



Introducción a la Ciencia de Datos 2024-2025

MASTER CIENCIA DE DATOS

UNIVERSIDAD DE GRANADA

Proyecto Final

MIGUEL GARCÍA LÓPEZ

Índice

1. Introducción	6
1.1. Tareas	6
1.1.1. Análisis de Datos	6
1.1.2. Regresión	7
1.1.3. Clasificación	7
2. EDA Regresión	8
2.1. Resumen de los datos	8
2.2. Normalidad	11
2.3. Outliers univariantes	11
2.4. Outliers multivariantes	13
2.5. Correlación	14
2.5.1. Construcción de características	15
2.6. PCA	16
3. EDA Clasificación	17
3.1. Resumen de los datos	17
3.2. Normalidad	17
3.3. Outliers	19
3.4. Distribución por clases	20
3.5. Correlación	22
3.6. Transformaciones	23
4. Regresión	24
4.1. Regresión simple	24

4.2. Regresión múltiple	26
4.2.1. Backward selection	26
4.3. Interacciones y no linealidad	26
4.4. KNN	28
4.5. Comparativas	28
5. Clasificación	31
5.1. KNN	31
5.2. LDA	32
5.2.1. Asunciones	32
5.2.2. Ajuste	34
5.3. QDA	37
5.4. Comparación	37
6. Apéndice	38
6.1. EDA Regresión	38
6.2. EDA Clasificación	44
6.3. Regresión	51
6.4. Clasificación	62
7. Bibliografía	71

Índice de figuras

1. Max temperature.	9
2. Min temperature.	9
3. Mean temperature.	9

4.	Dewpoint.	9
5.	Precipitation.	10
6.	Sea level pressure.	10
7.	Standard pressure.	10
8.	Visibility.	10
9.	Wind speed.	10
10.	Max wind speed.	10
11.	standardization formula (Z-score).	11
12.	QQplot for every variable in dataset.	12
13.	Shapiro test for normality.	12
14.	Boxplots for regression.	13
15.	Hz test for multivariate normality.	14
16.	Correlation between variables.	14
17.	Mean temperature value with sea level pressure.	15
18.	Biplot for Wankara.	16
19.	THS_value.	18
20.	T3resin.	18
21.	Thyroidstimulating.	18
22.	Thyroxin.	18
23.	Triiodothyronine.	18
24.	Qqplot for classification EDA.	19
25.	Boxplot for newthyroid.	20
26.	Class distribution for T3resin.	21
27.	Class distribution for Triiodothyronine.	21
28.	Class distribution for Thyroidstimulating.	21

29.	Class distribution for TSH _{value}	21
30.	Class distribution for Thyroxin.	21
31.	Biplot for classification EDA.	22
32.	Pairs plot EDA classification.	23
33.	Correlation plot EDA classification.	24
34.	Diff temperature vs Mean Temperature.	25
35.	Diff pressure vs Mean Temperature.	25
36.	Visibility vs Mean Temperature.	25
37.	Dewpoint vs Mean Temperature	25
38.	Wind speed vs Mean Temperature.	25
39.	K-fold cross validation results in linear regression.	29
40.	K-fold cross validation results in KNN	29
41.	K-fold cross validation results in M5	30
42.	confusion Matrix for $k = 1$	33
43.	confusion Matrix for $k = 5$	33
44.	confusion Matrix for $k = 15$	33
45.	confusion Matrix for $k = 40$	33
46.	confusion Matrix for $k = 100$	33
47.	QQPlot for T3resin grouped by class.	35
48.	QQPlot for Thyroxin grouped by class.	35
49.	QQPlot for Triiodothyronine grouped by class.	35
50.	QQPlot for Thyroidstimulating grouped by class.	35
51.	QQPlot for TSH_value grouped by class.	35
52.	Data transformation after LDA	36
53.	Confusion Matrix for QDA	37

Índice de cuadros

1.	Summary statistics for weather-related variables.	8
2.	Summary statistics for thyroid-related variables.	17
3.	Class distribution and proportions.	17
4.	Linear Simple Regression Results	26
5.	Summary of simple linear regression.	26
6.	Summary of results in k-fold validation.	28
7.	Friedman test + Holms (1: Linear Regression, 2: KNN, 3: M5) en <i>test</i> . . .	31
8.	Friedman test + Holms (1: Linear Regression, 2: KNN, 3: M5) en <i>training</i> . .	31
9.	K values and accuracies for KNN	32
10.	Normality test for every variable grouped by class.	36
11.	Summary of results in k-fold validation.	38

1. Introducción

En el presente proyecto final de la asignatura de **Introducción a la Ciencia de Datos**, se llevarán a cabo tres trabajos fundamentales en el ámbito de la ciencia de datos, centrados en el análisis de datos y la creación de modelos predictivos. Estos trabajos son los siguientes:

- **Análisis exploratorio de datos (EDA):** Análisis de dos conjuntos de datos, uno correspondiente a un problema de **regresión** y otro a un problema de **clasificación**.
- **Regresión:** Aplicación de modelos de regresión lineal simple y la realización de pruebas y análisis correspondientes.
- **Clasificación:** Ajuste de diversos modelos de clasificación, incluyendo kNN , y realización de comparaciones entre los diferentes algoritmos.

Para cada uno de los problemas mencionados, se dispondrán de dos conjuntos de datos diferentes: uno para la clasificación y otro para la regresión.

El documento se planteará con la siguiente estructura, primero los **EDAs** correspondientes a regresión y clasificación y después las tareas correspondientes al modelado de los predictores para cada problema junto con sus sub-tareas.

1.1. Tareas

1.1.1. Análisis de Datos

En este apartado, el estudiante deberá llevar a cabo un análisis preliminar de los dos conjuntos de datos asignados, denominados **datasetR** y **datasetC**. Este análisis deberá abarcar los siguientes puntos:

- **A-1 Descripción del tipo de datos de entrada:** Es necesario especificar el formato de los datos (por ejemplo, lista, data frame, etc.), el número de filas y columnas, así como los tipos de datos atómicos presentes.
- **A-2 Cálculo de medidas estadísticas:** Se deben calcular medidas estadísticas como la media, la desviación estándar y otras métricas descriptivas relevantes.
- **A-3 Gráficos:** Se deberán generar gráficos adecuados que faciliten la visualización de los datos.
- **A-4 Descripción del conjunto de datos:** A partir de los análisis anteriores, se debe proporcionar una descripción detallada del conjunto de datos.

1.1.2. Regresión

En este apartado, el estudiante deberá utilizar el conjunto de datos `datasetR` para llevar a cabo las siguientes tareas:

- **R-1 Regresión lineal simple:** Se aplicará el algoritmo de regresión lineal simple a cada regresor (variable de entrada) para obtener los modelos correspondientes. En el caso de que `datasetR` contenga más de cinco regresores, se deberá seleccionar de manera justificada los cinco más relevantes. Posteriormente, se elegirá el modelo que se considere más adecuado según las medidas de calidad conocidas.
- **R-2 Regresión lineal múltiple:** Se aplicará el algoritmo de regresión lineal múltiple. Es necesario justificar si el modelo obtenido representa una mejora respecto al modelo elegido en el paso anterior. En este análisis se deberán tener en cuenta posibles interacciones entre las variables y la no linealidad.
- **R-3 Algoritmo k-NN para regresión:** Se aplicará el algoritmo k-NN para la regresión.
- **R-4 Comparación de resultados:** Se realizará una comparación entre los resultados obtenidos con los dos algoritmos de regresión múltiple. Además, se realizarán comparativas con un tercer modelo, el modelo de regresión M5', cuyos resultados ya están disponibles en las tablas de resultados.

Nota: Al final de las transparencias de las clases de laboratorio sobre regresión, el estudiante encontrará aclaraciones adicionales sobre los apartados R-1 a R-4. Estas aclaraciones serán más comprensibles una vez que se haya realizado el trabajo práctico correspondiente.

1.1.3. Clasificación

En este apartado, el estudiante deberá utilizar el conjunto de datos `datasetC` asignado para realizar las siguientes tareas:

- **C-1 Algoritmo k-NN para clasificación:** Se aplicará el algoritmo k-NN utilizando diferentes valores de k. Se deberá seleccionar el valor de k más adecuado para el conjunto de datos. Además, se analizará el comportamiento de la precisión en los conjuntos de entrenamiento y prueba con los distintos valores de k.
- **C-2 Algoritmo LDA para clasificación:** Se aplicará el algoritmo de Análisis Discriminante Lineal (LDA) para la clasificación. Es fundamental verificar que se cumplen las asunciones necesarias para la correcta aplicación del algoritmo.

- **C-3 Algoritmo QDA para clasificación:** Se aplicará el algoritmo de Análisis Discriminante Cuadrático (QDA) para la clasificación, verificando también las asunciones necesarias.
- **C-4 Comparación de algoritmos:** Finalmente, se realizará una comparación entre los resultados obtenidos con los tres algoritmos mencionados.

2. EDA Regresión

El *dataset* asignado para el problema de regresión es el de **wankara**. Este conjunto de datos incluye una serie de variables relacionadas con el tiempo de la ciudad Turca de Ankara.

2.1. Resumen de los datos

Se obtienen algunos de los estadísticos relacionados con el *dataset*.

Variable	Min.	1st Qu.	Median	Mean	3rd Qu.	Max.
Max_temperature	23.00	46.40	60.80	61.56	77.40	100.00
Min_temperature	-7.10	26.60	36.00	37.08	48.20	65.50
Dewpoint	-3.10	28.50	36.80	36.29	45.30	57.60
Precipitation	0.0000	0.0000	0.0000	0.0541	0.0000	4.0000
Sea_level_pressure	29.46	29.83	29.96	29.98	30.11	30.60
Standard_pressure	26.30	26.69	26.77	26.78	26.87	27.18
Visibility	0.200	7.400	8.300	7.718	8.600	11.500
Wind_speed	0.000	3.110	5.060	5.393	7.250	18.000
Max_wind_speed	2.19	10.20	12.70	13.32	16.10	57.40
Mean_temperature	7.90	36.70	48.50	49.56	63.30	81.80

Cuadro 1: Summary statistics for weather-related variables.

El conjunto de datos cuenta con 1609 filas (observaciones) y 9 columnas (características). Todas las características del conjunto son de tipo numérico (son números reales) y no hay ningún valor faltante, por lo que no es necesario realizar ningún tipo de imputación sobre los datos o eliminar observaciones (si fuese plausible).

Dado los datos obtenidos en la tabla 1 se puede observar que los datos son bastante simples. No hay gran desviación en la mayoría de variables. La variable **Precipitación** si tiene datos más extremos, ya que desde su mínimo hasta el tercer cuartil sus valores son cero y su máximo cuatro, lo que significa un salto muy grande.

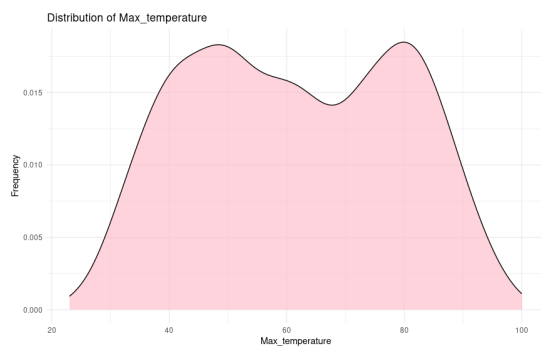


Figura 1: Max temperature.

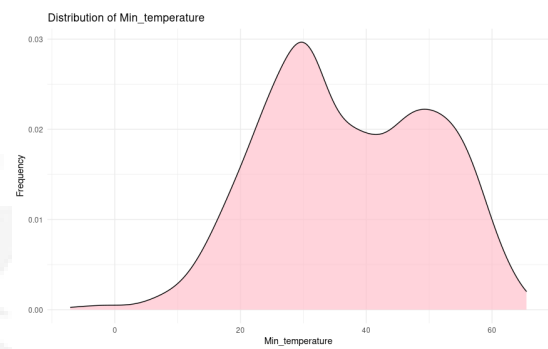


Figura 2: Min temperature.

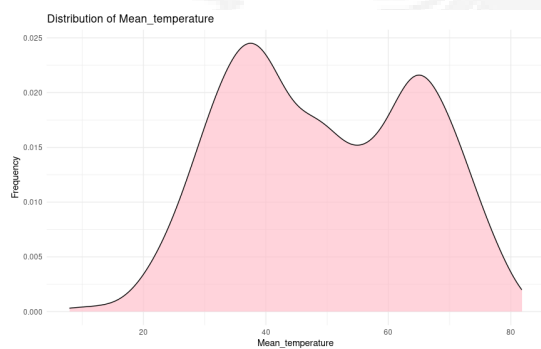


Figura 3: Mean temperature.

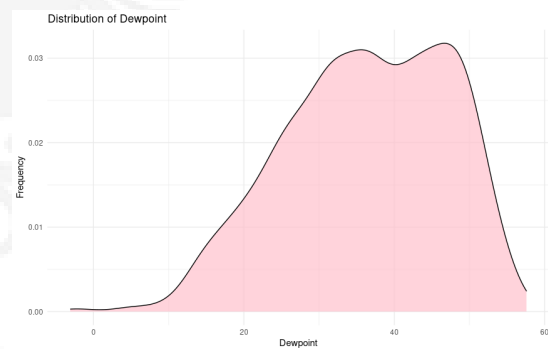


Figura 4: Dewpoint.

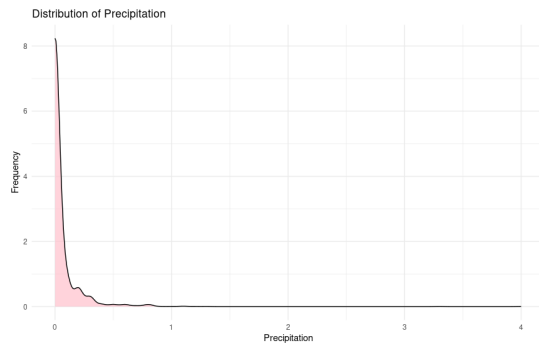


Figura 5: Precipitation.

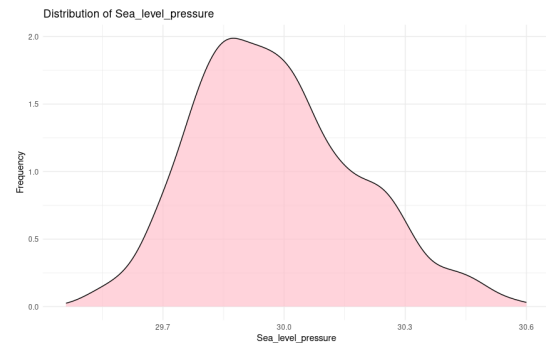


Figura 6: Sea level pressure.

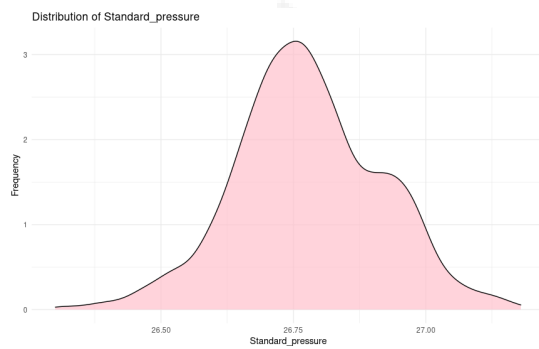


Figura 7: Standard pressure.

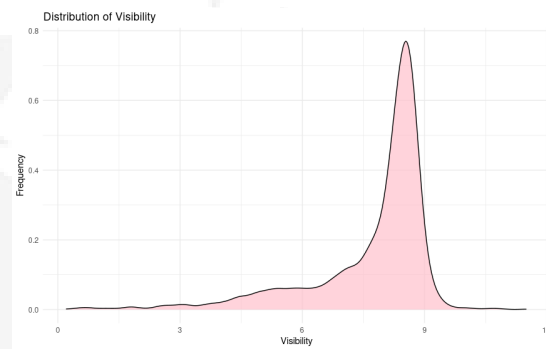


Figura 8: Visibility.

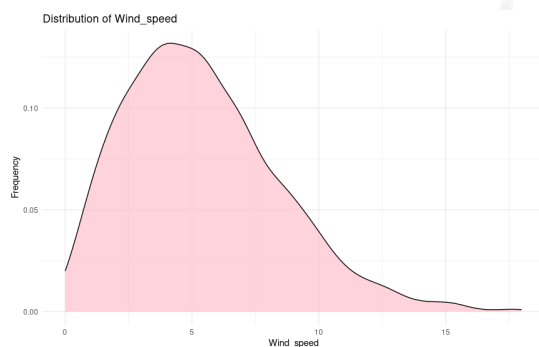


Figura 9: Wind speed.

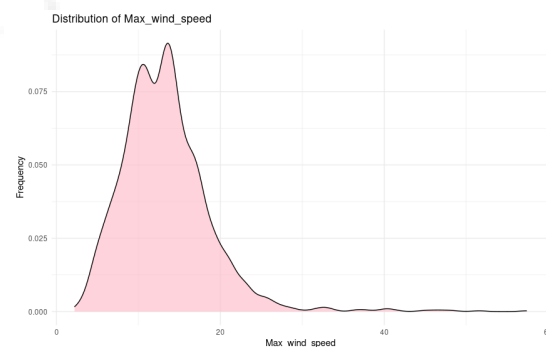


Figura 10: Max wind speed.

2.2. Normalidad

Puede observarse en las figuras que hay algunas variables con dos picos, las relacionadas con la temperatura concretamente (figuras 1 2 3). Esto suele indicar una distribución de tipo bimodal. Hay otras variables, como la precipitación 5, visibilidad 8 y las relativas al viento, que tienen asimetría. Algunas de ellas son muy extremas, como la precipitación. Aquellas un poco más leves y que son asimétricas hacia la derecha pueden “corregirse” con algunas transformaciones, como la logarítmica.

La bimodalidad en las temperaturas puede darse por varios motivos, pero el que parece más probable es la toma de datos en distintas zonas geográficas de Ankara, dando lugar a distintas modas. El calentamiento climático podría ser factible, pero no sería tan evidente y se necesitarían tomas de muchos años.

Se estandarizan los datos para que estos obtengan una media de 1 y varianza 0. Escalar permite que todas las variables contribuyan equitativamente a la distancia calculada, evitando que una variable desproporcionada influya en la detección de *outliers*.

$$z = \frac{x - \mu}{\sigma}$$

Figura 11: standardization formula (Z-score).

Se procede a obtener gráficos de *QQPlot* para todas las variables. Intuitivamente se puede saber que hay algunas que no se van a adecuar a la línea de la normalidad representada por los cuantiles de la distribución normal, pero así es posible descartarlas visualmente.

Como puede observarse en la figura 12, tan solo las relativas a las presiones parecen adecuarse. Es interesante con la precipitación, que tiene una asimetría extrema, tiene la cola derecha muy larga en el *QQPlot*. Las variables de presión al nivel del mar y presión estándar tienen una forma más cercana a la normal, por ello se les aplica un test de *Shapiro* para comprobar su normalidad. El test de *Shapiro* rechaza siempre y cuando haya evidencia estadística de que la distribución no sigue una normal.

Como se observa en la figura 13, se obtienen valores muy significativos, por lo que se rechaza la hipótesis nula y no se pueden considerar las distribuciones de las variables como normales.

2.3. Outliers univariantes

Se procede a analizar las gráficas de *Boxplots* para visualizar posibles valores anómalos en los datos.

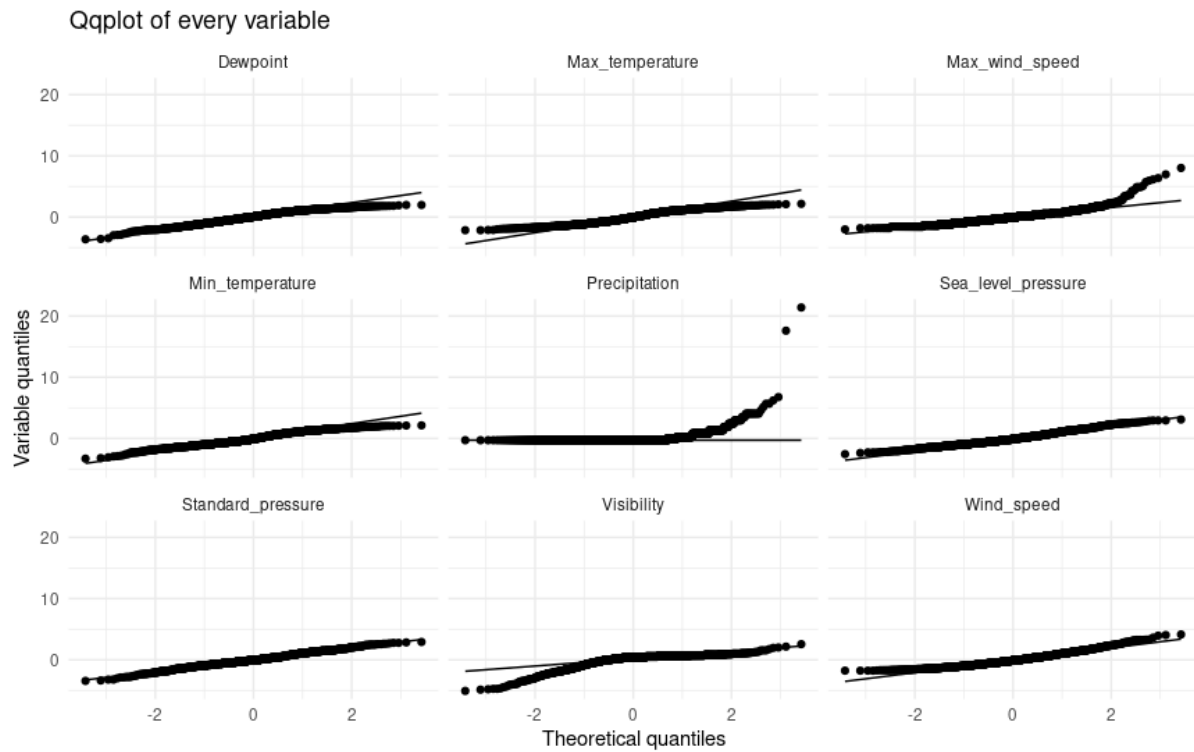


Figura 12: QQplot for every variable in dataset.

```
Shapiro-Wilk normality test
data: data$Sea_level_pressure
W = 0.98532, p-value = 9.76e-12

> shapiro.test(data$Standard_pressure)

Shapiro-Wilk normality test
data: data$Standard_pressure
W = 0.99516, p-value = 4.617e-05
```

Figura 13: Shapiro test for normality.

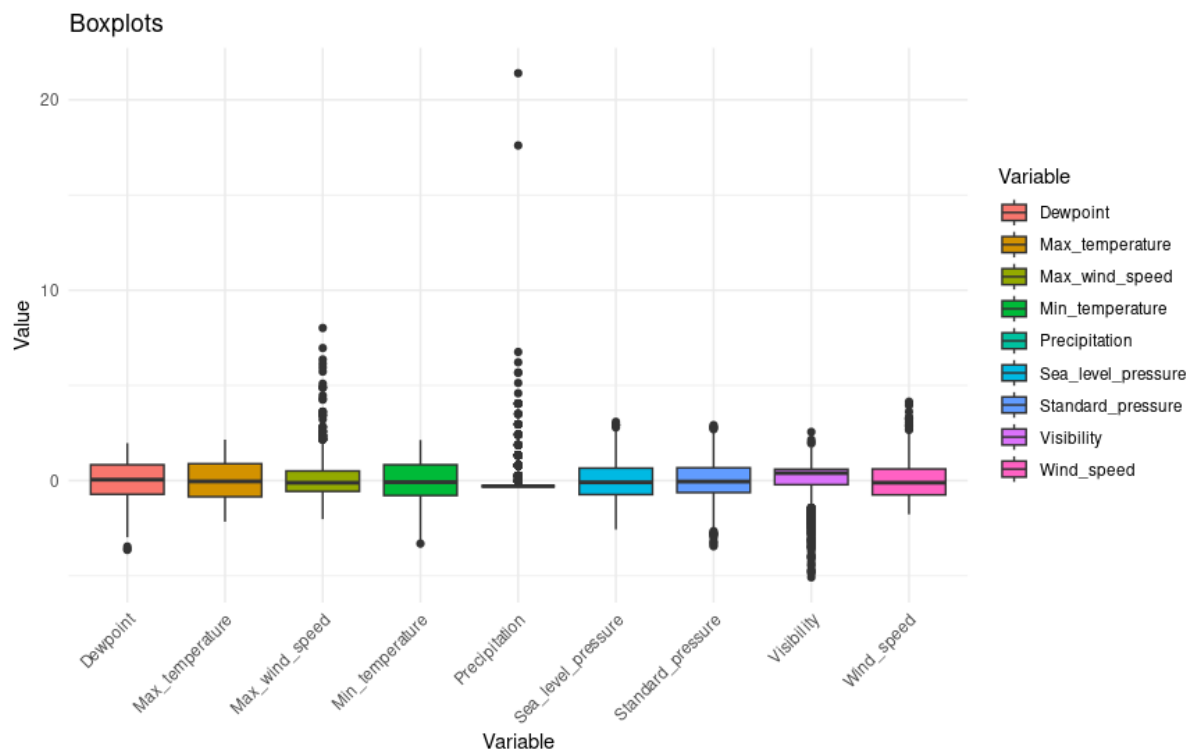


Figura 14: Boxplots for regression.

Puede observarse en la figura 14 como la variable de precipitación tiene lo que podría considerarse valores anómalos extremos. En el resumen obtenido anteriormente de los datos ya fue detectado. Ankara es una zona de pocas precipitaciones, por tanto la mayoría de valores son 0, pero eso no implica que el día que llueva sea un valor irreal. Es anómalo por su rareza, pero debería persistir en el conjunto de datos, pues es un valor real.

Observando el resto de las variables y sus valores máximos y mínimos no se podría considerar que existan *outliers* univariantes, ya que son valores posibles y nada irreales. Con ello se pretende decir que aunque salgan datos anómalos usando métodos como la distancia intercuartil, se decide dejar los datos tal y como están. Además no se cuenta con conocimiento experto para poder asegurar la eliminación de datos.

2.4. Outliers multivariantes

Se procede a analizar la normalidad del conjunto para todas sus dimensiones, esto con el objetivo de encontrar *outliers multivariantes*. Para ello se utiliza primero un *test* como el **HZ** o *Henze-Zirkler*. El test de HZ evalúa si un conjunto de datos sigue una distribución normal multivariada. Es decir, extiende la idea de pruebas de normalidad univariadas (como Shapiro-Wilk) a dimensiones superiores.

Esto se hace para poder utilizar la distancia *Mahalanobis* en la detección de *outliers*.

```
> mvn(data = data, mvnTest = "hz")
$multivariateNormality
      Test      HZ p value MVN
1 Henze-Zirkler 3.87392      0 NO
```

Figura 15: Hz test for multivariate normality.

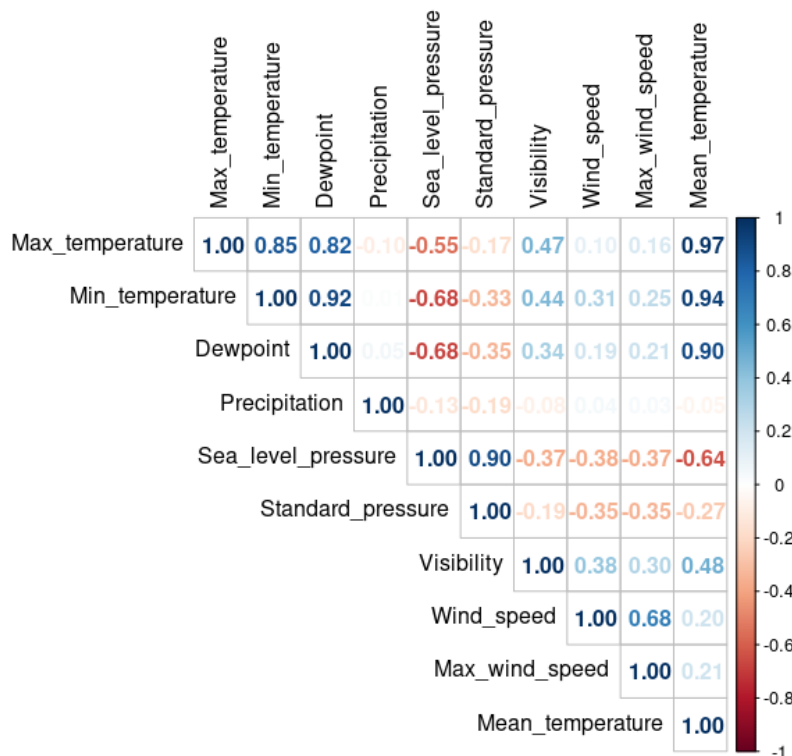


Figura 16: Correlation between variables.

Como indica la figura 15, no puede considerarse que el conjunto sea normal. Esto es obvio, si las variables por separado no lo eran, raramente lo serán en conjunto. Dados estos resultados, no es conveniente usar la distancia de *Mahalanobis*, ya que se asume normalidad multivariante. Se podrían utilizar métodos como **LOF**, pero se exceden el dominio de la asignatura y además no sería coherente con el análisis descrito anteriormente sobre *outliers* univariantes, ya que no sería posible por parte del estudiante saber si los datos son realmente irreales como para eliminarlos.

2.5. Correlación

Se analiza la correlación entre variables del *dataset*.

Según se observa en la figura 16 hay correlación (de tipo lineal) entre algunas variables, sobre todo entre aquellas relativas a la temperatura (mínima, máxima y media). El punto de condensación también está altamente correlacionado con la temperatura media

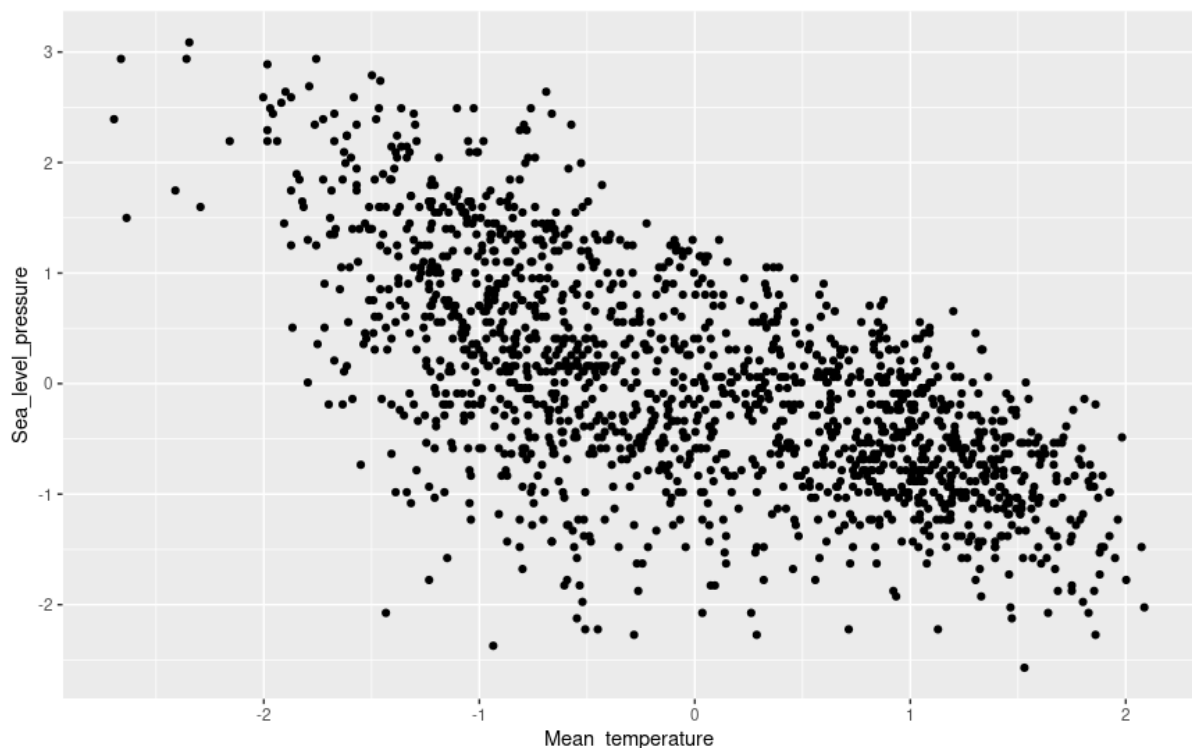


Figura 17: Mean temperature value with sea level pressure.

(variable objetivo). El hecho de que el punto de condensación o *dewpoint* también esté correlacionado con la variable objetivo refuerza la idea de que las condiciones atmosféricas, como la humedad y la temperatura, están relacionadas con el comportamiento de la variable de interés. En particular, el *dewpoint* es un indicador importante de la humedad en el aire, lo que podría influir directamente en el fenómeno que se está modelando.

Las variables de presión a nivel del agua también tienen correlación con las temperaturas. La presión atmosférica disminuye con el aumento de la elevación, por lo que podría intuirse que la temperatura media disminuye en lugares con poca altitud en este conjunto de datos, como puede observarse en la figura 17.

2.5.1. Construcción de características

Variables como *max_temperature* y *min_temperature* parecen redundantes, al igual que las de presión. Si bien correlación no implica que dos variables aporten la misma información, pues pueden ser ambas necesarias, quizá se pueda crear una nueva variable o varias nuevas. Pasar de máximo y mínimo de temperatura a *diff_temperature*. Esta nueva característica sería la suma de las dos anteriormente mencionadas, lo mismo puede realizarse con las relativas a la presión.

Realizando esta simple construcción, se elimina posible multicolinealidad entre las variables de temperatura con *mean_temperature*. Se obtiene una nueva variable que tiene muy

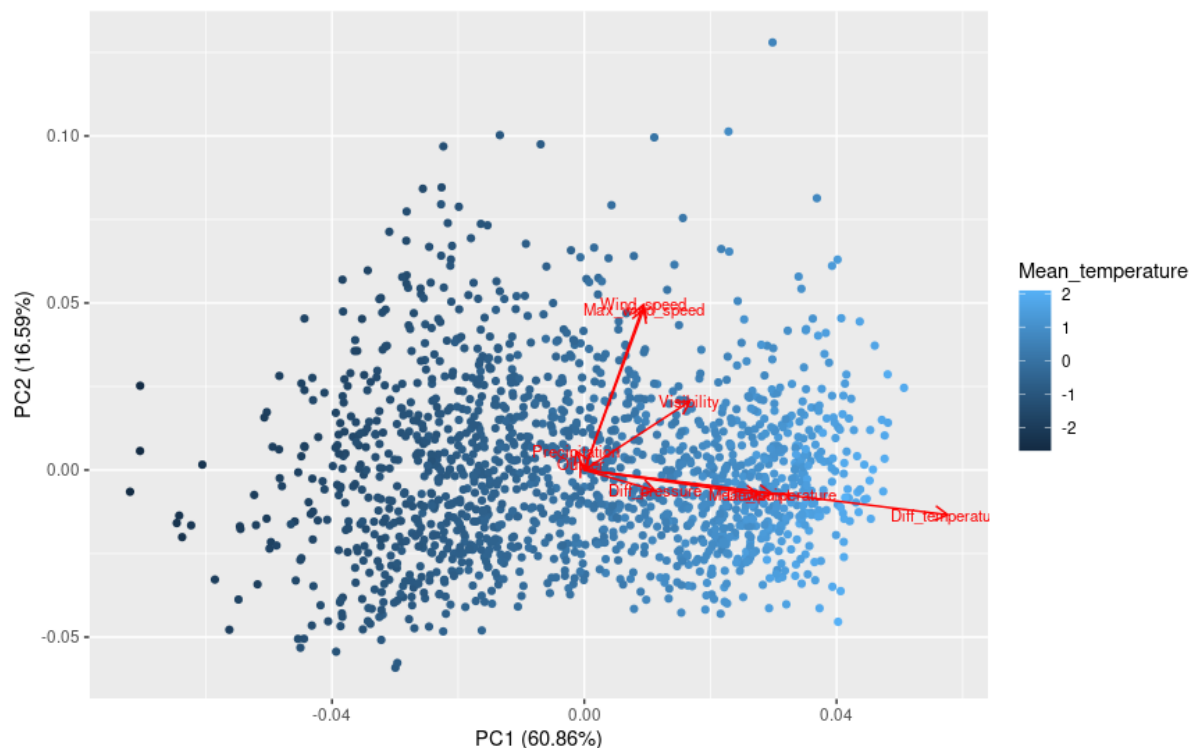


Figura 18: Biplot for Wankara.

buena correlación y no se pierde apenas información.

2.6. PCA

Si bien no es necesario reducir características en este conjunto de datos, pues son solo diez, es interesante para poder analizar ciertos gráficos como el *Biplot*.

Como se puede observar en el *Biplot* 18, allí donde el punto de condensación aumenta, lo hace la temperatura media, así como la presión. La nueva característica de diferencia de temperatura, al ser la resta entre mínimos y máximos de temperatura, crece en dirección perpendicular a la temperatura media. Calculando de nuevo la matriz de correlación, se podría observar que esta nueva característica ya no está correlacionada con la variable objetivo (linealmente al menos).

3. EDA Clasificación

3.1. Resumen de los datos

Se obtiene el siguiente resumen de los datos usando el comando de *summary* de **R**.

Variable	Min.	1st Qu.	Median	Mean	3rd Qu.	Max.
T3resin	65.0	103.0	110.0	109.6	117.5	144.0
Thyroxin	0.500	7.100	9.200	9.805	11.300	25.300
Triiodothyronine	0.20	1.35	1.70	2.05	2.20	10.00
Thyroidstimulating	0.10	1.00	1.30	2.88	1.70	56.40
TSH_value	-0.700	0.550	2.000	4.199	4.100	56.300
Class	1.000	1.000	1.000	1.442	2.000	3.000

Cuadro 2: Summary statistics for thyroid-related variables.

Hay solo 5 variables en todo el *dataset* y un total de 215 observaciones. No hay variables faltantes, por lo que no es necesario utilizar ningún tipo de técnica de imputación. En cuanto a las variables, según lo visto en la tabla 2, se puede decir que en general es necesario realizar un escalado de los datos. Variables como *TSH_value* o *Thyroidstimulating* se distribuyen muy ampliamente y con valores en escalas muy diferentes. Además ha de remarcarse que es necesario transformar la variable de clase (etiqueta *Class*) a un factor para poder tratarlas como categorías.

Existen tres clases, 1, 2 y 3. La proporción de clases en el conjunto de datos es la descrita en la tabla 3.

Class	Total Items	Proportion
1	150	0.6977
2	35	0.1628
3	30	0.1395
Total	215	1.0000

Cuadro 3: Class distribution and proportions.

Claramente existe un desbalance de clases a favor de la clase 1, las otras dos están proporcionadas entre sí. Es interesante analizar este desbalance más a fondo.

3.2. Normalidad

Como se observa en las figuras 19,20,21,22,23 estas presentan un alto grado de asimetría. Queda claro con un primer vistazo que no cumplen la normalidad, de hecho no pasan los test de *Shapiro*.

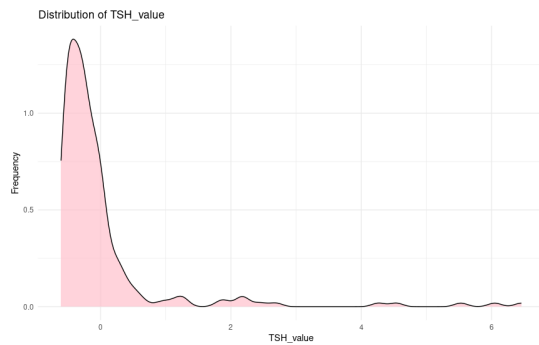


Figura 19: TSH_value.

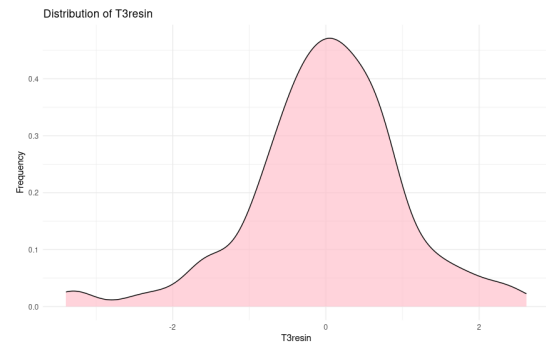


Figura 20: T3resin.

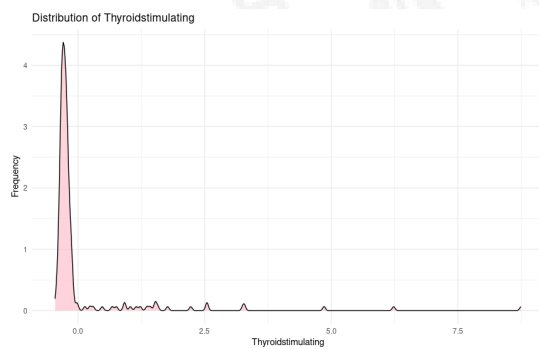


Figura 21: Thyroidstimulating.

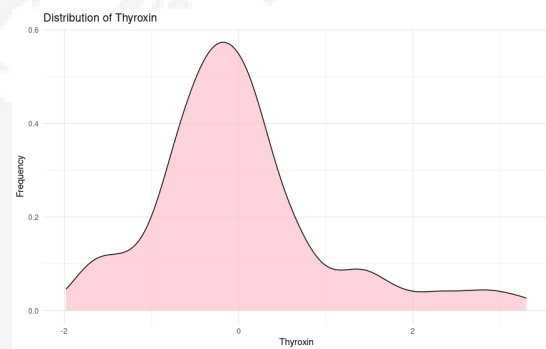


Figura 22: Thyroxin.

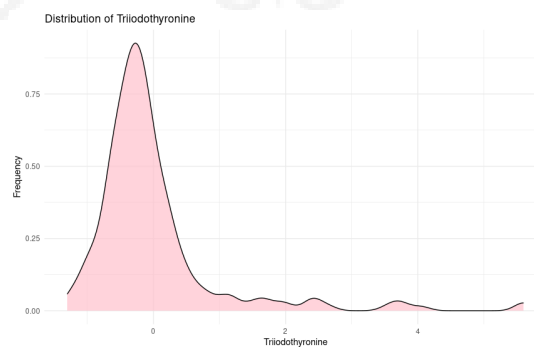


Figura 23: Triiodothyronine.

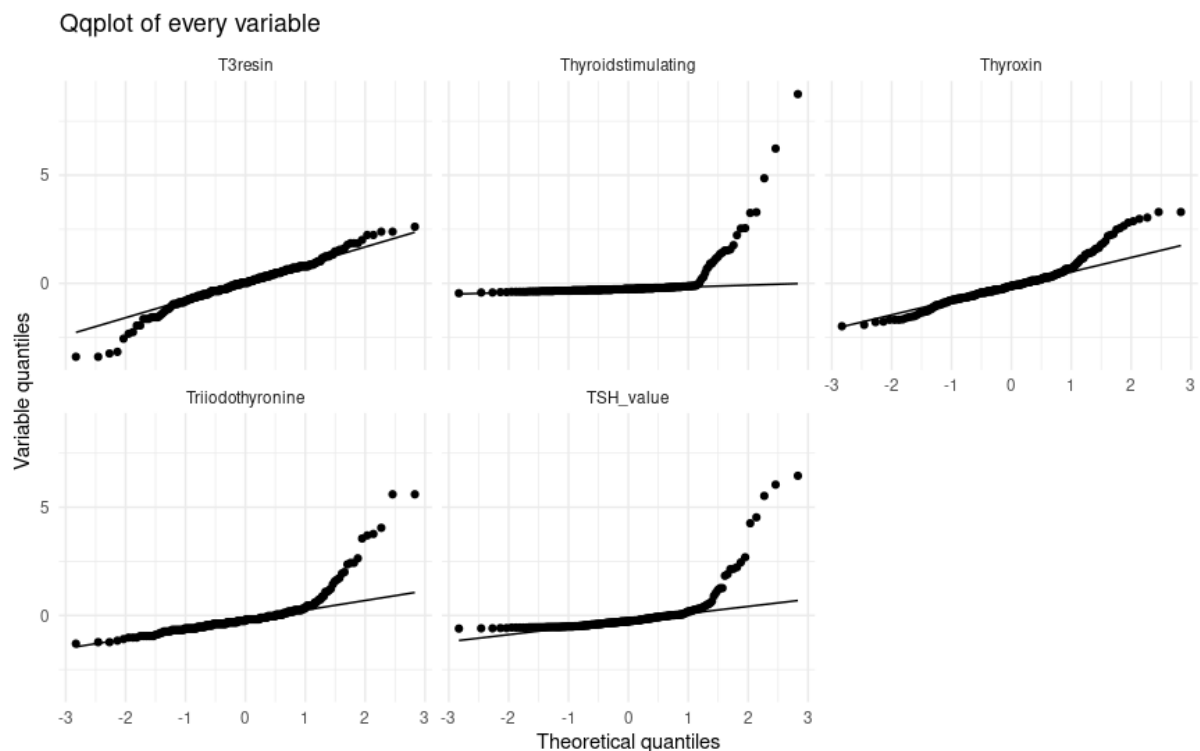


Figura 24: Qqplot for classification EDA.

En la figura 24 se puede ver más claro como las colas de las distribuciones son bastante grandes. En la mayoría, se presenta una asimetría a la derecha.

Si se les hace el test de *Agostino* (para la asimetría) se rechaza y por tanto no son normales y tienen alto grado de asimetría.

3.3. Outliers

En la sección de regresión se analizó tanto anomalías univariantes (en una sola característica) como variables multivariante. Se analizarán solo los posibles *outliers* univariantes, ya que si las variables por separado son uniformes, no lo serán en conjunto y por tanto métodos como *Mahalanobis* no serían tan efectivos (aunque podría usarse *LOF*, pero se excede de lo que ha de realizarse en esta asignatura). La normalidad multivariante implica la normalidad de las distribuciones marginales, por ejemplo, de las distribuciones univariantes de sus componentes [1].

En la figura de *boxplots* 25 se puede observar como hay variables (sobre todo *Thyroidstimulating*) que tienen valores muy alejados de la media. Aplicando la distancia intercuartil se obtiene que el 8 % de los datos en el conjunto son anómalos. El problema con la distancia intercuartil es que no es adecuada en distribuciones muy asimétricas, como es nuestro caso. Además, el estudiante no dispone de conocimiento experto para

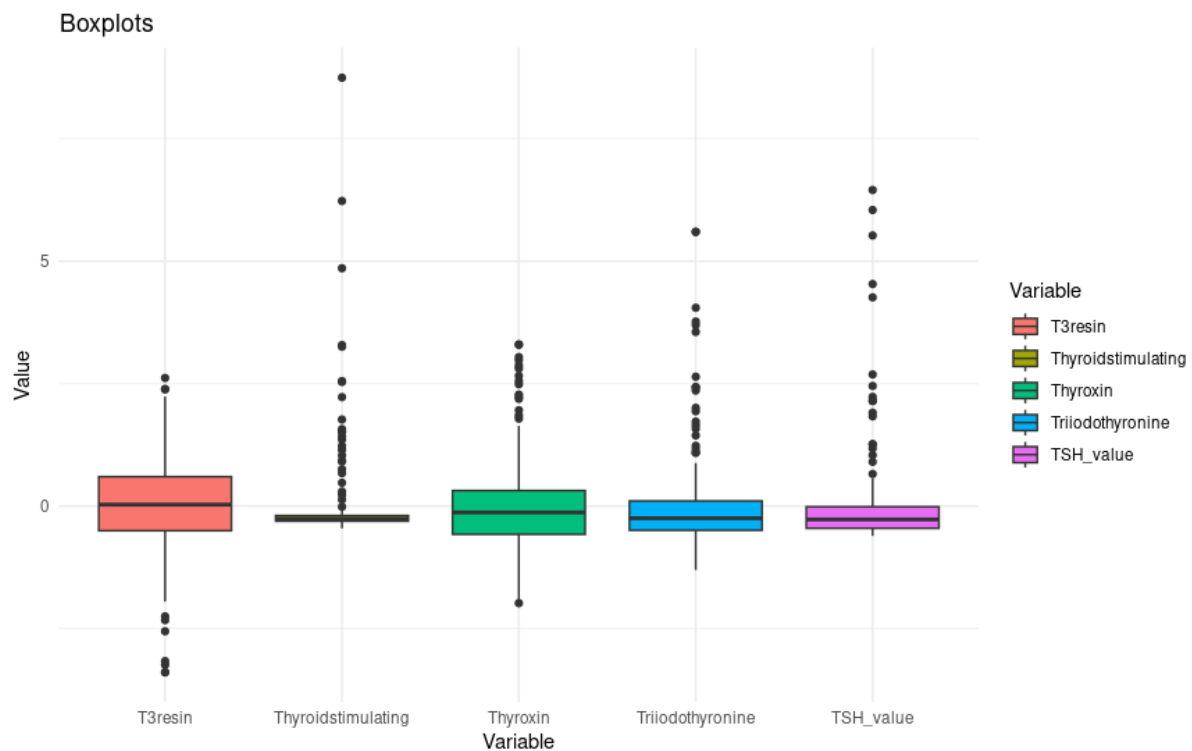


Figura 25: Boxplot for newthyroid.

valorar si las observaciones clasificadas como anómalas por este método son en realidad *outliers*. Por estas razones se decide dejar los valores.

3.4. Distribución por clases

Anteriormente se expuso que el conjunto de datos sufría de un desbalance de clases, por ello es interesante observar como se distribuyen las observaciones pertenecientes a cada clase en cada característica.

Dadas las gráficas de densidad por clase en 26,27,28,29,30 se observa como para las variables *T3resin*, *Thyroxin* y *Triiodothyronine* se suele clasificar en la clase 2 para valores más bajos y algo más dispersos. Valores altos de estas variables se relacionan más con las clases 1 y 2. Las distribuciones de clase para estas variables son muy parecidas, en cambio, en la variable *Thyroidstimulating* se puede percibir que las distribuciones son más discriminatorias para las clases 1 y 2, esto es interesante de cara a saber qué podría afectar a que se diese una clasificación u otra, ya que normalmente se observa que las clases 1 y 2 o crecen juntas o decrecen juntas mientras la clase 3 hace lo contrario. En esta variable se da un comportamiento nuevo.

Otra forma de observar la contribución de las variables a la clasificación es con un **Biplot** 31. En este se observa mejor la contribución de las variables entre sí y no por

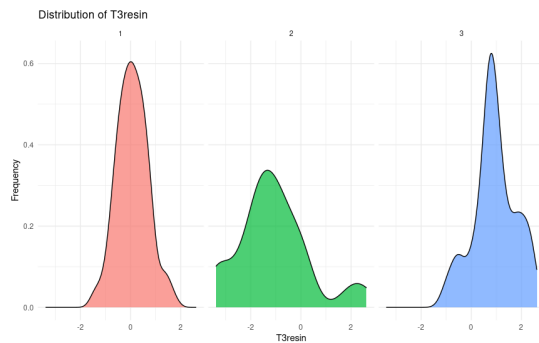


Figura 26: Class distribution for T3resin.

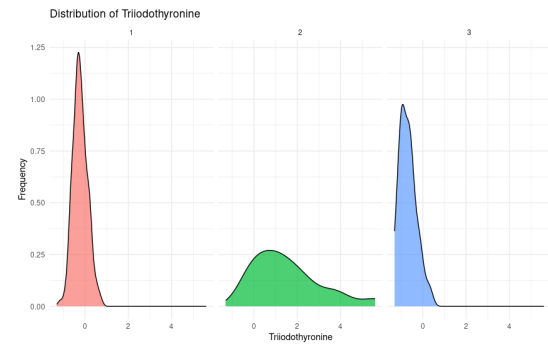


Figura 27: Class distribution for Triiodothyronine.

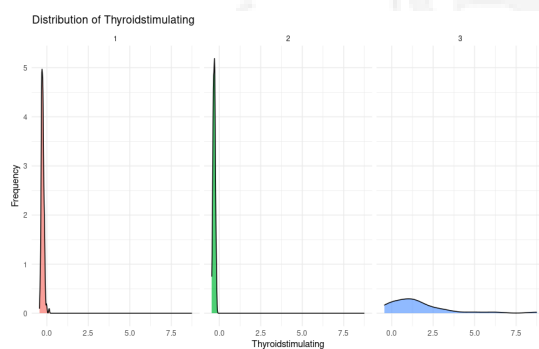


Figura 28: Class distribution for Thyroidstimulating.

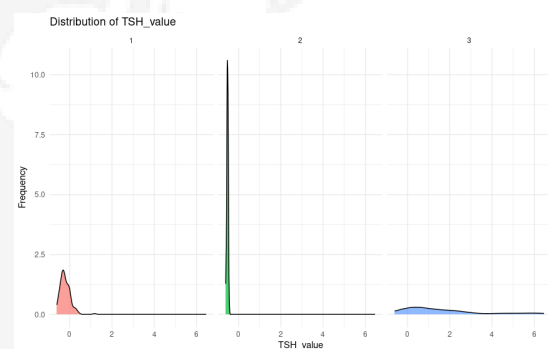


Figura 29: Class distribution for TSH_{value}.

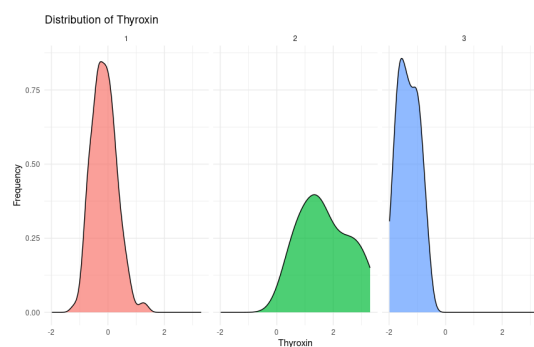


Figura 30: Class distribution for Thyroxin.

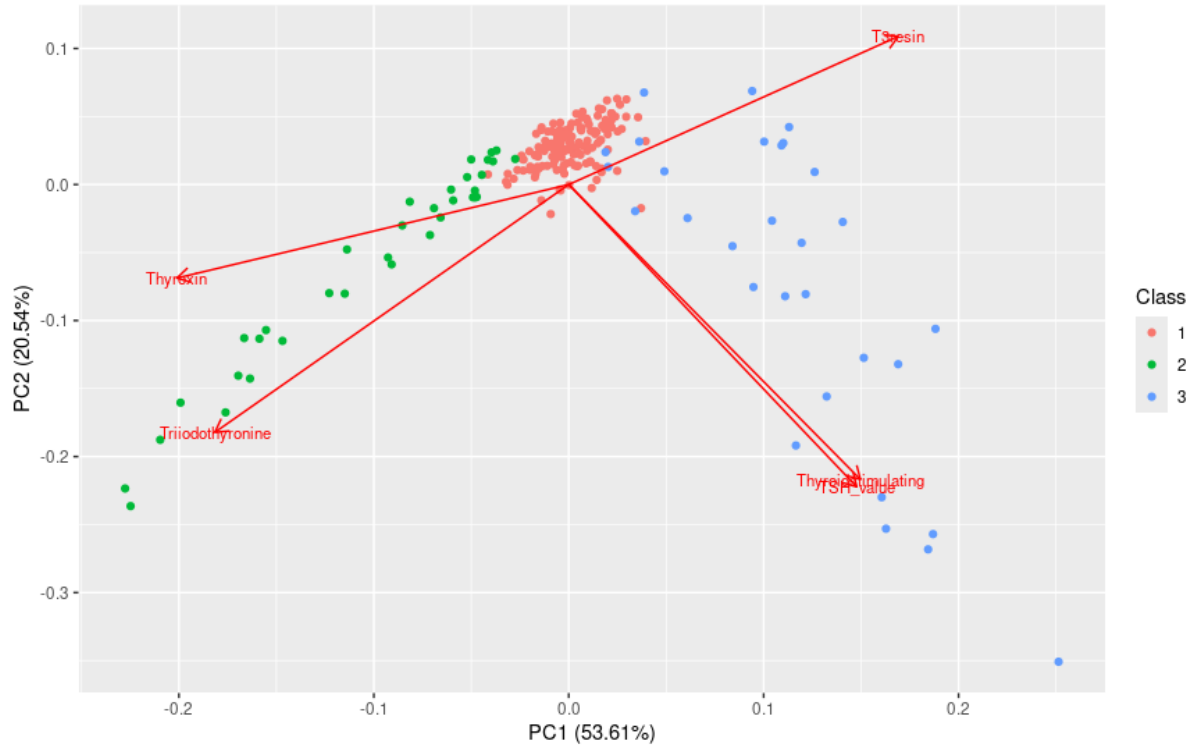


Figura 31: Biplot for classification EDA.

separado, como se ha analizado anteriormente. Combinaciones de valores altos en *Thyroxin* y *Triiodothyronine* parecen ser fuertes claves para que la observación sea clasificada como clase 2, por ejemplo.

3.5. Correlación

En este apartado se van a analizar las posibles relaciones entre variables del conjunto de datos. Para esto se pueden analizar gráficos visuales de puntos de variable *vs* variable y tablas de correlación.

Como puede observarse en el gráfico 32, existen algunas relaciones ligeramente lineales entre algunas variables, como es el caso de las variables *Thyroxin* y *Triiodothyronine*. El resto son menos obvias y no muy lineales. Si se obtiene la gráfica de la tabla de correlación, fig 16, el resultado que era obvio visualmente queda confirmado por un valor de correlación de 0,72.

Dados estos resultados, no se elimina ninguna característica por redundancia. Además, la dimensión del problema no es tal como para tener que optimizar el tamaño de las componentes/características que se tienen, pues hay muy pocas. De hecho muchas veces, dos características muy correladas no implican redundancia. Un ejemplo de esto aparece en [2] donde dos características con un ratio de correlación alto son totalmente necesarias

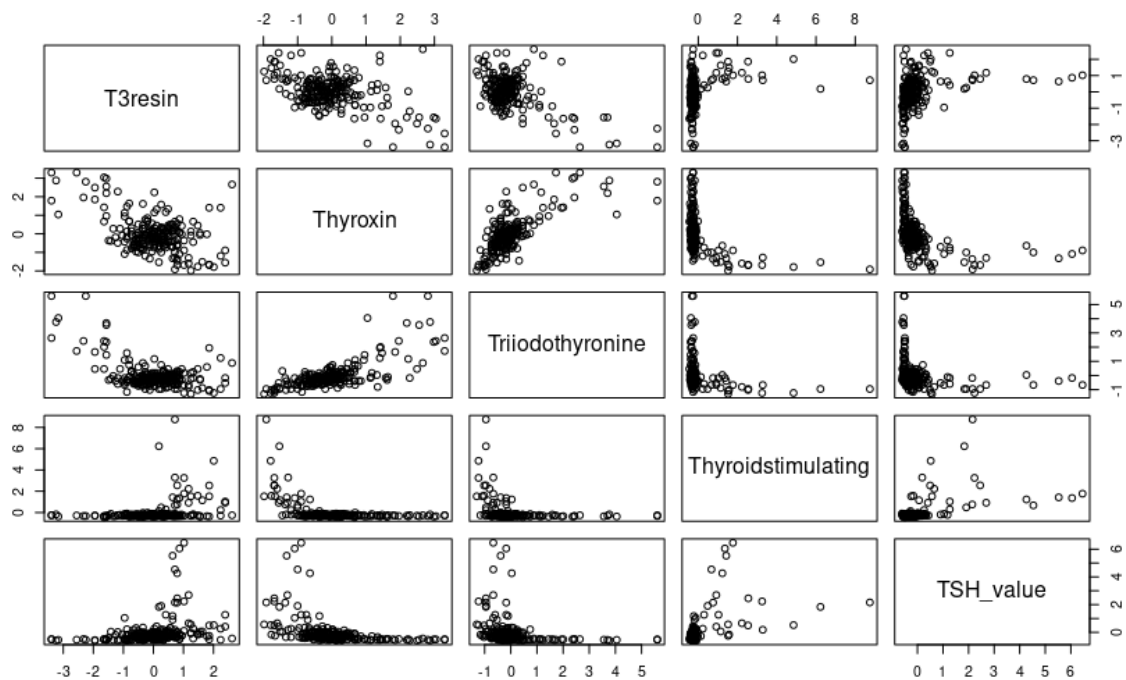


Figura 32: Pairs plot EDA classification.

(ambas) para la clasificación del problema.

3.6. Transformaciones

Si bien la normalidad es complicada alcanzarla en las variables que hay en el conjunto de datos (debido a grandes asimetrías) es posible intentar corregir algunas. Al aplicar la transformación logarítmica, los valores extremadamente altos se “comprimen”, reduciendo el impacto de los *outliers* positivos y acercando la distribución de los datos a una forma más simétrica, pero no garantiza la normalidad. Además, en clasificación, las clases pueden volverse más separables en el espacio de características.

Se ha optado por transformar tres variables, aquellas más alejadas a las normales, como prueba para ver si se reduce la asimetría. Aplicando los test correspondientes se observa una reducción de la asimetría muy notable, aunque eso no hace que las distribuciones de las variables sean normales.

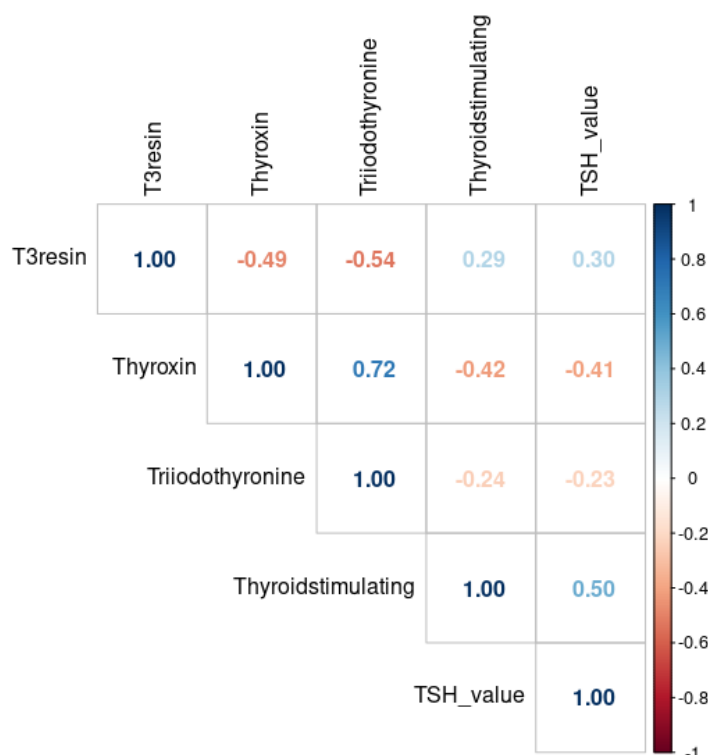


Figura 33: Correlation plot EDA classification.

4. Regresión

Se procede a describir el trabajo realizado en la parte de regresión con el conjunto de datos *wankara*.

4.1. Regresión simple

En este apartado se describirá el proceso que se ha seguido para obtener 5 regresores y cuál ha sido seleccionado como el mejor.

En el apartado del *EDA* de regresión se explicó que este conjunto de datos tiene 9 características, por lo que se debe elegir entre el total los que parezcan más prometedores. Según lo visto en el *EDA*, en lo relativo a correlación y al propio sentido común, la variable más prometedora es *diff_temperature* (es la combinación lineal de *max_temperature* y *min_temperature*). Seguido de *dewpoint*, *diff_pressure* (combinación lineal de variables originales de presión), *visibility* y *wind_speed*.

Se ha tenido en cuenta las variables que parecen más lineales con respecto a la variable objetivo, o que al menos parezcan tener una relación en algún sentido (por ejemplo, *visibility* es algo curva, algo logarítmica).

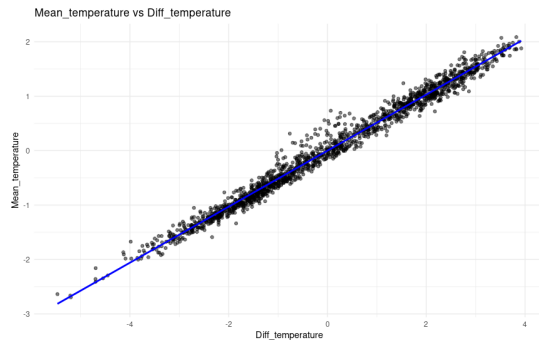


Figura 34: Diff temperature vs Mean Temperature.

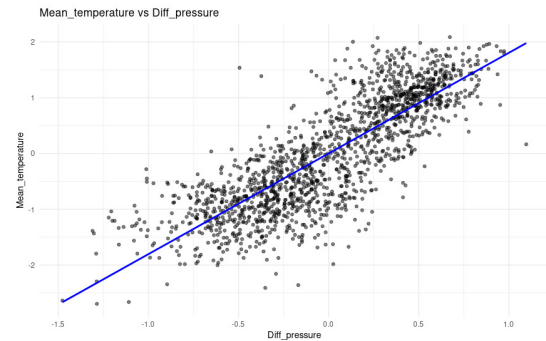


Figura 35: Diff pressure vs Mean Temperature.

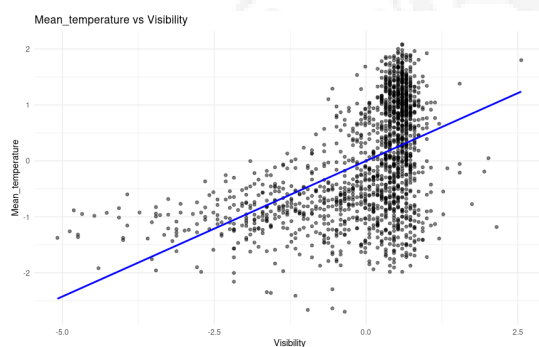


Figura 36: Visibility vs Mean Temperature.

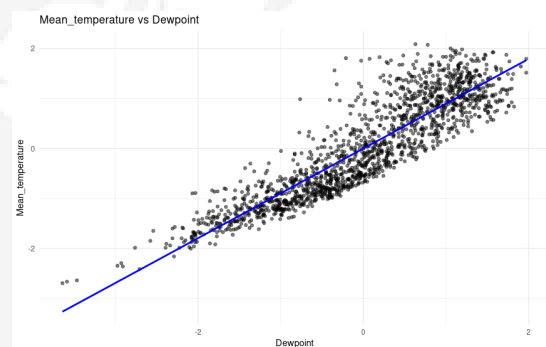


Figura 37: Dewpoint vs Mean Temperature

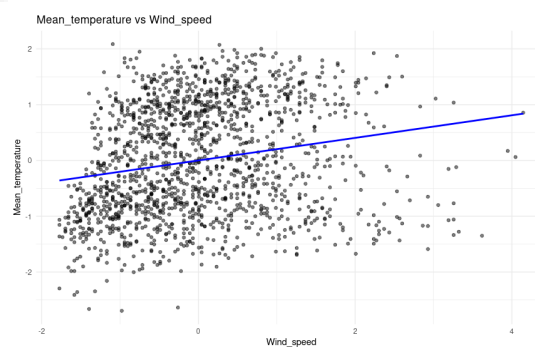


Figura 38: Wind speed vs Mean Temperature.

Cuadro 4: Linear Simple Regression Results

	Residual Std. Error	t value	p -value	R^2
Diff_temperature	0.132	301.2	$< 2 \times 10^{-16}$	0.9826
Dewpoint	0.4433	81.09	$< 2 \times 10^{-16}$	0.8036
Diff_pressure	0.5656	58.47	$< 2 \times 10^{-16}$	0.6803
Visibility	0.8749	22.22	$< 2 \times 10^{-16}$	0.235
Wind_speed	0.9796	8.29	$2,37 \times 10^{-16}$	0.0410

Cuadro 5: Summary of simple linear regression.

Como puede observarse en la tabla 5 la mejor variable es sin duda alguna *diff_temperature*, obtiene un error mucho mejor y el valor R^2 es casi de 1, lo que indica que prácticamente solo con esa variable se explica la varianza total de la variable objetivo. Se añaden gráficas de cada variable por separado y como estas se relacionan con la variable objetivo (figuras 34,35,36,37,38).

4.2. Regresión múltiple

En esta apartado se consideran múltiples variables así como las interacciones entre ellas y posibles no linealidades para la mejora del modelo.

4.2.1. Backward selection

Inicialmente se utiliza la técnica de selección de características *backward*. Esta consiste en un enfoque iterativo que comienza con el conjunto completo de características disponibles y elimina progresivamente aquellas que son menos relevantes para el modelo (por medio del p -valor), hasta que se alcanza un subconjunto “óptimo” de características.

Realizando este proceso se eliminan las variables *max_wind_speed* y *precipitation*. El resto de variables obtienen valores de p -valor por debajo de 0,05, por lo que con un 95 % de confianza se podría decir que están aportando suficiente al problema.

4.3. Interacciones y no linealidad

Se realizan pruebas añadiendo variables. Si bien es cierto que en la selección hacia atrás solo se eliminaron dos variables, como se ha analizado anteriormente solo *diff_temperature* parece ser suficiente, por ello se prueba primero esa variable aislada y luego se añaden más.

```

1 fit1 <- data %>%
2   lm(formula = Mean_temperature ~ Diff_temperature, data = .)
3

```

```
4 summary(fit1)
5
6
7 fit2 <- data %>%
8   lm(formula = Mean_temperature ~ Diff_temperature + Dewpoint + Diff_
9     pressure,
10     data = .)
11 summary(fit2)
12
13 fit3 <- data %>%
14   lm(
15     formula = Mean_temperature ~ Diff_temperature + Dewpoint + Diff_
16     pressure + Visibility,
17     data = .
18   )
19 summary(fit3)
```

El resultado de ello es que no se mejora apenas. Se prueban interacciones entre variables con sentido, como puede ser *dewpoint* o punto de condensación y la variable de la temperature. Se consideran también interacciones no lineales entre estas dos variables.

```
1 fit4 <- data %>%
2   lm(formula = Mean_temperature ~ Diff_temperature * Dewpoint,
3     data = .)
4
5 summary(fit4)
6
7 fit5 <- data %>%
8   lm(formula = Mean_temperature ~ Diff_temperature * I(Dewpoint ^ 2),
9     data = .)
10
11 summary(fit5)
12
13 fit6 <- data %>%
14   lm(formula = Mean_temperature ~ Diff_temperature * I(Dewpoint ^ 2) *
15     Dewpoint,
16     data = .)
17 summary(fit6)
18
19 fit7 <- data %>%
20   lm(
21     formula = Mean_temperature ~ Diff_temperature * I(Dewpoint ^ 2) *
22     Dewpoint * I(Diff_temperature ^
23
24     2),
25     data = .
26   )
27 summary(fit7)
```

Las métricas mejoran, se pasa de un error estándar de 0,13 a 0,11. Aunque es poco, es el mejor resultado hasta el momento.

Se han probado algunas combinaciones más, pero no tienen demasiado sentido (a priori) y no parecen mejorar en absoluto el modelo. También hay que tener en cuenta que un ajuste demasiado fino puede resultar en un **sobreajuste**.

4.4. KNN

Se realizan pruebas igual que con el algoritmo de regresión lineal, pero usando *KNN*. En este caso se obtienen resultados (a priori) mejores que en regresión lineal.

4.5. Comparativas

En esta sección se realizan comparativas entre los algoritmos de regresión lineal múltiple, *KNN* y *M5*. Para ello se implementa una función de *k-fold cross validation* con $k = 5$. Para una comparación justa entre algoritmos se ajustan los modelos con todas las variables.

Algorithm	Mean accuracy
Multiple Linear Regression	0.01034
KNN	0.01777
M5	0.02576

Cuadro 6: Summary of results in k-fold validation.

Como puede observarse en la tabla 6, el algoritmo de regresión lineal ha obtenido mejores resultados que los otros dos. Esto se debe a que la naturaleza de este problema es lineal y por tanto el algoritmo que mejor se ajusta a esto es el de regresión lineal. La temperatura media tiene relaciones lineales con algunas variables como son temperatura máxima, mínima y punto de condensación. Existen otras, pero las variables más importantes, como se analizó anteriormente, son las mencionadas, y sus correlaciones con la variable objetivo son casi perfectas.

Se muestran gráficas de las predicciones en cada *fold* de los diferentes algoritmos en 39,40,41. Los valores que más se acercan son los del modelo de regresión lineal, pero los otros algoritmos obtienen un ajuste muy bueno igualmente. El algoritmo de regresión lineal obtiene mejor varianza con respecto a las observaciones actuales y las observaciones predichas.

Para las comparaciones entre algoritmos con test estadísticos es necesario realizar ciertas normalizaciones, ya que en regresión los errores no están en las mismas escalas. Se usa la fórmula 1.

Para comparar los algoritmos, se realizan los test de *Wilcoxon*. El test de *Wilcoxon* es una

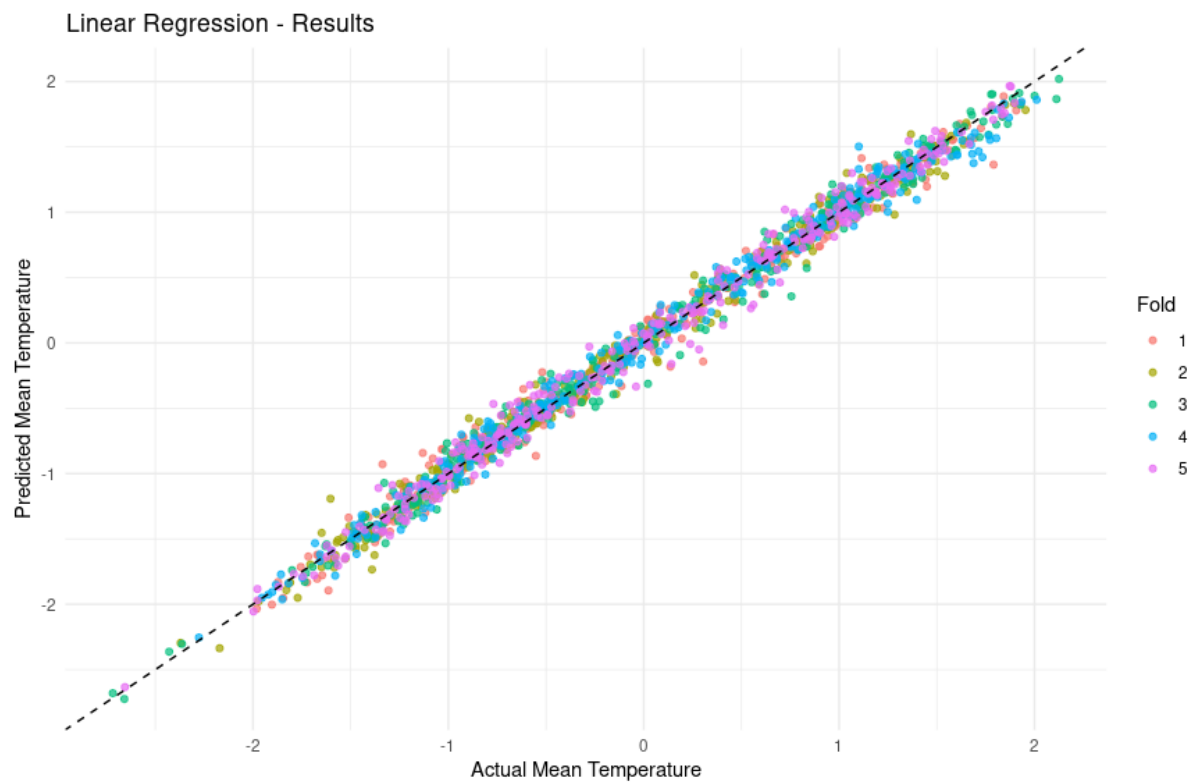
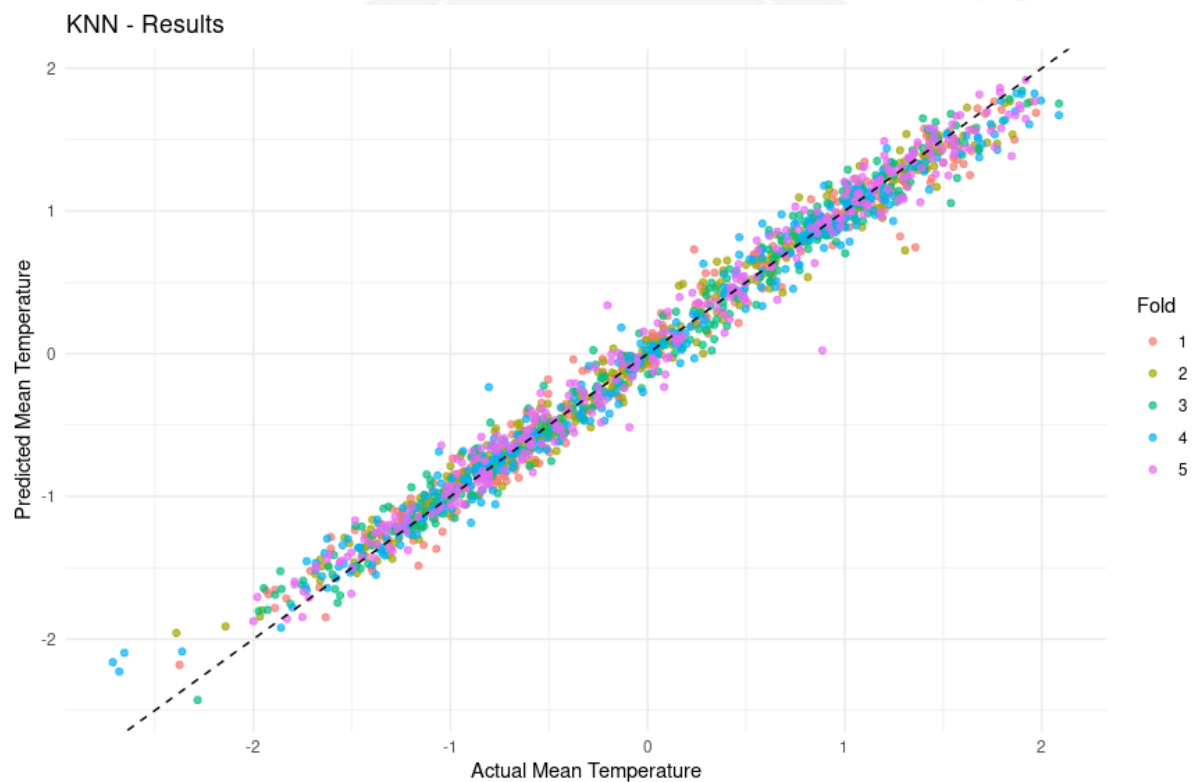


Figura 39: K-fold cross validation results in linear regression.

Figura 40: K-fold cross validation results in **KNN**.

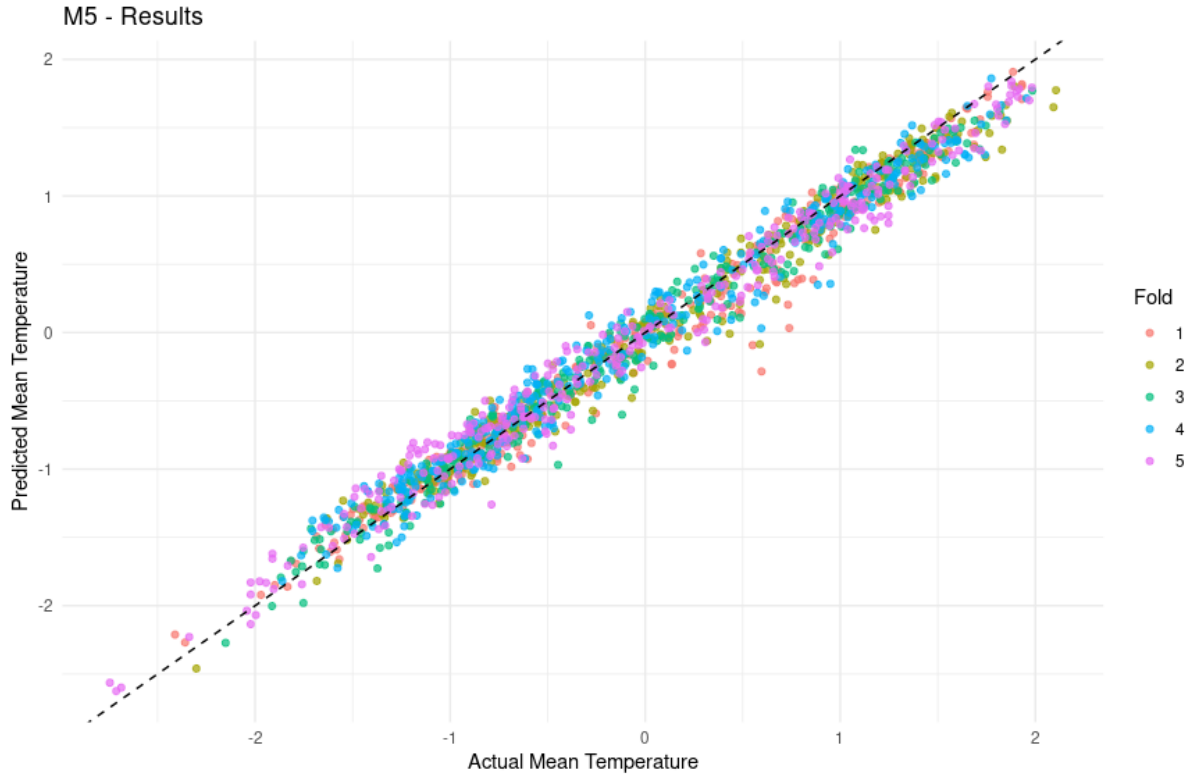


Figura 41: K-fold cross validation results in **M5**.

prueba estadística no paramétrica utilizada para comparar dos conjuntos de datos emparejados y determinar si sus distribuciones difieren significativamente. Se utiliza cuando no se puede asumir que los datos siguen una distribución normal.

$$\text{DIFF} = \frac{\text{Mean}(\text{Other}) - \text{Mean}(\text{Reference Algorithm})}{\text{Mean}(\text{Other})} \quad (1)$$

Los resultados en los test son los obtenidos en 2. Como puede verse, comparando los algoritmos de **KNN** vs los de regresión lineal múltiple, no se obtienen resultados estadísticamente significativos. Por ello, no se puede decir que un algoritmo sea mejor que otro, ya que las mejoras podrían haberse dado por pura aleatoriedad. No hay suficiente evidencia.

$$\text{p-value} = 0,701, \quad R+ = 76, \quad R- = 95 \quad (2)$$

Para comparar los tres algoritmos es necesario hacer el test de **Friedman** y realizar una corrección *post-hoc*, ya que solo con **Friedman** se acumulan errores que pueden invalidar las conclusiones. Se utiliza **Friedman** ya que es una alternativa no paramétrica al **ANOVA** de medidas repetidas, pues no asume normalidad en los datos.

Tras la corrección se obtienen los resultados en la tabla 7.

	1	2
2	0.44	-
3	0.16	0.28

Cuadro 7: Friedman test + Holms (1: Linear Regression, 2: KNN, 3: M5) en *test*.

Tampoco se obtienen resultados significativos.

Ahora se realizan las mismas comparaciones para el conjunto de *training*.

$$\text{p-value} = 0,00025, \quad R_+ = 9, \quad R_- = 162 \quad (3)$$

	1	2
2	0.00057	-
3	0.00475	0.00107

Cuadro 8: Friedman test + Holms (1: Linear Regression, 2: KNN, 3: M5) en *training*.

En las tablas 3, 8 se obtienen resultados estadísticamente significativos. De hecho si se comparan las medianas, se obtienen resultados a favor de los algoritmos **M5** y **KNN** frente a regresión lineal. Concretamente, la mediana para **M5** en *train* es de 2,8, la de **KNN** es de 1,8 y regresión lineal de 5. Por ello, las diferencias son significativas para cada comparación entre pares y concretamente, **KNN** sale parado como el mejor (en *train*). En *test* sin embargo, obtiene mejor mediana **M5**, pero los test no rechazan.

El hecho de que los test rechacen en el conjunto de entrenamiento y no en el de evaluación puede significar un posible sobreajuste. Si los resultados no son significativos en *test*, significa que los modelos no logran reproducir esos mismos resultados con los datos no vistos.

5. Clasificación

Se procede a describir el trabajo realizado en la parte de clasificación con el conjunto de datos *newthyroid*.

5.1. KNN

Se realiza un pequeño estudio para comprobar que valor para k es mejor en este problema concreto. Para ello se han aplicado las transformaciones necesarias definidas en el **EDA** de clasificación y se ha definido una función de particionamiento en conjuntos de datos de entrenamiento y de evaluación. Se ha omitido transformar la variable de *Triiodothyronine* con el logaritmo ya que esta transformación produce muchos nulos (hay

valores no definidos para el logaritmo dentro de la variable). El resto de transformaciones logarítmicas se realizan para reducir asimetrías.

K	Accuracy
1	0.95
5	0.95
15	0.93
40	0.81
100	0.7

Cuadro 9: K values and accuracies for **KNN**.

Como puede observarse en la tabla 9, se obtienen muy buenos resultados para valores de k pequeños, incluso para $k = 1$, que es un valor muy dado al sobreajuste debido a que hace al modelo mucho más sensible a cambios cuando se añaden nuevos datos, añade mucha varianza. Para valores más seguros, como $k = 5$ los resultados siguen siendo muy buenos. Para valores superiores a 15, el modelo tiende a equivocarse cada vez más, está aumentando el sesgo del modelo y por tanto generaliza cada vez peor.

Se añaden las gráficas de matrices de confusión (fig 42,43,44,45,46) donde puede verse que los primeros modelos con k bajos ajustan a la perfección las clases predichas y mientras este valor de k aumenta, suele clasificar erróneamente con la clase 1. Según se explicó en el **EDA**, había un gran desbalance con la clase 1 frente a las otras dos, por lo que es normal que según el modelo sea peor, este tienda a clasificar con la clase más mayoritaria.

Se escoge finalmente el modelo con $k = 5$ debido a que tiene una métrica de precisión igual y se reduce la complejidad del modelo.

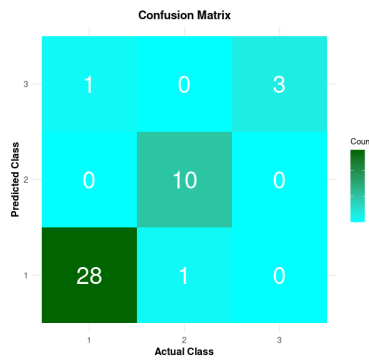
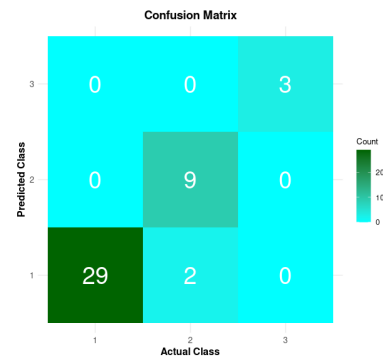
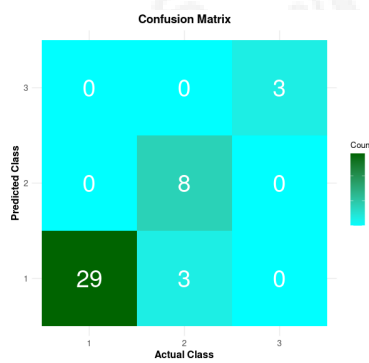
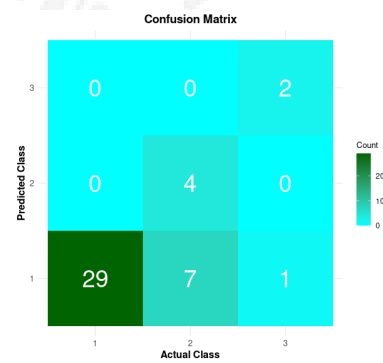
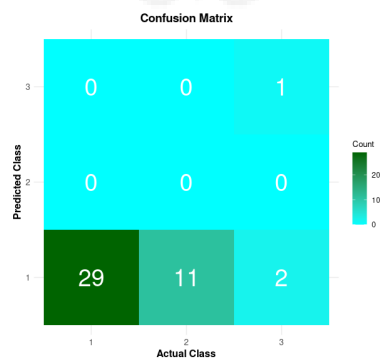
5.2. LDA

Se procede a realizar un ajuste de los datos con el algoritmo **LDA** (*Linear Discriminant Analysis*). **LDA** es una técnica supervisada utilizada principalmente para reducción de dimensionalidad y clasificación. Funciona buscando maximizar la separabilidad entre clases proyectando los datos en un espacio de menor dimensión. Este algoritmo asume ciertas propiedades sobre los datos, por ello, es necesario realizar una serie de comprobaciones sobre los datos.

5.2.1. Asunciones

El algoritmo **LDA** asume que:

- Cada característica del conjunto de datos es normal para cada clase, es decir, los datos dentro de cada clase forman una campana simétrica alrededor de la media.

Figura 42: confusion Matrix for $k = 1$.Figura 43: confusion Matrix for $k = 5$.Figura 44: confusion Matrix for $k = 15$.Figura 45: confusion Matrix for $k = 40$.Figura 46: confusion Matrix for $k = 100$.

- Se asume que la matriz de covarianza de las características es la misma para todas las clases. Esto implica que las clases tienen la misma dispersión y varianza en todas las dimensiones del espacio de características.
- Los datos obtenidos para la creación del *dataset* se han obtenido como una muestra aleatoria.

Para comprobar la primera asunción se realiza una serie de test sobre cada variable agrupada por clase, concretamente el test de *Shapiro*. Se utiliza la siguiente función, la cual también crea una gráfica *QQPlot*.

```

1 test_variable <- function(var_name) {
2   resultados <- data %>%
3     group_by(Class) %>%
4     summarize(
5       p_value = if (is.numeric(.data[[var_name]]) &&
6                     all(!is.na(.data[[var_name]]))) {
7         shapiro.test(.data[[var_name]])$p.value
8       } else {
9         NA
10      },
11      normal = ifelse(p_value > 0.05, "Puede ser normal", "No es normal
12      "),
13      .groups = "drop"
14    )
15 plot <- ggplot(data, aes(sample = .data[[var_name]])) +
16   stat_qq() + stat_qq_line() +
17   facet_wrap(~ Class)
18 print(plot)
19 print(resultados)
20 }
```

Como muestra la tabla 10, algunas variables rechazan el test, por lo que no pueden ser consideradas normales con una garantía estadística del 95 %.

Para comprobar la igualdad de matrices de covarianza-varianza se realizan test de *Bartlett* para cada variable. Los resultados son que todas las variables rechazan con una diferencia estadística muy significativa, lo que significa que las varianzas de los grupos no son iguales y, por lo tanto, no se cumple el supuesto de homocedasticidad (es decir, las varianzas iguales entre grupos).

5.2.2. Ajuste

Al ajustar el modelo con **LDA** se obtiene que la primera variable discriminante aporta un 83 % de la variable explicada, mientras que la segunda aporta un 17 %. Esto indica que la primera variable es la que consigue la mayor parte de la separación. Dentro de los datos recogidos se observa que la variable que más aporta a **LD1** es *TSH_value* con un coeficiente de 1,42.

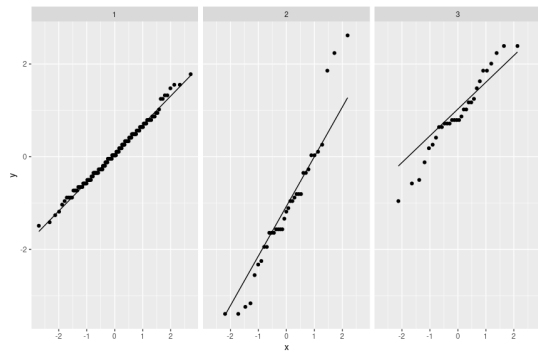


Figura 47: **QQPlot** for T3resin grouped by class.

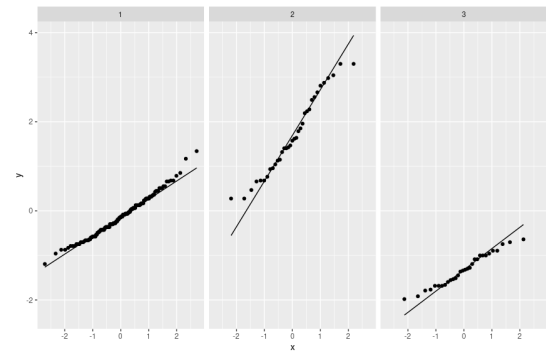


Figura 48: **QQPlot** for Thyroxin grouped by class.

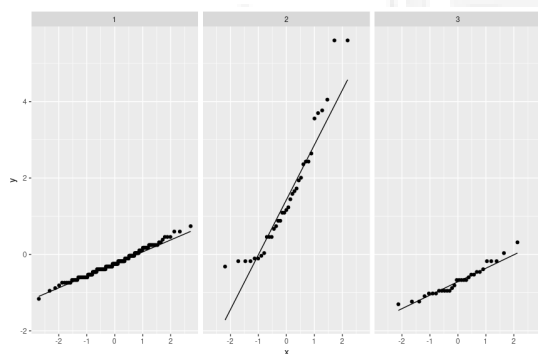


Figura 49: **QQPlot** for Triiodothyronine grouped by class.

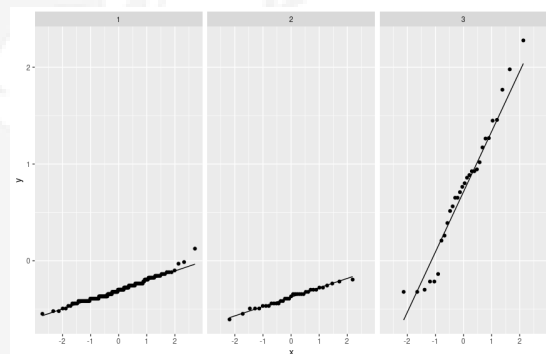


Figura 50: **QQPlot** for Thyroidstimulating grouped by class.

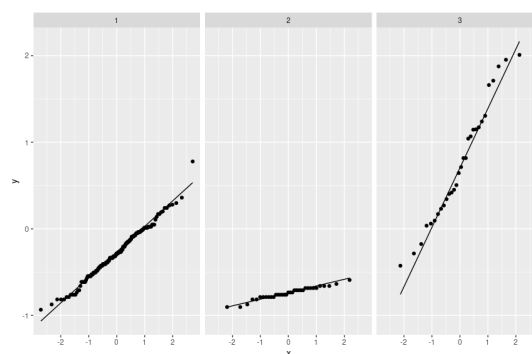
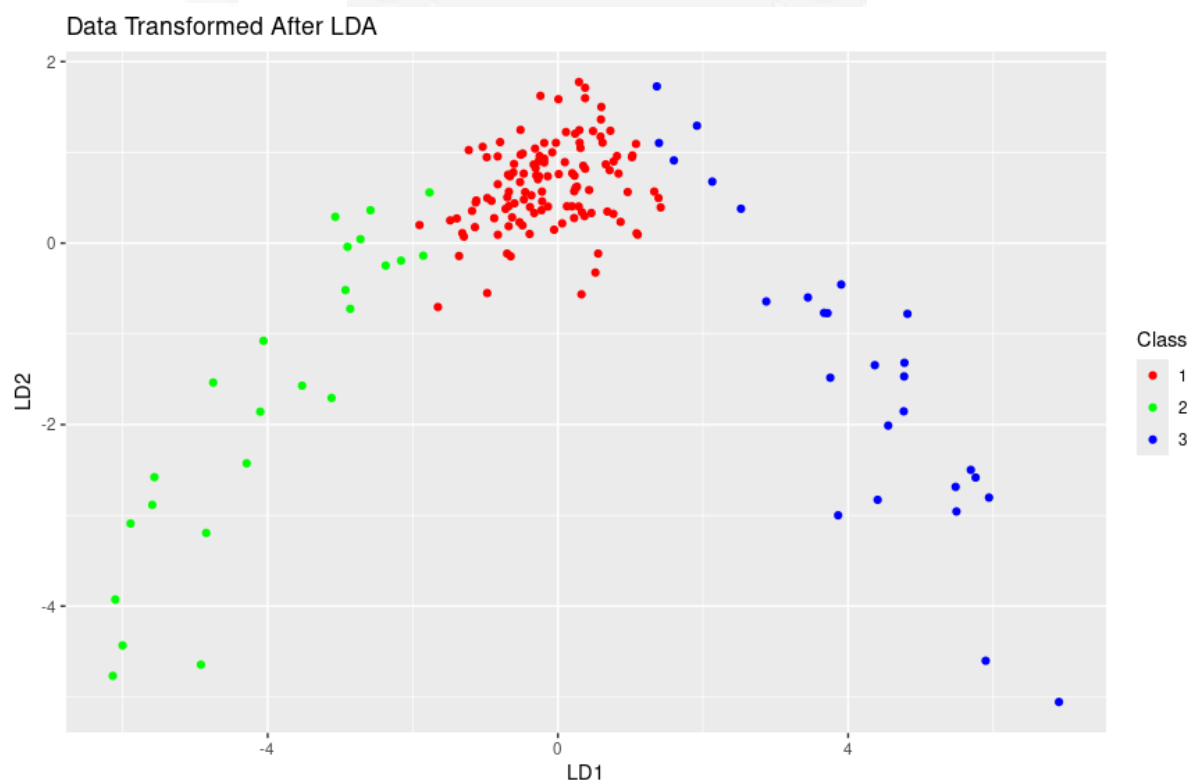


Figura 51: **QQPlot** for TSH_value grouped by class.

Variable	Class	p_value	normal
T3resin	1	0.736	Puede ser normal
	2	0.0479	No es normal
	3	0.414	Puede ser normal
Thyroxin	1	0.143	Puede ser normal
	2	0.200	Puede ser normal
	3	0.448	Puede ser normal
Triiodothyronine	1	0.257	Puede ser normal
	2	0.00357	No es normal
	3	0.189	Puede ser normal
Thyroidstimulating	1	0.00683	No es normal
	2	0.865	Puede ser normal
	3	0.352	Puede ser normal
TSH_value	1	0.251	Puede ser normal
	2	0.293	Puede ser normal
	3	0.405	Puede ser normal

Cuadro 10: Normality test for every variable grouped by class.

Figura 52: Data transformation after **LDA**.

En la figura 52 se puede ver como se transforman los datos de *train* tras aplicar el algoritmo. Se obtiene un *accuracy* del 90,6 %. Pese a no cumplir la segunda asunción y la primera en parte, el algoritmo es capaz de ajustar los datos muy bien.

5.3. QDA

A diferencia de **LDA**, **QDA** no asume que todas las clases compartan la misma matriz de covarianza. Cada clase puede tener su propia matriz de covarianza. Como consecuencia, **QDA** genera fronteras de decisión cuadráticas (curvas) entre los grupos. Esto permite que **QDA** sea más flexible, ya que puede modelar relaciones no lineales entre las variables.

A priori, se podría suponer que este algoritmo debería ir mejor que **LDA** ya que ninguna de las variables cumplía la segunda asunción. De hecho, al ajustar el modelo, se obtiene un 97 % de precisión. En la matriz de confusión (fig 53) puede observarse la frecuencia de clasificación por clases.

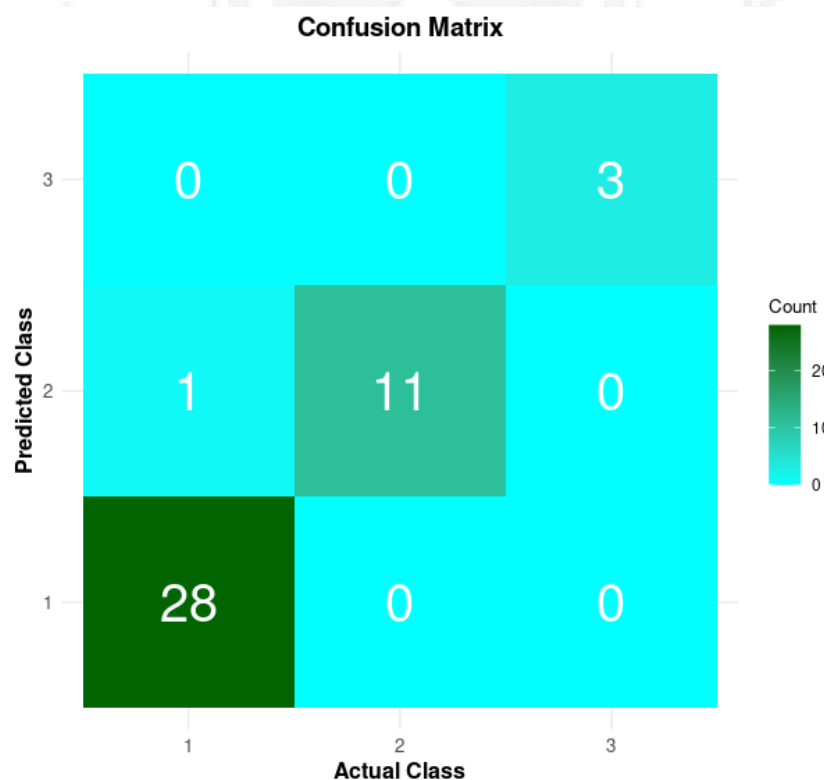


Figura 53: Confusion Matrix for **QDA**.

5.4. Comparación

Al realizar una comparación de los tres algoritmos usando *10-fold cross validation*, se obtienen los resultados descritos en la tabla 11. **KNN** obtiene unos resultados muy

buenos. **LDA** obtiene una precisión muy alta, pero se queda por detrás de los otros dos algoritmos, además es más variante en comparación a **KNN** y **QDA** (aunque sigue siendo un valor de desviación típica muy bajo). Es de esperar que este algoritmo fuese a ser el peor, pues sus asunciones básicas no eran cumplidas en algunos casos. Pese a ello los resultados siguen siendo muy buenos.

El algoritmo que mejores resultados a logrado para este conjunto de datos es **QDA**, con un asombroso 97 % de precisión y la desviación típica más baja.

Algorithm	Mean accuracy
KNN	$0,96 \pm 0,039$
LDA	$0,933 \pm 0,053$
QDA	$0,976 \pm 0,033$

Cuadro 11: Summary of results in k-fold validation.

6. Apéndice

6.1. EDA Regresión

```

1 library(tidyverse)
2 library(moments)
3 library("corrplot")
4 library(ggfortify)
5 library(MVN)
6
7 "Leemos los datos saltandonos las cabeceras iniciales y despu s las
   parseamos
8 a mano. Seguido, vamos a ver unos pocos datos con head para hacernos
   una primera
9 idea."
10
11 data <- read.csv("wankara/wankara.dat", skip = 14, header = FALSE)
12 colnames(data) <- c(
13   "Max_temperature",
14   "Min_temperature",
15   "Dewpoint",
16   "Precipitation",
17   "Sea_level_pressure",
18   "Standard_pressure",
19   "Visibility",
20   "Wind_speed",
21   "Max_wind_speed",
22   "Mean_temperature"
23 )
24 head(data)
25
26 "Con la funci n de summary obtenemos un resumen estad stico general
   de cada

```

```
27 columna. Adem s vamos a obtener otros datos informativos como la
    dimensi n,
28 y la estructura."
29
30 summary(data)
31 dim(data)
32 str(data)
33 colSums(is.na(data))
34
35 "Los datos parecen bastante buenos a priori. No hay escalas num ricas
    muy
36 grandes, no hay datos faltantes y en principio parece que las
    distribuciones
37 se centran en la media (a excepci n de algunas variables como
    Precipitation y
38 alguna m s). En principio lo m s intuitivo ser a analizar
    visualmente las
39 variables."
40
41 plot_distribution <- function(data,
42                               var_name,
43                               binwidth = 0.1,
44                               fill_color = "pink") {
45   ggplot(data, aes_string(x = var_name)) +
46     geom_density(fill = fill_color,
47                  color = "black",
48                  alpha = 0.7) +
49     labs(
50       title = paste("Distribution of", var_name),
51       x = var_name,
52       y = "Frequency"
53     ) +
54     theme_minimal()
55 }
56
57 for (column in names(data)) {
58   p <- plot_distribution(data, column)
59   print(p)
60 }
61
62 "Las variables de temperatura m xima y m nima parecen adecuarse a la
    forma de
63 una pseudo-bimodal (Tienen dos picos y una depresi n claramente
    diferenciados).
64 La variable de precipitaci n es asim trica extrema, con valores
    pr cticamente
65 nulos (entiendase por nulo el cero). Las relativas a la presi n son
    distribuciones que si bien no son normales exactas, se asimilan mucho (
    ambas
67 poseen la peculiaridad de un peque o \"bulto\" para presiones de entre
    26-30. En
68 cuanto al resto, hay variedad de distribuciones con asimetr as."
69
```



```
70 "La bimodalidad en las temperaturas puede darse por varios motivos,
    pero el que
71 considero m s probable es la toma de datos en distintas zonas
    geogr ficas de
72 Ankara, dando lugar a distintas modas. El calentamiento clim tico
    podr a ser
73 factible, pero no ser a tan evidente y se necesitar an tomas de
    muchos a os."
```

74

```
75 "Vamos a escalar los datos, ya que as se reduzcan las diferencias
    entre
76 variables. Escalar permite que todas las variables contribuyan
    equitativamente a
77 la distancia calculada, evitando que una variable desproporcionada
    influya en
78 la detecci n de outliers."
```

79

```
80 data <- as.data.frame(scale(data))
81 summary(data)
```

82

```
83 "A continuaci n vamos a aplicar test de normalidad a aquellas
    variables m s
84 prometedoras y ver si la cumplen. Adem s obtenemos los qqplots de
    todas las
85 variables para ver como de cerca est n las distribuciones de la normal
    te rica."
```

86

```
87 long_data <- data %>%
88   pivot_longer(
89     cols = setdiff(names(data), "Mean_temperature"),
90     names_to = "Variable",
91     values_to = "Value"
92   )
```

93

```
94 ggplot(long_data, aes(sample = Value)) +
95   stat_qq() +
96   stat_qq_line() +
97   facet_wrap(~ Variable) +
98   labs(title = "Qqplot of every variable", y = "Variable quantiles", x
    = "Theoretical quantiles") +
99   theme_minimal()
```

100

```
101 shapiro.test(data$Sea_level_pressure)
102 shapiro.test(data$Standard_pressure)
```

103

```
104 "Vamos a analizar los boxplots de las variables"
```

105

```
106 ggplot(long_data, aes(x = Variable, y = Value, fill = Variable)) +
107   geom_boxplot() +
108   labs(title = "Boxplots", x = "Variable", y = "Value") +
109   theme_minimal() +
110   theme(axis.text.x = element_text(angle = 45, hjust = 1))
```

111

```
112 "Observando las variables y sus valores máximos y mínimos no
    considerar a que
113 existan outliers univariantes, ya que son valores posibles y nada
    irreales.
114 Las precipitaciones es la que tiene valores más extremos, pero por lo
    general,
115 al ser una zona seca, es normal que la mayoría de valores sean cero.
    Más tarde
116 comprobar si pueden existir outliers multivariantes."
117
118 "Si bien es cierto que el test rechaza (no es una distribución normal)
    . Los
119 qqplot de las variables concernientes a la presión son muy buenos. De
    hecho en
120 las gráficas de densidad hemos observado que no presentan formas
    tan malas. Dados
121 estos resultados, pese a no pasar los test, podremos tomar una
    postura un poco
122 más relajada en cuanto a restricciones de la normal en lo que se
    refiere a esas dos
123 variables, sobre todo en algunas técnicas (si se llegan a usar) como
    ANOVA, que
124 son muy robustas a violaciones de normalidad.
125 El resto de variables son también bastante normales (algo asimétricas
    ).
126 Se dan excepciones, por ejemplo, La visibilidad tiene una cola muy
    larga a la
127 derecha, al igual que la máxima velocidad del viento (aunque ya hemos
    visto que
128 las variables de viento parecen bimodales)."
129
130 "Vamos a ver la correlación de las características"
131
132 cor_matrix <- cor(data)
133 corrrplot(cor_matrix,
134           method = "number",
135           type = "upper",
136           tl.col = "black")
137
138 "Existen fuertes coeficientes de correlación entre variables de
    mínimo, máximo y
139 temperatura media. Estas correlaciones son muy grandes y positivas,
    cuando crece
140 una variable lo hace la variable objetivo. Además el punto de
    condensación o
141 dewpoint también tiene correlación con la variable objetivo. Estas
    correlaciones
142 tienen sentido. El hecho de que el punto de condensación o dewpoint
    también est
143 correlacionado con la variable objetivo refuerza la idea de que las
    condiciones
144 atmosféricas, como la humedad y la temperatura, están relacionadas
    con el
```

```
145 comportamiento de la variable de inter s. En particular, el dewpoint
    es un
146 indicador importante de la humedad en el aire, lo que podr a influir
147 directamente en el fen meno que se est modelando."
148
149 ggplot(data, aes(x = Dewpoint, y = Mean_temperature)) +
150   geom_point()
151
152 "Las variables de presi n a nivel del agua tambi n tienen
    correlaci n con las
153 temperaturas. La presi n atmosf rica disminuye con el aumento de la
    elevaci n,
154 por lo que podr a intuirse que la temperatura media disminuye en
    lugares con
155 poca altitud en este conjunto de datos."
156
157 ggplot(data, aes(x = Mean_temperature, y = Sea_level_pressure)) +
158   geom_point()
159
160 "Vamos a crear dos nuevas variables que resuman la informaci n de max_
    temperature,
161 min_temperature y las variables relativas a presi n."
162
163 data <- data %>%
164   mutate(
165     Diff_temperature = Max_temperature + Min_temperature,
166     Diff_pressure = Standard_pressure - Sea_level_pressure
167   ) %>%
168   select(-Min_temperature, -Max_temperature, -Sea_level_pressure, -
    Standard_pressure)
169
170 "Vamos a realizar un an lisis de comprobaci n de outliers
    multivariantes, es
171 decir, en combinaci n con m ltiples variables. Primero, comprobamos
172 la normalidad con un test multivariante."
173
174 mvn(data = data, mvnTest = "hz")
175
176 "No lo pasa, por lo que los datos no siguen una normal con un alto
    grado de
177 confianza en altas dimensiones. Vamos a usar Mahalanobis igualmente
    para ver
178 la proporci n de potenciales outliers"
179
180 mahal_dist <- mahalanobis(data, colMeans(data), cov(data))
181 threshold <- qchisq(0.975, df = dim(data)[2]) # 97.5% confidence level
182
183 data <- as.data.frame(data) %>% mutate(Outlier = mahal_dist > threshold
    )
184
185 # % of potential outliers
186 (nrow(data %>% filter(Outlier)) / nrow(data) * 100)
187
188 ggplot(as.data.frame(mahal_dist), aes(x = mahal_dist)) +
```

```

189 geom_histogram(
190   binwidth = 0.5,
191   fill = "lightblue",
192   color = "black",
193   alpha = 0.7
194 ) +
195 labs(title = "Mahalanobis distances distribution", x = "Mahalanobis
196   Dist", y = "Freq") +
197 theme_minimal()
198 "Tal y como se observa en el gráfico, la distribución de las
199 distancias tampoco
200 sigue una normal. Es un indicativo más de que los datos no son
201 normales. Por ello
202 no voy a eliminar las variables encontradas por este método."
203 "Por último vamos a usar PCA para obtener un resumen de los datos y
204 gráficos que
205 poder analizar más que por reducir dimensionalidad, ya que solo
206 tenemos 10
207 variables."
208 pca_res <- prcomp(data)
209 pca_df <- data.frame(pca_res$x)
210 ggplot(pca_df, aes(x = PC1, y = PC2)) +
211   geom_point(color = "blue") +
212   labs(title = "PCA - First Two Principal Components", x = "Principal
213     Component 1", y = "Principal Component 2") +
214   theme_minimal()
215 # Get the variance of each component and visualize it
216 explained_variance <- summary(pca_res)$importance[, 1]
217 cumulative_variance <- cumsum(explained_variance)
218 var_explained_df <- data.frame(
219   Component = 1:length(explained_variance),
220   Variance_Explained = explained_variance,
221   Cumulative_Variance = cumulative_variance
222 )
223 # We visualize the variance explained by each component
224 ggplot(var_explained_df, aes(x = Component)) +
225   geom_bar(
226     aes(y = Variance_Explained),
227     stat = "identity",
228     fill = "blue",
229     alpha = 0.7
230   ) +
231   geom_line(aes(y = Cumulative_Variance * max(Variance_Explained)),
232     color = "red",
233     size = 1) +
234   labs(title = "Variance Explained by PCA Components", x = "Principal
235     Component", y = "Variance Explained") +

```

```

235 scale_y_continuous(sec.axis = sec_axis(~ . / max(var_explained_df$
236   Variance_Explained), name = "Cumulative Variance")) +
237 theme_minimal()
238 "Dada esta gr fica podemos observar que la varianza acumulada seg n
239   se a aden
240   m s componentes deja de aumentar significativamente a partir de cuatro
241   m s o menos.
242 De todas formas, con solo dos (para poder visualizarlo) ya se obtienen
243 una varianza
244 alta y suficiente para obtener cierta informaci n."
245
246 # Biplot
247 autoplot(
248   pca_res,
249   data = data,
250   colour = 'Mean_temperature',
251   loadings = TRUE,
252   loadings.label = TRUE,
253   loadings.label.size = 3
254 )
255 "En este biplot se puede observar como temperaturas muy bajas se dan en
256   lugares
257   con presiones altas y vientos altos, ya que los valores que se
258   encuentran en
259   esa zona son combinaciones de valores altos para esas variables.
260   Obviamente
261   valors bajos de las variables de temperatura (max, min) contribuyen
262   much simo
263   a la temperatura media, pero por ser ta obvios son menos interesantes."

```

6.2. EDA Clasificación

```

1 library(tidyverse)
2 library(moments)
3 library("corrplot")
4 library(ggfortify)
5
6 "Leemos los datos saltandonos las cabeceras iniciales y despu s las
7   parseamos
8   a mano. Seguido, vamos a ver unos pocos datos con head para hacernos
9   una primera
10  idea."
11
12 data <- read.csv("newthyroid/newthyroid.dat",
13                 skip = 10,
14                 header = FALSE)
15 colnames(data) <- c(
16   "T3resin",
17   "Thyroxin",
18   "Triiodothyronine",

```

```
17 "Thyroidstimulating",
18 "TSH_value",
19 "Class"
20 )
21 head(data)
22
23 "Con la funci n de summary obtenemos un resumen estad stico general
  de cada
24 columna. Adem s vamos a obtener otros datos informativos como la
  dimensi n,
25 y la estructura."
26
27 summary(data)
28 dim(data)
29 str(data)
30 colSums(is.na(data))
31
32 "Por lo que podemos ver, existen 6 variables o columnas (o
  caracter sticas) para
33 este conjunto. Existe gran variabilidad en caracter sticas como TSH_
  value y
34 Thyroidstimulating, donde valores muy bajos son cercanos al cero y muy
  altos
35 superan los 50.
36
37 Adem s, la mayor a de valores en estas variables parecen ser que
  rondan valores
38 bajos, esto nos lo puede indicar el tercer cuartil, que es en ambas
  variables
39 muy bajo en comparaci n al valor m ximo que toman.
40
41 Las variables son todas n mericas flotantes, excepto T3resin que es
  entera. No
42 hay valores faltantes en ninguna columna."
43
44 "Las clases a predecir en este problema de clasificaci n son las
  siguientes:"
45
46 table(data$Class)
47 prop.table(table(data$Class))
48
49 "Se puede observar un desbalanceo en el n mero de ejemplos para cada
  clase. La
50 clase 1 es muy mayoritaria."
51
52 "Vamos a normalizar todas las variables. Las que m s necesitan esta
  transformaci n son las anteriormente mencionadas, que tienen valores
  muy
54 separados en el espacio. Las clases adem s, deben transformarse en
  factores
55 para poder tratarlas como categor as."
56
57 data <- data %>%
58   mutate(Class = as.factor(Class)) %>%
```

```
59 mutate(across(where(is.numeric), ~ (. - mean(., na.rm = TRUE)) / sd
60 (.)))
61 selected_columns <- c("T3resin",
62                       "Thyroxin",
63                       "Triiodothyronine",
64                       "Thyroidstimulating",
65                       "TSH_value")
66
67 long_data <- data %>%
68   pivot_longer(cols = selected_columns,
69               names_to = "Variable",
70               values_to = "Value")
71
72 ggplot(long_data, aes(x = Variable, y = Value, fill = Variable)) +
73   geom_boxplot() +
74   labs(title = "Boxplots", x = "Variable", y = "Value") +
75   theme_minimal()
76
77 "Dada esta gráfica de boxplots podemos observar aquellos valores que
78 se alejan
79 demasiado de la distribución de los datos. De hecho, el método de
80 distancia
81 intercuartil es muy útil para eliminar outliers."
82
83 q3 <- quantile(data$TSH_value, 0.75)
84 iqr <- IQR(data$TSH_value)
85 upper_bound <- q3 + 1.5 * iqr
86
87 # Porcentaje de "outliers" en TSH_value
88 length(data$TSH_value[data$TSH_value > upper_bound]) / length(data$TSH_
89 value)
90
91 "Usando el método de distancia intercuartil (solo por arriba, ya que
92 en el
93 boxplot podemos ver que no hay ninguno por abajo) encontramos
94 potenciales
95 outliers. El porcentaje de ellos para la característica seleccionada
96 es muy
97 pequeño. En el caso de que estuviese seguro (conocimiento de dominio)
98 de que
99 esos valores son irreales, los quitaría, pero realmente no lo sé y el
100 quitarlos ahora me evitaría un posible futuro análisis con el modelo
101 de
102 predicción usado, ya que a priori no sé si van a afectar severamente
103 el
104 rendimiento de mi modelo."
105
106 "Creamos una pequeña función para poder mostrar un gráfico de barras
107 , de esta
108 forma podremos ver fácilmente las posibles distribuciones de cada
109 variable."
110 plot_distribution <- function(data,
```

```
101         var_name ,
102         binwidth = 0.1,
103         fill_color = "pink") {
104   ggplot(data, aes_string(x = var_name)) +
105     geom_density(fill = fill_color,
106                 color = "black",
107                 alpha = 0.7) +
108     labs(
109       title = paste("Distribution of", var_name),
110       x = var_name,
111       y = "Frequency"
112     ) +
113     theme_minimal()
114 }
115
116 plot_distribution(data, "TSH_value")
117
118 "Claramente tiene una asimetría a la izquierda."
119
120 plot_distribution(data, "T3resin")
121
122 "T3resin tiene un aspecto que parece bastante normal. Ligeramente
123   asimétrica
124   hacia la derecha."
125
126 plot_distribution(data, "Thyroxin")
127
128 "Thyroxin también tiene un aspecto que parece bastante normal."
129
130 plot_distribution(data, "Triiodothyronine")
131
132 "Triiodothyronine tiene una asimetría a la izquierda."
133
134 plot_distribution(data, "Thyroidstimulating")
135
136 "Thyroidstimulating tiene una asimetría a la izquierda muy acentuada."
137
138 "Dadas estas gráficas de densidades de cada variable, podemos observar
139   como se
140   distribuyen sus valores. Vemos que hay unas cuantas que cuentan con
141   asimetría
142   extrema. Estos que presentan esta asimetría son aquellas variables que
143   anteriormente habíamos identificado con pocos valores extremos y
144   muchos
145   valores bajos."
146
147 ggplot(long_data, aes(sample = Value)) +
148   stat_qq() +
149   stat_qq_line() +
150   facet_wrap(~ Variable) +
151   labs(title = "Qqplot of every variable", y = "Variable quantiles", x
152         = "Theoretical quantiles") +
153   theme_minimal()
```



```
150 "T3resin y Thyroxin son las que se parecen m s a una normal, pero a n
    as tienen
151 variaciones en las colas. El resto de variables son muy asim tricas ,
    ya se hab a
152 diagnosticado, pero es una confirmaci n m s.
153
154 Vamos a hacer un test de Shapiro sobre las dos variables m s
    prometedoras para
155 ver si podemos rechazar la hip tesis nula de que no siguen una
    distribuci n
156 normal."
157
158 shapiro.test(data$T3resin)
159 shapiro.test(data$Thyroxin)
160
161 "El test de shapiro es muy significativo para las dos variables m s
162 prometedoras, lo que indica que podemos rechazar las asumpciones de
    normalidad
163 que hemos estado construyendo por medio de las gr ficas. Ninguna
    variable de las
164 que tenemos sigue una distribuci n normal, esto lo tendremos en cuenta
    para la
165 fase de clasificaci n."
166
167 "Vamos a hacer test de aimetr a y de curtosis para las variables que
    sabemos
168 no son normales."
169
170 agostino.test(data$T3resin)
171 agostino.test(data$Thyroxin)
172 anscombe.test(data$T3resin)
173 anscombe.test(data$Thyroxin)
174
175 "Entendemos pues que los datos tienen asimetr a y curtosis, pues se
    rechazan en
176 ambos casos."
177
178
179 "Creamos otra funci n para ver la distribuci n, pero por clase. De
    esta forma
180 podemos observar que distribuciones de valores parecen asociarse a
    ciertas
181 clases."
182
183 plot_distribution_for_every_class <- function(data, var_name, binwidth
    = 0.6) {
184   ggplot(data, aes(x = !!sym(var_name), fill = Class)) +
185     geom_density(color = "black", alpha = 0.7) +
186     labs(
187       title = paste("Distribution of", var_name),
188       x = var_name,
189       y = "Frequency"
190     ) +
191     facet_wrap(~ Class) +
```

```
192   theme_minimal() +
193   theme(legend.position = "none")
194 }
195
196 plot_distribution_for_every_class(data, "Thyroidstimulating")
197 plot_distribution_for_every_class(data, "TSH_value")
198 plot_distribution_for_every_class(data, "T3resin")
199 plot_distribution_for_every_class(data, "Thyroxin")
200 plot_distribution_for_every_class(data, "Triiodothyronine")
201
202 "Parece observarse que para todas las variables, valores altos de la
    misma
203 suelen agruparse en la clase 1."
204
205 pca_res <- prcomp(data %>% select(-Class))
206
207 # Biplot
208 autoplot(
209   pca_res,
210   data = data,
211   colour = 'Class',
212   loadings = TRUE,
213   loadings.label = TRUE,
214   loadings.label.size = 3
215 )
216
217 "Vamos a mirar si existe alguna correlaci n entre variables, para ello
    podemos
218 usar dos m todos. Con pairs mostramos cada variable en relaci n a
    otra variable.
219 Lo bueno es que podemos de manera intuitiva identificar relaciones que
    vayan
220 m s all de la linealidad. Con la librer a de correlaci n, podemos
    mostrar
221 gr ficamente junto al coeficiente las correlaciones entre variables."
222
223 numeric_data <- select(data, -Class)
224 pairs(numeric_data)
225
226 "Una de las relaciones m s notables parece estar entre TSH_value y
    Thyroidstimulating, donde se observa una distribuci n bastante
    concentrada cerca
227 de los valores m s bajos, con algunos valores at picos.
228 Hay patrones interesantes entre T3resin y Thyroxin, que muestran una
    dispersi n
229 no aleatoria, sugiriendo alg n tipo de relaci n entre estas hormonas
    tiroideas.
230 Triiodothyronine tambi n muestra patrones de agrupaci n interesantes
    con otras
231 variables, especialmente visible en algunas de las dispersiones."
232
233
234 cor_matrix <- cor(numeric_data)
235 corrplot(cor_matrix,
236           method = "number",
```

```
237     type = "upper",
238     tl.col = "black")
239
240 "En el caso de este dataset, no encuentro ninguna correlaci n
    demasiado grande a
241 excepci n de Triiodothyronine y Thyroxin, que presenten una
    correlaci n positiva
242 y algo lineal bastante interesante.
243 Dados estos resultados, yo no eliminar a ninguna por redundancia.
    Adem s , la
244 dimensi n del problema no es tal como para tener que optimizar el
    tama o
245 de las componentes / caracter sticas que tenemos, pues hay muy pocas.
    De hecho
246 muchas veces, dos caracter sticas muy correladas no implican
    redundancia. Un
247 ejemplo de esto aparece en Feature Extraction - Foundations and
    Applications
248 by I. Guyon et al. (p.10, figure 2 (e)), donde dos caracter sticas con
249 un ratio de correlaci n alto son totalmente necesarias (ambas) para la
250 clasificaci n del problema."
251
252 "En cuanto a transformaciones, es posible que fuese muy til en
    aquellas
253 aquellas variables con asimetr as largas a la derecha, en ese caso
    podr amos
254 aplicar transformaciones logar tmicas."
255
256 data <- data %>%
257   mutate(
258     log_Triiodothyronine = log1p(Triiodothyronine),
259     log_Thyroidstimulating = log1p(Thyroidstimulating),
260     log_TSH_value = log1p(TSH_value)
261   )
262
263 plot_distribution(data, "log_Triiodothyronine")
264 plot_distribution(data, "Triiodothyronine")
265 plot_distribution(data, "log_Thyroidstimulating")
266 plot_distribution(data, "Thyroidstimulating")
267 plot_distribution(data, "log_TSH_value")
268 plot_distribution(data, "TSH_value")
269
270 shapiro.test(data$log_Triiodothyronine)
271 shapiro.test(data$Triiodothyronine)
272 agostino.test(data$log_Triiodothyronine)
273 agostino.test(data$Triiodothyronine)
274
275 shapiro.test(data$log_Thyroidstimulating)
276 shapiro.test(data$Thyroidstimulating)
277 agostino.test(data$log_Thyroidstimulating)
278 agostino.test(data$Thyroidstimulating)
279
280 shapiro.test(data$log_TSH_value)
281 shapiro.test(data$TSH_value)
```

```
282 agostino.test(data$log_TSH_value)
283 agostino.test(data$TSH_value)
284
285 "Se ha reducido la asimetría notablemente, aunque eso no hace que las
286 distribuciones de las variables sean normales. "
```

6.3. Regresión

```
1 library(tidyverse)
2 library(kknn)
3 library(RWeka)
4
5 "Leemos los datos saltandonos las cabeceras iniciales y después las
6 parseamos
7 a mano."
8
9 data <- read.csv("wankara/wankara.dat", skip = 14, header = FALSE)
10 colnames(data) <- c(
11   "Max_temperature",
12   "Min_temperature",
13   "Dewpoint",
14   "Precipitation",
15   "Sea_level_pressure",
16   "Standard_pressure",
17   "Visibility",
18   "Wind_speed",
19   "Max_wind_speed",
20   "Mean_temperature"
21 )
22 str(data)
23 "Vamos a aplicar las transformaciones realizadas en la parte del EDA."
24
25 data <- as.data.frame(scale(data))
26 data$Diff_temperature <- data$Max_temperature + data$Min_temperature
27 data$Diff_pressure <- data$Standard_pressure - data$Sea_level_pressure
28
29 data <- data[, !names(data) %in% c("Min_temperature",
30   "Max_temperature",
31   "Sea_level_pressure",
32   "Standard_pressure")]
33
34 summary(data)
35 pairs(data)
36
37 "Ahora vamos a utilizar las 5 variables regresoras que más sentido
38 tengan según
39 nuestro criterio. Según lo visto en el EDA, en lo relativo a
40 correlación y al
41 propio sentido común, la variable más prometedora es Diff_temperature
42 (es la
```

```

40 combinaci n lineal de Max_temperature y Min_temperature). Seguido de
    Dewpoint,
41 Diff_pressure (combinaci n lineal de variables originales de presi n)
    , Visibility
42 y Wind_speed.
43
44 Se ha tenido en cuenta las variables que parecen m s lineales con
    respecto a la
45 variable objetivo, o que al menos parezcan tener una relaci n en
    alg n sentido (por
46 ejemplo, Visibility es algo curva, algo logaritmica)."
47
48 perform_linear_regression <- function(data, target_var = "Mean_
    temperature", predictors) {
49   regression_results <- list()
50
51   for (predictor in predictors) {
52     formula <- as.formula(paste(target_var, "~", predictor))
53
54     fit <- lm(formula, data = data)
55
56     plot_data <- data.frame(x = data[[predictor]],
57                             y = data[[target_var]],
58                             fitted = predict(fit))
59
60     p <- ggplot(plot_data, aes(x = x, y = y)) +
61       geom_point(alpha = 0.5) +
62       geom_line(aes(y = fitted), color = "blue", size = 1) +
63       labs(
64         title = paste(target_var, "vs", predictor),
65         x = predictor,
66         y = target_var
67       ) +
68       theme_minimal()
69
70     regression_results[[predictor]] <- list(model = fit,
71                                             summary = summary(fit),
72                                             plot = p)
73
74     # Print summary and plot
75     cat("\n--- Regression Results for", predictor, "---\n")
76     print(regression_results[[predictor]]$summary)
77     print(regression_results[[predictor]]$plot)
78   }
79
80   return(regression_results)
81 }
82
83 predictors <- c("Diff_temperature",
84                 "Dewpoint",
85                 "Diff_pressure",
86                 "Visibility",
87                 "Wind_speed")
88 regression_analysis <- perform_linear_regression(data, predictors =

```

```
predictors)
89
90 "Claramente el mejor predictor de la temperatura media es diff_
    temperature, seguido
91 de Dewpoint. El RSE es extremadamente bajo usando ese predictor,
    adem s tiene un
92 R cuadrado ajustado del 0.98, lo que significa que con solo esa
    variable es posible
93 explicar el 98% de la variabilidad de la variable objetivo."
94
95 "Vamos a eliminar variables usando backwards selection."
96
97 fit1 <- data %>%
98   lm(formula = Mean_temperature ~ ., data = .)
99
100 summary(fit1)
101
102 fit2 <- data %>%
103   lm(formula = Mean_temperature ~ . - Max_wind_speed, data = .)
104
105 summary(fit2)
106
107 fit3 <- data %>%
108   lm(formula = Mean_temperature ~ . - Max_wind_speed - Precipitation,
109       data = .)
110
111 summary(fit3)
112
113 "A continuaci n vamos a ajustar un modelo de regresi n lineal
    m ltiple utilizando
114 varias variables y teniendo en cuenta relaciones no linales e
    interacciones. Aunque
115 los resultados son dif cilmente mejorables"
116
117 fit1 <- data %>%
118   lm(formula = Mean_temperature ~ Diff_temperature, data = .)
119
120 summary(fit1)
121
122
123 fit2 <- data %>%
124   lm(formula = Mean_temperature ~ Diff_temperature + Dewpoint + Diff_
125       pressure,
126       data = .)
127
128 summary(fit2)
129
130 fit3 <- data %>%
131   lm(
132     formula = Mean_temperature ~ Diff_temperature + Dewpoint + Diff_
133     pressure + Visibility,
134     data = .
135   )
```

```
135 summary(fit3)
136
137 "A adir poco a poco todas las variables mejora el resultado, pero es
    m nimo. No creo
138 que sea un camino a seguir. Voy a probar con interacciones entre
    variables que
139 tengan sentido."
140
141 fit4 <- data %>%
142   lm(formula = Mean_temperature ~ Diff_temperature * Dewpoint,
143       data = .)
144
145 summary(fit4)
146
147 fit5 <- data %>%
148   lm(formula = Mean_temperature ~ Diff_temperature * I(Dewpoint ^ 2),
149       data = .)
150
151 summary(fit5)
152
153 fit6 <- data %>%
154   lm(formula = Mean_temperature ~ Diff_temperature * I(Dewpoint ^ 2) *
155       Dewpoint,
156       data = .)
157
158 summary(fit6)
159
160 fit7 <- data %>%
161   lm(
162     formula = Mean_temperature ~ Diff_temperature * I(Dewpoint ^ 2) *
163     Dewpoint * I(Diff_temperature ^
164
165     2),
166     data = .
167   )
168
169 summary(fit7)
170
171 "Parece que la interacci n de las dos variables m s importantes junto
    a no linealidad
172 mejora m s el resultado."
173
174 "Voy a probar con Visibilidad al cuadrado, ya que pareciese por la
    gr fica que
175 sigue una relaci n no lineal con respecto a la variable objetivo."
176
177 fit8 <- data %>%
178   lm(formula = Mean_temperature ~ Diff_temperature + I(Visibility ^ 2),
179       data = .)
180
181 summary(fit8)
182
183 "Probando varias combinaciones, lo que he podido ver es que el modelo
    parece
```

```
181 ajustar mejor a añadiendo interacciones entre varias variables en vez de
    solo
182 a adirlas. Alguna no linealidad (sobre todo con Dewpoint) parece
    a adir algo
183 de mejora. Voy a usar el modelo final del m todo backwards por la
    interpretabilidad
184 y porque el mejor modelo que usa interacciones y no linealidad (fit7)
    no mejora tanto
185 como para plantearse sacrificar interpretabilidad."
186
187 "Ahora vamos a aplicar regresión con el algoritmo kNN."
188
189 fitknn1 <- kknk(Mean_temperature ~ ., data, data)
190
191 "Visualizamos los datos originales vs los puntos ajustados del knn"
192 ggplot(data = data, aes(x = Diff_temperature, y = Mean_temperature)) +
193   geom_point() +
194   geom_point(aes(y = fitknn1$fitted.values),
195             color = "blue",
196             shape = 20) +
197   labs(x = "Diff_temperature", y = "Mean_temperature", title = "kNN") +
198   theme_minimal()
199
200 "Calculamos el RMSE"
201 yprime <- fitknn1$fitted.values
202 sqrt(sum((data$Mean_temperature - yprime) ^ 2) / length(yprime))
203
204 "Utilizando todos los datos podemos obtener un RMSE prácticamente
    igual que en
205 regresión usando solo la variable de diff_temperature. Este problema
    concreto
206 parece bastante apto para una regresión lineal, de todas formas
    probar
207 algunas combinaciones más."
208
209 fitknn2 <- kknk(Mean_temperature ~ Diff_temperature, data, data)
210
211 "Calculamos el RMSE"
212 yprime <- fitknn2$fitted.values
213 sqrt(sum((data$Mean_temperature - yprime) ^ 2) / length(yprime))
214
215 "Usando solo la variable más importante es capaz de conseguir un mejor
    ajuste
216 que la regresión lineal."
217
218 ggplot(data = data, aes(x = Diff_temperature, y = Mean_temperature)) +
219   geom_point() +
220   geom_point(aes(y = fitknn2$fitted.values),
221             color = "blue",
222             shape = 20) +
223   labs(x = "Diff_temperature", y = "Mean_temperature", title = "kNN") +
224   theme_minimal()
225
226 "Podemos apreciar mucho menos ruido que usando todas las variables."
```



```

227
228 fitknn3 <- kknk(Mean_temperature ~ Diff_temperature + Dewpoint, data,
229               data)
230
231 "Calculamos el RMSE"
232 yprime <- fitknn3$fitted.values
233 sqrt(sum((data$Mean_temperature - yprime) ^ 2) / length(yprime))
234
235 ggplot(data = data, aes(x = Diff_temperature, y = Mean_temperature)) +
236   geom_point() +
237   geom_point(aes(y = fitknn3$fitted.values),
238             color = "blue",
239             shape = 20) +
240   labs(x = "Diff_temperature", y = "Mean_temperature", title = "kNN") +
241   theme_minimal()
242
243 "A añadiendo Dewpoint se consigue un ajuste con menor ruido incluso. Con
244   diff_temperature
245 se pod a modelar una l nea muy buena que ajustase los datos, pero
246   utilizando adem s
247   dewpoint parece que la varianza se modela mejor."
248
249 fitknn4 <- kknk(Mean_temperature ~ Diff_temperature * Dewpoint * I(
250               Dewpoint ~
251               ,
252               data,
253               data)
254
255 "Calculamos el RMSE"
256 yprime <- fitknn4$fitted.values
257 sqrt(sum((data$Mean_temperature - yprime) ^ 2) / length(yprime))
258
259 ggplot(data = data, aes(x = Diff_temperature, y = Mean_temperature)) +
260   geom_point() +
261   geom_point(aes(y = fitknn4$fitted.values),
262             color = "blue",
263             shape = 20) +
264   labs(x = "Diff_temperature", y = "Mean_temperature", title = "kNN") +
265   theme_minimal()
266
267 "Si en vez de sumar las variables a adimos interacciones entre ellas y
268   adem s tenemos
269 en cuenta un poco de no linealidad mejora el ajuste."
270
271 "He probado con otras combinaciones, pero no parece que encuentre nada
272   mejor. Adem s
273 hay que tener en cuenta que quiz estemos ajustando demasiado los
274   resultados a la
275 muestra, lo cu l en un entrenamiento real llevar a al sobreajuste."
276
277 "Vamos a realizar 5-fold cross validation con el mejor modelo obtenido
278   de regresión
279 lineal m ltiple y a partir de ah , intentar obtener un mejor modelo

```

```
con kNN. Después
272 compararemos todos los resultados entre sí y con el algoritmo M5."
273
274 run_k_fold_cv <- function(data, model = "lm", k = 5) {
275   columns <- c(
276     "Dewpoint",
277     "Precipitation",
278     "Max_temperature",
279     "Min_temperature",
280     "Sea_level_pressure",
281     "Standard_pressure",
282     "Visibility",
283     "Wind_speed",
284     "Max_wind_speed",
285     "Mean_temperature"
286   )
287
288   colnames(data) <- columns
289
290   train_mse_list <- numeric()
291   test_mse_list <- numeric()
292
293   data$Diff_temperature <- data$Max_temperature + data$Min_temperature
294   data$Diff_pressure <- data$Standard_pressure - data$Sea_level_
     pressure
295
296   data <- data[, !names(data) %in% c("Min_temperature",
297                                     "Max_temperature",
298                                     "Sea_level_pressure",
299                                     "Standard_pressure")]
300
301   idx <- sample(nrow(data))
302   data <- data[idx, ]
303
304   fold_size <- round(nrow(data) / k)
305
306   results <- data.frame(
307     Actual = numeric(),
308     Predicted = numeric(),
309     Fold = numeric(),
310     MSE = numeric(),
311     Train = numeric(),
312     Test = numeric()
313   )
314
315   for (i in 0:(k - 1)) {
316     start <- 1 + i * fold_size
317     end <- min((i + 1) * fold_size, nrow(data))
318
319     val <- data[start:end, ]
320     train <- data[-(start:end), ]
321
322     train_scaled <- scale(train)
323     mean_train <- attr(train_scaled, "scaled:center")
324     sd_train <- attr(train_scaled, "scaled:scale")
325   }
326 }
```

```
324
325   train <- as.data.frame(train_scaled)
326
327   val <- as.data.frame(scale(val, center = mean_train, scale = sd_
train))
328
329   if (model == "lm") {
330     fit <- lm(Mean_temperature ~ ., data = train)
331     train_pred <- predict(fit, newdata = train)
332     yprime <- predict(fit, newdata = val)
333   } else if (model == "kkn") {
334     fit <- knn(Mean_temperature ~ ., train = train, test = val)
335     fit_train <- knn(Mean_temperature ~ ., train = train, test =
train)
336     train_pred <- fit_train$fitted.values
337     yprime <- fit$fitted.values
338   } else if (model == "m5") {
339     fit <- M5P(Mean_temperature ~ ., data = train)
340     train_pred <- predict(fit, newdata = train)
341     yprime <- predict(fit, newdata = val)
342   }
343
344   train_mse <- mean((train$Mean_temperature - train_pred) ^ 2)
345   test_mse <- mean((val$Mean_temperature - yprime) ^ 2)
346
347   train_mse_list <- c(train_mse_list, train_mse)
348   test_mse_list <- c(test_mse_list, test_mse)
349
350   results <- rbind(
351     results,
352     data.frame(
353       Actual = val$Mean_temperature,
354       Predicted = as.numeric(yprime),
355       Fold = i + 1,
356       MSE = test_mse,
357       Train = tail(train_mse_list, 1),
358       Test = tail(test_mse_list, 1)
359     )
360   )
361 }
362
363 results
364 }
365
366 data <- read.csv("wankara/wankara.dat", skip = 14, header = FALSE)
367 lm_results <- run_k_fold_cv(data, model = "lm")
368 knn_results <- run_k_fold_cv(data, model = "kkn")
369 m5_results <- run_k_fold_cv(data, model = "m5")
370
371 # Summary of results
372 lm_mse <- mean(lm_results$MSE)
373 knn_mse <- mean(knn_results$MSE)
374 m5_mse <- mean(m5_results$MSE)
375
```

```
376 lm_mse_train <- mean(lm_results$Train)
377 knn_mse_train <- mean(knn_results$Train)
378 m5_mse_train <- mean(m5_results$Train)
379
380 print(paste("Linear Regression MSE:", lm_mse))
381 print(paste("KNN MSE:", knn_mse))
382 print(paste("M5 MSE:", m5_mse))
383
384 plot_predictions <- function(results, model) {
385   ggplot(results, aes(
386     x = Actual,
387     y = Predicted,
388     color = as.factor(Fold)
389   )) +
390   geom_point(alpha = 0.7) +
391   geom_abline(
392     slope = 1,
393     intercept = 0,
394     linetype = "dashed",
395     color = "black"
396   ) +
397   labs(
398     title = paste(model, "-", "Results"),
399     x = "Actual Mean Temperature",
400     y = "Predicted Mean Temperature",
401     color = "Fold"
402   ) +
403   theme_minimal()
404 }
405
406 plot_lm <- plot_predictions(lm_results, "Linear Regression")
407 plot_knn <- plot_predictions(knn_results, "KNN")
408 plot_m5 <- plot_predictions(m5_results, "M5")
409
410 print(plot_lm)
411 print(plot_knn)
412 print(plot_m5)
413
414 "Por los resultados de test y lo visto en las gráficas, parece que el
415 mejor se ajusta a este dataset es regresión lineal múltiple seguido
416 muy de cerca por knn."
417
418 "Vamos a comparar los algoritmos."
419
419 resultados_test <- read.csv("regr_test_alumnos.csv")
420 resultados_train <- read.csv("regr_train_alumnos.csv")
421 tablatst <- cbind(resultados_test[, 2:dim(resultados_test)[2]])
422 colnames(tablatst) <- names(resultados_test)[2:dim(resultados_test)[2]]
423 rownames(tablatst) <- resultados_test[, 1]
424 tablatra <- cbind(resultados_train[, 2:dim(resultados_train)[2]])
425 colnames(tablatra) <- names(resultados_train)[2:dim(resultados_train)
426 [2]]
426 rownames(tablatra) <- resultados_train[, 1]
```

```
427
428 tablatst[17, 1] <- lm_mse
429 tablatst[17, 2] <- knn_mse
430 tablatst[17, 3] <- m5_mse
431
432 tablatra[17, 1] <- lm_mse_train
433 tablatra[17, 2] <- knn_mse_train
434 tablatra[17, 3] <- m5_mse_train
435
436 "Normalizamos utilizando las diferencias relativas entre los resultados
    de los
437 algoritmos. Despues, generamos una tabla con valores ajustados y
    procedemos al
438 test de Wilcoxon, que es un test por pares."
439
440 difs <- (tablatst[, 1] - tablatst[, 2]) / tablatst[, 1]
441 wilc_1_2 <- cbind(ifelse (difs < 0, abs(difs) + 0.1, 0 + 0.1),
442                  ifelse (difs > 0, abs(difs) + 0.1, 0 + 0.1))
443 colnames(wilc_1_2) <- c(colnames(tablatst)[1], colnames(tablatst)[2])
444 head(wilc_1_2)
445
446 LMvsKNNtst <- wilcox.test(wilc_1_2[, 1],
447                          wilc_1_2[, 2],
448                          alternative = "two.sided",
449                          paired = TRUE)
450 Rmas <- LMvsKNNtst$statistic
451 pvalue <- LMvsKNNtst$p.value
452 LMvsKNNtst <- wilcox.test(wilc_1_2[, 2],
453                          wilc_1_2[, 1],
454                          alternative = "two.sided",
455                          paired = TRUE)
456 Rmenos <- LMvsKNNtst$statistic
457 Rmenos
458 Rmas
459 pvalue
460
461 "Dado un p-valor de menos de 0.7, no se puede rechazar la hipotesis
    nula, por lo que
462 no podemos asegurar que existen diferencias estadisticamente
    significativas entre
463 KNN y LM."
464
465 "Ahora realizamos test multiples para comparar todos los algoritmos
    entre s. Para
466 ello usamos el test de friedman."
467
468 test_friedman <- friedman.test(as.matrix(tablatst))
469 test_friedman
470
471 "Se aplica una correccion para evitar errores acumulados."
472
473 tam <- dim(tablatst)
474 groups <- rep(1:tam[2], each = tam[1])
475 pairwise.wilcox.test(as.matrix(tablatst),
```

```
476         groups,
477         p.adjust = "holm",
478         paired = TRUE)
479
480 "Esto indica que el algoritmo 3 (M5) no tiene evidencia estadística de
481     ser mejor sobre
482 los otros dos algoritmos. En este test además se puede ver como el
483 algoritmo 1 y el
484 2 (KNN vs LM) tampoco tienen diferencias significativas entre ellos,
485 como se pudo
486 ver en el anterior análisis."
```

```
487 "Vamos a hacer lo mismo con los resultados de train."
488
489 difs <- (tablatra[, 1] - tablatra[, 2]) / tablatra[, 1]
490 wilc_1_2 <- cbind(ifelse (difs < 0, abs(difs) + 0.1, 0 + 0.1),
491                   ifelse (difs > 0, abs(difs) + 0.1, 0 + 0.1))
492 colnames(wilc_1_2) <- c(colnames(tablatra)[1], colnames(tablatra)[2])
493 head(wilc_1_2)
494
495 LMvsKNNtst <- wilcox.test(wilc_1_2[, 1],
496                           wilc_1_2[, 2],
497                           alternative = "two.sided",
498                           paired = TRUE)
499 Rmas <- LMvsKNNtst$statistic
500 pvalue <- LMvsKNNtst$p.value
501 LMvsKNNtst <- wilcox.test(wilc_1_2[, 2],
502                           wilc_1_2[, 1],
503                           alternative = "two.sided",
504                           paired = TRUE)
505 Rmenos <- LMvsKNNtst$statistic
506 Rmenos
507 Rmas
508 pvalue
509
510 apply(tablatra, median)
511 apply(tablatst, median)
512
513 test_friedman <- friedman.test(as.matrix(tablatra))
514 test_friedman
515
516 tam <- dim(tablatra)
517 groups <- rep(1:tam[2], each = tam[1])
518 pairwise.wilcox.test(as.matrix(tablatra),
519                       groups,
520                       p.adjust = "holm",
521                       paired = TRUE)
522
523 "Lo esperado en este apartado era que se obtuvieran resultados más
524     significativos.
525 Si los algoritmos iban bien en test, en training deberían haber ido
526     incluso mejor.
527 Aunque esta premisa no es cierta siempre, en este caso puede observarse
528     como
```

```
524 los p-valores son m s bajos incluso. Por ello, parece que los mismos
    algoritmos siguen
525 siendo los mejores y no parece que se haya realizan un sobreajuste."
```

6.4. Clasificación

```
1 library(caret)
2 library(kknn)
3 library("MASS")
4 library(tidyverse)
5
6 "Leemos los datos saltandonos las cabeceras iniciales y despu s las
   parseamos
7 a mano. Seguido, vamos a ver unos pocos datos con head para hacernos
   una primera
8 idea."
9
10 data <- read.csv("newthyroid/newthyroid.dat",
11                  skip = 10,
12                  header = FALSE)
13 colnames(data) <- c(
14   "T3resin",
15   "Thyroxin",
16   "Triiodothyronine",
17   "Thyroidstimulating",
18   "TSH_value",
19   "Class"
20 )
21 summary(data)
22
23 "Normalizamos los datos usando el escalado z-score. Esto hace que los
   datos sigan
24 una distribuci n con media 0 y desviaci n t pica 1. Esto es muy
   necesario cuando
25 se tienen datos en distintas escalas, como es el caso, y se utilizan
   algoritmos
26 que usan distancias."
27
28 data <- data %>% mutate(Class = as.factor(Class))
29 class_column = data$Class
30 data <- as.data.frame(scale(data[1:ncol(data) - 1]))
31 data$Class <- class_column
32
33 data <- data %>%
34   mutate(
35     Thyroidstimulating = log1p(Thyroidstimulating),
36     TSH_value = log1p(TSH_value)
37   )
38 summary(data)
39
40 "Vamos a dividir el conjunto d datos en conjuntos de train y test. Para
   eso definimos"
```

```
41 la siguiente función"
42
43 train_test_split <- function(data,
44                               test_percentage = 0.2,
45                               seed = 42) {
46   set.seed(seed)
47   data <- data[sample(1:nrow(data)), ]
48
49   n <- round((1 - test_percentage) * nrow(data))
50
51   train <- data[1:n, 1:ncol(data)]
52   test <- data[(n + 1):nrow(data), 1:ncol(data)]
53
54   list(train = train, test = test)
55 }
56
57 results <- train_test_split(data)
58 train <- results$train
59 test <- results$test
60 str(results)
61
62 "Una vez divididos los conjuntos podemos empezar a experimentar con KNN
63 ."
64 fit1 <- knn(Class ~ .,
65             train = train,
66             test = test,
67             k = 1)
68
69 make_confusion_matrix <- function(knn_model) {
70   predicted <- fitted(knn_model)
71   conf_matrix <- table(Predicted = predicted, Actual = test$Class)
72   conf_matrix_df <- as.data.frame(as.table(conf_matrix))
73
74   plot <- ggplot(conf_matrix_df, aes(x = Actual, y = Predicted, fill =
75   Freq)) +
76     geom_tile() +
77     geom_text(aes(label = sprintf("%d", Freq)),
78               color = "white",
79               size = 10) +
80     scale_fill_gradient(low = "cyan", high = "darkgreen") +
81     theme_minimal() +
82     theme(
83       axis.title = element_text(size = 12, face = "bold"),
84       axis.text = element_text(size = 10),
85       legend.title = element_text(size = 10),
86       plot.title = element_text(
87         size = 14,
88         face = "bold",
89         hjust = 0.5
90       )
91     ) +
92     labs(
93       title = "Confusion Matrix",
94       x = "Actual Class",
```



```
93     y = "Predicted Class",
94     fill = "Count"
95   ) +
96   coord_equal()
97
98   accuracy <- sum(diag(conf_matrix)) / sum(conf_matrix)
99   print(plot)
100   print(paste("Accuracy:", round(accuracy, 2)))
101 }
102
103 make_confusion_matrix(fit1)
104
105 "Dadas estas m tricas parece que el knn est  generalizando muy bien
106   con k=1, lo
107   cual es bastante asombroso, pues k=1 es un valor muy dado al
108   sobreajuste ya que
109   el modelo se hace mucho m s sensible a cambios cuando a adimos nuevos
110   datos, se
111   hace muy variable."
112
113 "Vamos a probar varios valores de k para ver como var an los
114   resultados."
115
116 fit2 <- kknn(Class ~ .,
117               train = train,
118               test = test,
119               k = 5)
120
121 make_confusion_matrix(fit2)
122
123
124 fit3 <- kknn(Class ~ .,
125               train = train,
126               test = test,
127               k = 15)
128
129 make_confusion_matrix(fit3)
130
131
132 fit4 <- kknn(Class ~ .,
133               train = train,
134               test = test,
135               k = 40)
136
137 make_confusion_matrix(fit4)
138
139
140 fit5 <- kknn(Class ~ .,
141               train = train,
142               test = test,
143               k = 100)
144
145 make_confusion_matrix(fit5)
146
147 "El modelo con k cada vez mayor pierde precisi n , porque al hacer esto
148   , el
```

```
142 modelo se suaviza y se vuelve m s general. Esto puede reducir el
    sobreajuste
143 al hacer que las predicciones dependan de un n mero mayor de vecinos,
    lo que
144 ayuda a promediar los datos y a reducir la sensibilidad a las
    fluctuaciones
145 peque as. Sin embargo, si k es demasiado grande, el modelo comienza a
    perder
146 capacidad para captar patrones espec ficos , lo que puede llevar a un
    subajuste
147 (underfitting), ya que los puntos de datos m s cercanos pueden no
    tener tanta
148 influencia sobre la predicci n final. "
```

149

```
150 "En mi caso escoger a k=5. Aunque k=1 sea mejor y me haya dado mejores
    resultados
151 en test, a adir un poco de regularizaci n en este caso, considero que
    no deval a
152 demasiado la calidad del modelo y previene de posibles ajustes
    demasiado finos. Si
153 bien es verdad que el resultado de test es mejor para k=1 y por tanto
    el modelo,
154 incluso con k=1 parece no sobreajustar nada en datos nunca vistos,
    prefiero ser
155 precavido."
```

156

```
157 "Ahora vamos a usar LDA para ajustar un modelo al conjunto de datos.
    Pero primero
158 han de comprobarse ciertas asunciones"
```

159

```
160 "Los datos son normales para cada clase."
161 str(data)
162 test_variable <- function(var_name) {
163     resultados <- data %>%
164         group_by(Class) %>%
165         summarize(
166             p_value = if (is.numeric(.data[[var_name]]) &&
167                           all(!is.na(.data[[var_name]]))) {
168                 shapiro.test(.data[[var_name]])$p.value
169             } else {
170                 NA
171             },
172             normal = ifelse(p_value > 0.05, "Puede ser normal", "No es normal
173 "),
174             .groups = "drop"
175         )
176     plot <- ggplot(data, aes(sample = .data[[var_name]])) +
177         stat_qq() + stat_qq_line() +
178         facet_wrap(~ Class)
179     print(plot)
180     print(resultados)
181 }
182 test_variable("T3resin")
```

```
183 test_variable("Thyroxin")
184 test_variable("Triiodothyronine")
185 test_variable("Thyroidstimulating")
186 test_variable("TSH_value")
187
188 "Parece que el test de Shapiro no rechaza para algunas variables dentro
189 de algunas
190 clases. Si bien no se puede decir que no son normales, tampoco se
191 pueden rechazar.
192 Dada esta evidencia y la de los qqplots, no se verifica la primera
193 asunci n de LDA."
194
195 "Para comprobar si cada clase tiene matrices de varianze-covarianza
196 id nticas, se
197 puede utilizar el test de Bartlett"
198
199 bartlett.test(T3resin ~ Class, data)
200 bartlett.test(Thyroxin ~ Class, data)
201 bartlett.test(Triiodothyronine ~ Class, data)
202 bartlett.test(Thyroidstimulating ~ Class, data)
203 bartlett.test(TSH_value ~ Class, data)
204
205 "Se rechaza el test de Bartlett para todas las variables, lo que
206 significa que
207 las varianzas de los grupos o muestras comparadas no son iguales, es
208 decir, no
209 se cumple el supuesto de homogeneidad de varianzas (homocedasticidad)."
```

```
210
211 "Clasificamos con LDA"
212 lda_model <- lda(Class ~ ., data = train)
213 lda_model
214 lda.pred.train <- predict(lda_model, train)
215 lda.pred.test <- predict(lda_model, test)
216
217 plot_data <- lda.pred.train$x %>%
218   as_tibble() %>%
219   mutate(Class = train$Class)
220
221 "Mostramos el gr fico donde se muestra como se distribuyen las
222 observaciones de
223 las clases en el espacio generado por LDA utilizando las primeras dos
224 componentes
225 lineales discriminantes (LD1 y LD2)."
```

```
226
227 "Se puede observar que las clases est n bien separadas en el espacio
228 generado por
229 LDA. Una buena separaci n indica que el modelo ha logrado discriminar
230 correctamente entre clases."
```

```
231
232 ggplot(data = plot_data) +
233   geom_point(aes(x = LD1, y = LD2, color = Class)) +
234   scale_colour_manual(
235     name = "Class",
236     values = c("red", "green", "blue"),
```

```

228   labels = c("1", "2", "3")
229 ) +
230 labs(title = "Data Transformed After LDA")
231
232 t <- table(lda.pred.test$class, test$Class)
233 t
234
235 sum(diag(t)) / nrow(test)
236
237 plot_data <- lda.pred.test$x %>%
238   as_tibble() %>%
239   mutate(known = test$Class, # Replace with the correct column for
240          class labels
241          prediction = lda.pred.test$class) %>%
242   pivot_longer(c("prediction", "known"),
243               names_to = "Type",
244               values_to = "Class")
245
246 ggplot(data = plot_data) +
247   geom_point(aes(
248     x = LD1,
249     y = LD2,
250     shape = Type,
251     color = Class
252   )) +
253   scale_colour_manual(
254     name = "Species",
255     values = c("red", "green", "blue"),
256     labels = c("1", "2", "3")
257   ) +
258   scale_shape_manual(name = "Type", values = c(5, 3)) +
259   labs(title = "Validation Data + Predictions Transformed After LDA")
260
261 "Pese a no cumplirse las normalidades de datos por clase y la igualdad
262 de covarianza-varianza, el modelo es capaz de generalizar muy bien y de
263 encontrar
264 separabilidad entre las clases."
265
266 qda_model <- qda(Class ~ ., data = train)
267 qda_model
268 qda.pred.train <- predict(qda_model, train)
269 qda.pred.test <- predict(qda_model, test)
270
271 t <- table(Predicted = qda.pred.test$class, Actual = test$Class)
272 t
273
274 sum(diag(t)) / nrow(test)
275
276 conf_matrix_df <- as.data.frame(as.table(t))
277
278 ggplot(conf_matrix_df, aes(x = Actual, y = Predicted, fill = Freq)) +
279   geom_tile() +
280   geom_text(aes(label = sprintf("%d", Freq)), color = "white", size =

```

```
10) +
279 scale_fill_gradient(low = "cyan", high = "darkgreen") +
280 theme_minimal() +
281 theme(
282   axis.title = element_text(size = 12, face = "bold"),
283   axis.text = element_text(size = 10),
284   legend.title = element_text(size = 10),
285   plot.title = element_text(size = 14, face = "bold", hjust = 0.5)
286 ) +
287 labs(title = "Confusion Matrix",
288       x = "Actual Class",
289       y = "Predicted Class",
290       fill = "Count") +
291 coord_equal()
292
293 "Vamos a comparar los tres algoritmos."
294
295 k_fold_cross_validation <- function(data,
296                                     k = 10,
297                                     model_function,
298                                     metric_function,
299                                     seed = 42) {
300   set.seed(seed)
301   folds <- sample(1:k, size = nrow(data), replace = TRUE)
302   performance_metrics <- c()
303
304   for (i in 1:k) {
305     test_indices <- which(folds == i)
306     train_indices <- setdiff(1:nrow(data), test_indices)
307
308     train_data <- data[train_indices, ]
309     test_data <- data[test_indices, ]
310
311     model <- model_function(train_data, test_data)
312     if (inherits(model, "kkn")) {
313       predictions <- fitted(model)
314     } else {
315       predictions <- predict(model, test_data)
316     }
317
318     if (inherits(model, "lda") | inherits(model, "qda")) {
319       predictions_class <- predictions$class
320     } else {
321       predictions_class <- predictions
322     }
323
324     performance <- metric_function(predictions_class, test_data$Class)
325     performance_metrics <- c(performance_metrics, performance)
326   }
327
328   mean_performance <- mean(performance_metrics)
329   std_deviation <- sd(performance_metrics)
330
331   return(list(mean = mean_performance, std_deviation = std_deviation))
}
```

```
332 }
333
334 model_function_knn <- function(train_data, test_data) {
335   model <- kknn(Class ~ .,
336                 train = train_data,
337                 test = test_data,
338                 k = 5)
339   return(model)
340 }
341
342 model_function_lda <- function(train_data, test_data) {
343   model <- lda(Class ~ ., data = train_data)
344   return(model)
345 }
346
347 model_function_qda <- function(train_data, test_data) {
348   model <- qda(Class ~ ., data = train_data)
349   return(model)
350 }
351
352 metric_function_accuracy <- function(predictions, actual) {
353   confusion_matrix <- table(predicted = predictions, actual = actual)
354   accuracy <- sum(diag(confusion_matrix)) / sum(confusion_matrix)
355   return(accuracy)
356 }
357
358 result_knn <- k_fold_cross_validation(
359   data,
360   k = 10,
361   model_function = model_function_knn,
362   metric_function = metric_function_accuracy
363 )
364 print(paste(
365   "KNN Accuracy: ",
366   round(result_knn$mean, 4),
367   " ",
368   round(result_knn$std_deviation, 4)
369 ))
370
371 result_lda <- k_fold_cross_validation(
372   data,
373   k = 10,
374   model_function = model_function_lda,
375   metric_function = metric_function_accuracy
376 )
377 print(paste(
378   "LDA Accuracy: ",
379   round(result_lda$mean, 4),
380   " ",
381   round(result_lda$std_deviation, 4)
382 ))
383
384 result_qda <- k_fold_cross_validation(
385   data,
```

```
386 k = 10,  
387 model_function = model_function_qda,  
388 metric_function = metric_function_accuracy  
389 )  
390 print(paste(  
391     "QDA Accuracy: ",  
392     round(result_qda$mean, 4),  
393     " ",  
394     round(result_qda$std_deviation, 4)  
395 ))
```



7. Bibliografía

- [1] Y. Shao y M. Zhou, “A characterization of multivariate normality through univariate projections,” *Journal of Multivariate Analysis*, vol. 101, n.º 10, págs. 2637-2640, 2010, ISSN: 0047-259X. DOI: <https://doi.org/10.1016/j.jmva.2010.04.015>. URL: <https://www.sciencedirect.com/science/article/pii/S0047259X10001168>.
- [2] I. Guyon, S. Gunn, M. Nikravesh y L. Zadeh, *Feature extraction. Foundations and applications. Papers from NIPS 2003 workshop on feature extraction, Whistler, BC, Canada, December 11–13, 2003. With CD-ROM*. ene. de 2006, vol. 207, ISBN: 978-3-540-35487-1. DOI: 10.1007/978-3-540-35488-8.

