

# Taller 2 - Gestión de datos

Calidad del aire - Estación La Flora Santiago de Cali

*Ronald Fernando Rodríguez Barbosa*

*Maestría en Ingeniería de Sistemas y Computación*

*Maestría en Analítica para la Inteligencia de Negocios*

*Pontificia Universidad Javeriana*

*18 de Mayo de 2019*

## Introducción

La **calidad del aire**, se define como la cantidad general de polución presente en un area y como la pureza promedio atmosférica en relación a las medidas de descarga tomadas de una fuente de polución (Gooch 2007). La contaminación del aire, representa un importante riesgo medioambiental para la salud, bien sea en los países desarrollados o en los países en desarrollo ya que es evidenciado en casos en morbilidad por trastornos cerebrovasculares, cánceres de pulmón y neumopatías crónicas y agudas. Por lo tanto, cuanto más bajos sean los niveles de contaminación del aire, mejor será la salud cardiovascular y respiratoria de la población a largo y a corto plazo (OMS 2016).

De los 23 países de América Latina y el Caribe, 18 tienen sus propias regulaciones vigentes en la actualidad relacionadas con la calidad del aire, que son de acceso público en los sitios web oficiales (Morantes et al. 2016). La trazabilidad para tales regulaciones se establece para los contaminantes de criterio (PM10, PM2.5, SO2, NO2, O-3, CO), utilizando como referencia la secuencia histórica de estándares de la Agencia de Protección Ambiental de los Estados Unidos (USEPA 2019) y los valores de referencia de La Organización Mundial de la Salud (OMS-WHO 2019).

En Colombia, el Sistema de Vigilancia de la Calidad del Aire de Cali (SVCAC 2019), opera bajo la coordinación y administración del Departamento Administrativo de Gestión del Medio Ambiente (DAGMA 2019), Grupo de Calidad del Aire. El Sistema de Vigilancia de Calidad del Aire de Santiago de Cali SVCASC fue acreditado en la norma NTC-ISO/IEC 17025 del año 2005 por el IDEAM a través de la Resolución 1328 del 23 de junio de 2018. El SVCASC actualmente funciona con nueve (9) estaciones: La estación La Flora, ubicada en el barrio La Flora en la zona norte; La estación ubicada en el barrio obrero y la estación Ermita ubicada en el barrio San Pedro ambas en la zona centro; la estación transitoria EDB - Navarro ubicada en el barrio Poblado en la zona oriente; la estación base Aérea, ubicada en el acuparque de la Caña en zona nororiente; la estación Pance ubicada en la Zona Rural; la estación Univalle ubicada en el barrio Meléndez en la zona sur; la estación compartir ubicada en el Barrio Compartir en la zona oriente y la estación Cañaveralejo ubicada en la estación SITM del MIO en la zona suroccidente.

El presente trabajo, realizará un análisis de la captura de datos de la estación *La Flora*, la cual es una estación automática que reporta información horaria al centro de control del DAGMA. Esta estación, mide los niveles de Material Particulado Menor a 10 micrómetros (PM10), Dióxido de Azufre (SO2), Dióxido de Nitrógeno (NO2), Monóxido de Carbono (CO), Ozono (O3) y variables meteorológicas como velocidad del viento, dirección del viento, temperatura, humedad, radiación solar y precipitación. El proceso de preparación y análisis de los datos seguirá la siguiente estructura:

1. **Carga y exploración:** Se incluirá el archivo de datos y se identificarán las características del conjunto de datos
2. **Limpieza de datos:** A partir de la exploración, se definirán los procedimientos para la limpieza de datos para facilitar el análisis
3. **Creación de la vista minable:** Se establecerán los conjuntos de datos finales y su correspondiente análisis
4. **Conclusiones e infografía:** Se compilará el conocimiento adquirido y las cifras de interés adquiridas.

# 1. Carga y exploración

## Tablas de resumen

Para los procedimientos de carga y exploración, se emplearán las herramientas RStudio y RapidMiner. Inicialmente, se realiza la carga de archivos con el fin de resumir las características generales de los datos.

```
datos_base_flora<-read.csv(file="data/dataCAFlora.csv",
                           header = TRUE,sep = ",", dec = ".",fileEncoding = "latin1")
str(datos_base_flora,vec.len=0)

## 'data.frame':    84629 obs. of  12 variables:
##  $ Fecha...Hora      : Factor w/ 75869 levels "01/01/2011 01:00:00 AM",...: NULL ...
##  $ PM10...ug.m3.     : Factor w/ 1390 levels "0.1","0.2","0.3",...: NULL ...
##  $ SO2...ug.m3.      : Factor w/ 3968 levels "0.03","0.05",...: NULL ...
##  $ NO2...ug.m3.      : Factor w/ 5036 levels "0.02","0.13",...: NULL ...
##  $ CO...ug.m3.       : Factor w/ 8027 levels "1000.10","1000.35",...: NULL ...
##  $ O3...ug.m3.       : Factor w/ 6397 levels "0","0.0","0.02",...: NULL ...
##  $ Vel.Viento...m.s. : Factor w/ 63 levels "0","0.1","0.2",...: NULL ...
##  $ Dir.Viento..Grados : Factor w/ 3602 levels "0","0.1","0.2",...: NULL ...
##  $ Temperatura..CÃ.. : Factor w/ 172 levels "16.2","16.3",...: NULL ...
##  $ Humedad....       : Factor w/ 708 levels "100","100.3",...: NULL ...
##  $ Radiacion.Solar..Watt.M2.: Factor w/ 7176 levels "0","0.1","0.2",...: NULL ...
##  $ Lluvia..mm.       : Factor w/ 182 levels "#","0","0.1",...: NULL ...
```

El conjunto de datos contiene un total 84.629 observaciones con 12 variables. La descripción de las variables se relaciona a continuación:

Variable	Tipo de variable	Descripción
Fecha & Hora	Fecha	Fecha y hora de la captura del senso de polutantes
PM10 (ug/m3)	Contínua	Concentración de Material Particulado Menor a 10 micrómetros
SO2 (ug/m3)	Contínua	Concentración de Dióxido de Azufre
NO2 (ug/m3)	Contínua	Concentración de Dióxido de Nitrógeno
CO (ug/m3)	Contínua	Concentración de Monóxido de Carbono
O3 (ug/m3)	Contínua	Concentración de Ozono
Vel Viento (m/s)	Contínua	Velocidad del viento en metros por segundo
Dir Viento (Grados)	Contínua	Dirección del viento
Temperatura (C°)	Contínua	Temperatura en grados celsius
Humedad (%)	Contínua	Porcentaje de humedad
Radiacion Solar (Watt/M2)	Contínua	Radiación Solar
Lluvia (mm)	Contínua	Cantidad de precipitaciones

```
summary(datos_base_flora[1:3])
```

```
##           Fecha...Hora      PM10...ug.m3.      SO2...ug.m3.
## 01/01/2011 01:00:00 AM:    2      ND      :17551      ND      :64788
## 01/01/2011 01:00:00 PM:    2      38      : 421      9.56      : 25
## 01/01/2011 02:00:00 AM:    2      48      : 374      11.09     : 24
## 01/01/2011 02:00:00 PM:    2      27      : 344      5.26      : 24
## 01/01/2011 03:00:00 AM:    2      32      : 342      6.77      : 24
## 01/01/2011 03:00:00 PM:    2      39      : 342      8.57      : 24
## (Other)           :84617      (Other):65255      (Other):19720
```

```
summary(datos_base_flora[4:6])
```

```
## NO2...ug.m3. CO...ug.m3. O3...ug.m3.
## ND :58116 ND :69992 ND :53828
## 18.48 : 25 1956.56: 9 0.0 : 640
## 20.91 : 25 1100.01: 8 0.5 : 238
## 20.26 : 23 1245.79: 8 0.6 : 207
## 15.03 : 21 1493.70: 8 0.1 : 206
## 22.82 : 21 1526.33: 8 4.3 : 189
## (Other):26398 (Other):14596 (Other):29321
```

```
summary(datos_base_flora[7:9])
```

```
## Vel.Viento...m.s. Dir.Viento..Grados. Temperatura..CÃ..
## ND :39017 ND :39017 ND :39018
## 0.2 : 4434 0 : 52 22.5 : 702
## 0.3 : 4184 33.6 : 46 22.8 : 644
## 0.4 : 3811 48.8 : 41 21.9 : 642
## 0.1 : 3528 272.6 : 40 22.4 : 639
## 0.5 : 3342 87.1 : 39 22.2 : 634
## (Other):26313 (Other):45394 (Other):42350
```

```
summary(datos_base_flora[10:12])
```

```
## Humedad.... Radiacion.Solar..Watt.M2. Lluvia..mm.
## ND :39017 ND :26960 0 :72107
## 81.5 : 139 0 :26733 ND : 7018
## 81.9 : 138 665 : 821 0.25 : 1874
## 78.9 : 133 664 : 587 0.51 : 563
## 79.9 : 132 0.1 : 414 0.76 : 344
## 78.7 : 131 666 : 379 1.02 : 280
## (Other):44939 (Other):28735 (Other): 2443
```

El DAGMA mediante el grupo de calidad del aire en el transcurso de los años ha trabajado en robustecer el sistema de calidad del aire para garantizar la continuidad en los datos, sin embargo es normal que los sistemas de monitoreo de calidad del aire presenten discontinuidad en los datos (datos faltantes o espacios en blanco) debido a dos situaciones: la primera corresponde a las anomalías que se dan en las estaciones de monitoreo, tales como: Fallas en los equipos, falta de energía eléctrica en la zona, hurto de equipos o cableado, mantenimiento o cambio de equipos, etc y la segunda causa corresponde a la inclusión o exclusión de algunos contaminantes o variables meteorológicas, según criterio de los expertos y característica de la zona a monitorear.

```
## [1] 37914.01188 836.02325 3503.71392 4039.33929 7319.93788
## [6] 5094.69424 33.57069 2626.20164 118.98171 547.83351
## [11] 3504.55330 19.25865
```

Métricas de tendencia central

Métricas de dispersión

Visualización de métricas de dispersión

Análisis de las visualizaciones

Análisis de Correlación

Recuerde incluir análisis de todo lo visualizado y no solo el gráfico

## 2. Limpieza de datos

Reconocimiento y tratamiento de atributos con valores únicos o distintos

Reconocimiento y tratamiento de atributos con valores faltantes

Reconocimiento y tratamiento de atributos con valores atípicos

Reconocimiento y tratamiento de registros atípicos

Reconocimiento y tratamiento de atributos redundantes

## 3. Creación de la vista minable

Generación de variables derivadas tipo 1 y 2

Normalización de al menos un atributo

Discretización de al menos un atributo

Numerización 1 a n de al menos un atributo

## 4. Conclusiones e Infografía

## Referencias

DAGMA, Departamento Administrativo de Gestión del Medio Ambiente. 2019. “Sitio Oficial - Departamento Administrativo de Gestión Del Medio Ambiente.” <http://www.cali.gov.co/dagma/>.

Gooch, Jan W., ed. 2007. “Ambient Air Quality.” In *Encyclopedic Dictionary of Polymers*, 48–48. New York, NY: Springer New York. doi:10.1007/978-0-387-30160-0\_522.

Morantes, Gioberti, Narciso Perez, Rafael Santana, and Gladys Rincon. 2016. “A REVIEW OF THE REGULATORY INSTRUMENTS FOR AIR QUALITY AND ATMOSPHERIC MONITORING SYSTEMS:

LATIN AMERICA AND THE CARIBBEAN.” *INTERCIENCIA* 41 (4): 235–42.

OMS, Organizacion Mundial de la Salud. 2016. “Calidad Del Aire Ambiente Y Salud.” <http://origin.who.int/mediacentre/factsheets/fs313/es/>.

OMS-WHO, Organizacion Mundial de la Salud. 2019. “Sitio Oficial - Organización Mundial de La Salud.” <https://www.who.int/es/home/>.

SVCAC, Sistema de Vigilancia de Calidad del Aire de Cali. 2019. “Sitio Oficial - Sistema de Vigilancia de Calidad Del Aire de Cali.” [http://www.cali.gov.co/dagma/publicaciones/38365/sistema\\_de\\_vigilancia\\_de\\_calidad\\_del\\_aire\\_de\\_cali\\_svcac/](http://www.cali.gov.co/dagma/publicaciones/38365/sistema_de_vigilancia_de_calidad_del_aire_de_cali_svcac/).

USEPA, Agencia de Protección Ambiental de los Estados Unidos. 2019. “Sitio Oficial - Agencia de Protección Ambiental de Los Estados Unidos.” <https://www.epa.gov/>.