280e-03   12.544 < 2e-16 ***
## var.vmed.3   -5.048e-02  1.570e-03  -32.165 < 2e-16 ***
## diaLunes     -2.582e-02  1.440e-03  -17.933 < 2e-16 ***
## diaMartes    -3.023e-02  1.425e-03  -21.221 < 2e-16 ***
## diaMiercoles -2.949e-02  1.441e-03  -20.465 < 2e-16 ***
## diaJueves    -3.272e-02  1.448e-03  -22.593 < 2e-16 ***
## diaViernes   -2.923e-02  1.409e-03  -20.736 < 2e-16 ***
## n.festivo     8.753e-03  5.862e-04   14.931 < 2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.349 on 897415 degrees of freedom
## Multiple R-squared:  0.8665, Adjusted R-squared:  0.8665
## F-statistic: 3.067e+05 on 19 and 897415 DF,  p-value: < 2.2e-16
```

```
setwd(ruta)
saveRDS(lm.M30, './modelos/lmM30_vel_15min.RData')
```

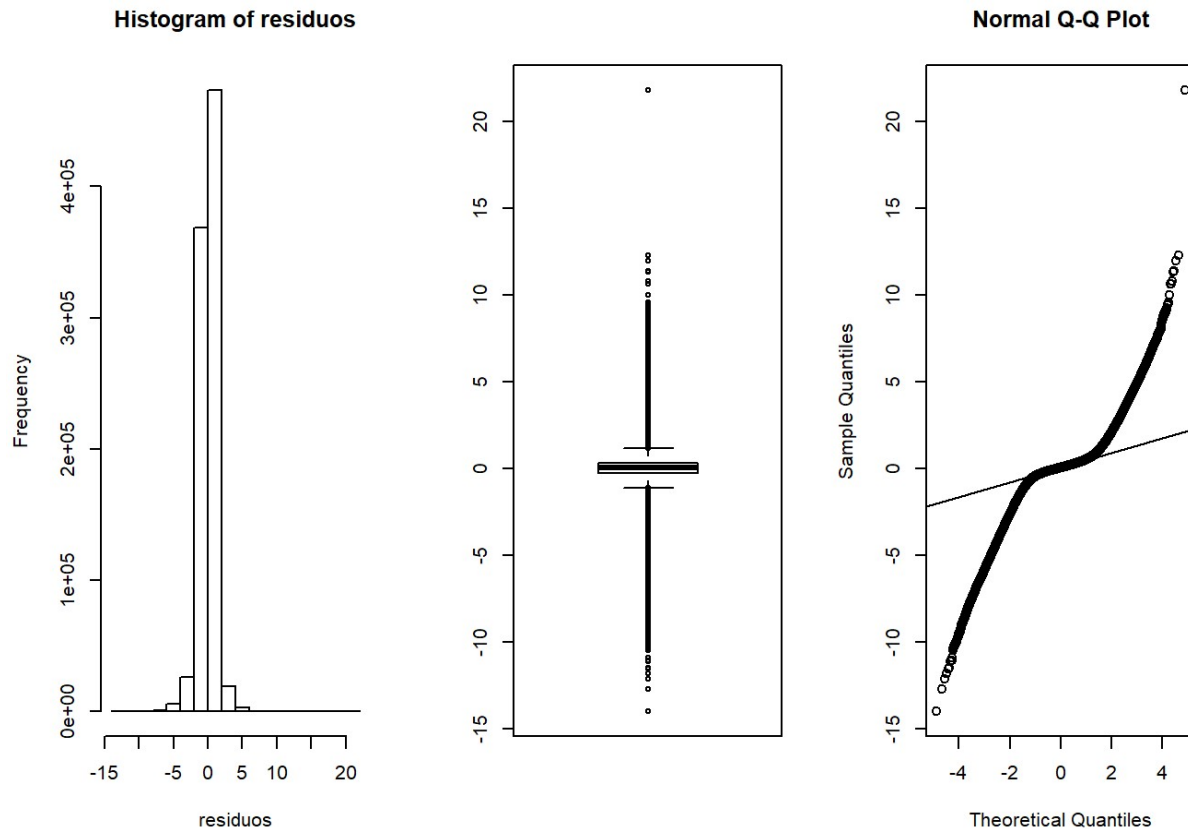
Análisis de los residuos

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:

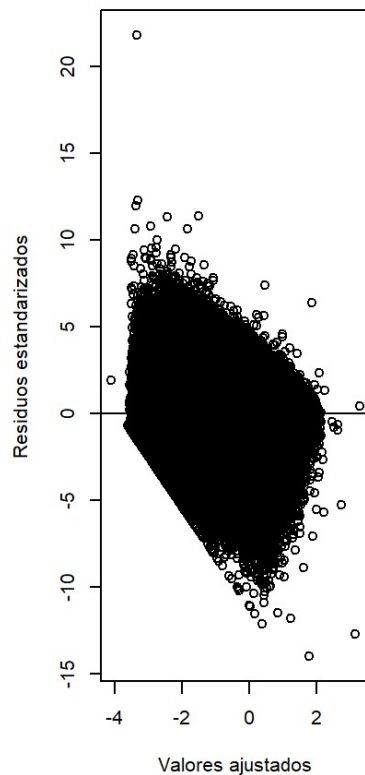
- No es constante, tiene tendencia lo que indica que hay una variable desconocida que impacta en el tráfico pero no la hemos detectado.

```

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: 63295.12
```

b. Test

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

```
##          error
## 1: 26444.68
```