

Introducción a la Ciencia de Datos

Trabajo Teórico/Práctico Integrador

Cristian González Guerrero

24 de agosto de 2017

Índice

1. Introducción	2
2. Presentación y visualización de los datos	2
2.1. abalone	2
2.2. tae (Teaching Assistant Evaluation)	5
3. Regresión sobre el conjunto de datos abalone	10
3.1. Regresión lineal simple	10
3.2. Regresión lineal múltiple	10
3.3. Regresión mediante el algoritmo k-NN	13
3.4. Comparación de los resultados	16
3.4.1. Comparación entre los modelos lineales y regresiones con k-NN para el caso estudiado	16
3.4.2. Comparación genérica entre los modelos lineal, k-NN y M5'	20
4. Clasificación sobre el conjunto de datos tae	20
4.1. Clasificación mediante el algoritmo k-NN	20
4.2. Clasificación mediante análisis discriminante lineal (LDA)	21
4.3. Clasificación mediante análisis discriminante cuadrático (QDA)	23
4.4. Comparación de los algoritmos	23
Referencias	25
A. Código fuente: Análisis de datos	26
A.1. Código fuente usado en la visualización de los datos <i>abalone</i>	26
A.2. Código fuente usado en la visualización de los datos <i>tae</i>	29
B. Código fuente: Regresión	32
B.1. Código fuente de las funciones usadas en los ejercicios de regresión	32
B.2. Código fuente usado en el ejercicio de regresión lineal	35
B.3. Código fuente usado en el ejercicio de regresión k-NN	40
B.4. Código fuente de la comparación entre los distintos algoritmos de regresión	45
C. Código fuente: Clasificación	48
C.1. Código fuente de de las funciones usadas en el ejercicio de clasificación	48
C.2. Código fuente usado en el ejercicio de clasificación	51

1. Introducción

Este es el Trabajo Teórico/Práctico Integrador de la asignatura *Introducción a la Ciencia de Datos*, impartida en el *Máster en Ciencia de Datos e Ingeniería de Computadores* (DATCOM). En el mismo, se desarrollan los ejercicios propuestos por los profesores, aplicando los conocimientos adquiridos en la asignatura.

La realización de este trabajo se ha llevado a cabo usando *RStudio* [1, 2], en su versión 0.99.903 - *R version 3.3.2 (2016-10-31)* [3]. La mayor parte de los gráficos han sido generados usando *ggplot2* [4, 5], en su versión 2.2.1, lo cual se ha considerado una ampliación del apartado de visualización.

Este trabajo, realizado de forma original por Cristian González Guerrero, se presenta para su valoración en la evaluación extraordinaria de Septiembre de 2017.

2. Presentación y visualización de los datos

Los conjuntos de datos asignados por los profesores son los siguientes.

abalone se trata de un conjunto de datos que presenta distintas variables fisiológicas del abulón, un marisco muypreciado en la cocina oriental, junto con el número de anillos que presenta, del cual puede deducirse directamente su edad [6, 7]. El ejercicio consiste en la estimación de la edad del abulón a partir de las otras variables fisiológicas medidas a través de la regresión.

tae (Teaching Assistant Evaluation) se trata de un conjunto de datos que recoge las evaluaciones del desempeño docente de 150 asistentes en el Departamento de Estadística de la Universidad de Wisconsin-Madison [8, 9, 6]. Las puntuaciones fueron divididas en tres grupos de tamaño similar. Estos grupos o clases son el objetivo del ejercicio de clasificación.

A continuación se presenta un estudio más pormenorizado de los citados datasets. El código fuente utilizado para obtener los datos aquí presentados puede consultarse en section §A.

2.1. abalone

El conjunto de datos *abalone* corresponde a un dataframe con una dimensión de 4176 filas y 9 columnas. Esto corresponde a 4176 muestras en las que se han medido 9 propiedades. Estas propiedades o campos se presentan en table 1, no apareciendo en ninguna valores omitidos (*missing values*). Como puede observarse en la tabla, todos los datos son de tipo numérico, salvo el sexo, que es una variable categórica que sólo toma tres valores. El tipo de datos de la variable de salida es un número entero, que se tratará como una variable continua a la hora de realizar la regresión.

Las medidas estadísticas de las distintas variables pueden observarse en table 2 y en table 3. Estas medidas estadísticas nos hablan de las distribuciones de las distintas variables, que pueden visualizarse mediante un diagrama de caja y bigotes como el mostrado en figure 1.

En el ejercicio de regresión, que usará estos datos, será conveniente hacerse una idea de la correlación entre las distintas variables. Para ello, se ha calculado el coeficiente de correlación de Pearson entre las distintas variables numéricas, obteniendo así la matriz de correlación de table 4. Esta matriz nos da una idea de la relación lineal que tienen las variables entre sí. Otra forma de visualizar esta relación es con un diagrama de dispersión. En figure 2 puede observarse la relación que existe entre las distintas variables cuantitativas de entrada y la variable de salida (fila en negrita en table 4).

En vista de las medidas de correlación realizadas, podemos concluir que este conjunto de datos presenta muchas variables correlacionadas entre sí. Esto quiere decir que muchas de las variables de entrada aportan información redundante que puede causar sobreajuste en los modelos de regresión lineal. Por otro lado, la correlación de las distintas variables de entrada con respecto a la salida es considerablemente alta, por lo que un modelo sencillo podría servir para predecir la edad del abulón con un error relativamente pequeño.

Attribute	Domain	Data type	Description
Sex	[1,3]	factor	1 = M 2 = F 3 = I (infant)
Length	[0.075, 0.815]	numeric	Longest shell measurement (mm)
Diameter	[0.055, 0.65]	numeric	Diameter perpendicular to length (mm)
Height	[0.0, 1.13]	numeric	Height with meat in shell (mm)
Whole_weight	[0.0020, 2.8255]	numeric	Weight of the whole abalone (grams)
Shucked_weight	[0.0010, 1.488]	numeric	Weight of meat (grams)
Viscera_weight	[0.0005, 0.76]	numeric	Gut weight after bleeding (grams)
Shell_weight	[0.0015, 1.005]	numeric	Shell weight after being dried (grams)
Rings	[1,29]	integer	+1.5 gives the age in years

Cuadro 1: Descripción de los atributos presentes en *abalone* [6, 10].

Variable	Min.	1st Qu.	Median	Mean	3rd Qu.	Max.	Var.	SD
Length	0.075	0.450	0.545	0.524	0.615	0.815	0.014	0.120
Diameter	0.055	0.350	0.425	0.408	0.480	0.650	0.010	0.099
Height	0.000	0.115	0.140	0.140	0.165	1.130	0.002	0.042
Whole_weight	0.002	0.442	0.800	0.829	1.153	2.826	0.240	0.490
Shucked_weight	0.001	0.186	0.336	0.359	0.502	1.488	0.049	0.222
Viscera_weight	0.001	0.094	0.171	0.181	0.253	0.760	0.012	0.110
Shell_weight	0.002	0.130	0.234	0.239	0.329	1.005	0.019	0.139
Rings	1.000	8.000	9.000	9.934	11.000	29.000	10.396	3.224

Cuadro 2: Medidas estadísticas sobre las variables cuantitativas de *abalone*.

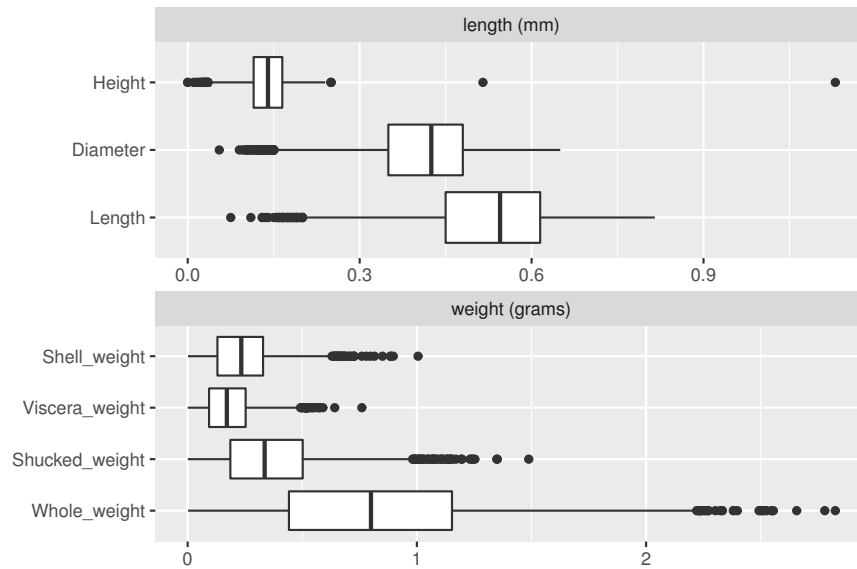
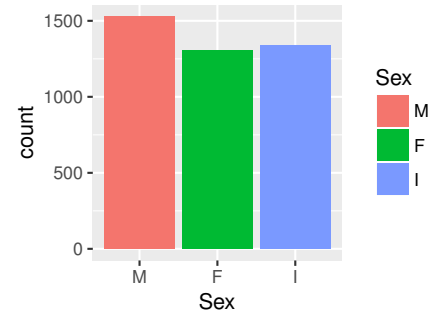


Figura 1: Distribución de las variables cuantitativas de entrada en *abalone*.

Sex	Freq.	Rel. Freq.
M	1528	0.366
F	1307	0.313
I	1341	0.321



Cuadro 3: Tabla de frecuencias de las variables categóricas de *abalone*, junto con su diagrama de barras.

Correlation	Length	Diameter	Height	Whole_weight	Shucked_weight	Viscera_weight	Shell_weight	Rings
Length	1.000	0.987	0.828	0.925	0.898	0.903	0.898	0.557
Diameter	0.987	1.000	0.834	0.925	0.893	0.900	0.905	0.575
Height	0.828	0.834	1.000	0.819	0.775	0.798	0.817	0.557
Whole_weight	0.925	0.925	0.819	1.000	0.969	0.966	0.955	0.540
Shucked_weight	0.898	0.893	0.775	0.969	1.000	0.932	0.883	0.421
Viscera_weight	0.903	0.900	0.798	0.966	0.932	1.000	0.908	0.504
Shell_weight	0.898	0.905	0.817	0.955	0.883	0.908	1.000	0.627
Rings	0.557	0.575	0.557	0.540	0.421	0.504	0.627	1.000

Cuadro 4: Correlación entre las variables cuantitativas en *abalone*.

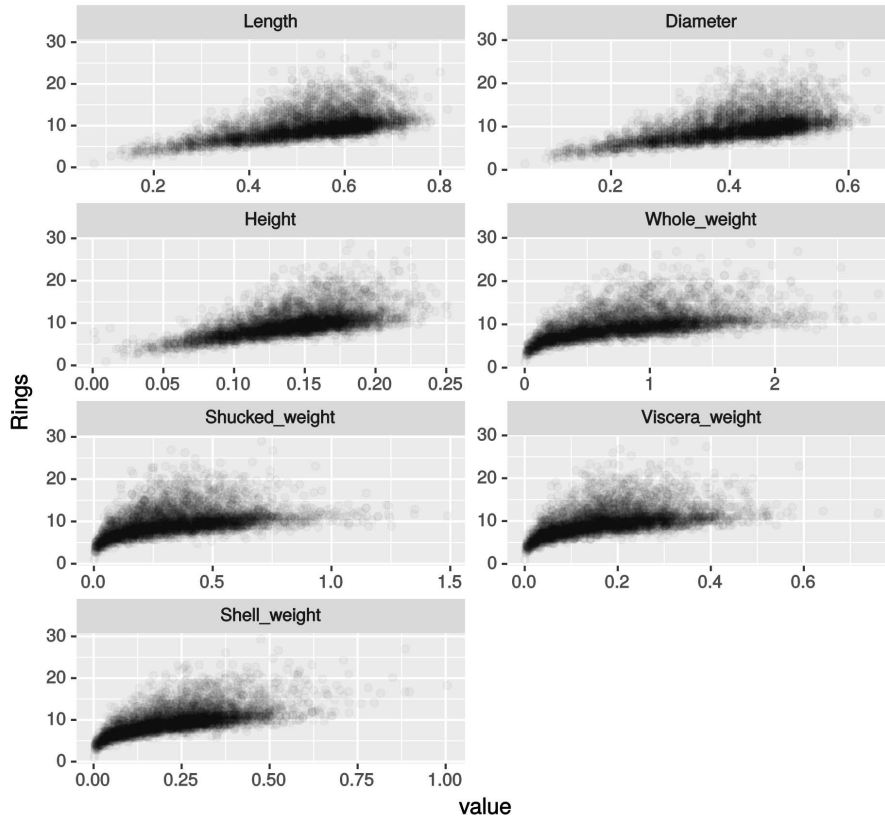


Figura 2: Nube de puntos de las variables de entrada con respecto a la salida en *abalone*.

Attribute	Domain	Data type	Description
Native	[1, 2]	factor	whether of not the TA is a native English speaker: 1 = English speaker 2 = non-English speaker
Instructor	[1, 25]	factor	course instructor (25 categories)
Course	[1, 26]	factor	26 categories
Semester	[1, 2]	factor	summer or regular semester: 1 = Summer 2 = Regular
Size	[3, 66]	integer	class size (numerical)
Class	{1, 2, 3}	factor	TA score: 1 = low 2 = medium 3 = high

Cuadro 5: Descripción de los atributos presentes en *tae* [6, 10].

Variable	Min.	1st Qu.	Median	Mean	3rd Qu.	Max.	Var.	SD
Size	3	19	27	27.93	37	66	166.83	12.92
Instructor	1	8	13	13.58	20	25	46.31	6.81
Course	1	3	4.5	8.14	15	26	49.49	7.03

Cuadro 6: Medidas estadísticas sobre la única variables cuantitativa de *tae*. Las dos últimas filas representan las mismas medidas estadísticas sobre dos variables categóricas, a las que se han asignado etiquetas numéricas empezando en el 1.

2.2. *tae* (Teaching Assistant Evaluation)

El conjunto de datos *tae* corresponde a un dataframe con una dimensión de 150 filas y 6 columnas que no presentan valores omitidos. Esto corresponde a 150 muestras en las que se han medido 6 propiedades tal y como se muestra en table 5. Como puede observarse en la tabla, todas las variables son categóricas, salvo el tamaño de la clase, que es un entero. El tipo de datos de la variable de salida es también categórico, pudiendo tomar tres valores distintos: *low*, *medium* y *high*, por lo que tendremos que realizar una clasificación con tres clases.

Las medidas estadísticas sobre la variable numérica *Size* pueden observarse en table 6. Junto a ellas, se han calculado los momentos estadísticos de las variables *Instructor* y *Course*, que, a pesar de ser categóricas, tienen asignados valores numéricos enteros. No obstante, esta representación numérica no es apropiada, ya que los números han sido dados de forma arbitraria y no tienen un significado concreto. Esto se pone de manifiesto en la figura figure 3, en la que puede observarse que la distribución de estas variables no es normal, ni sigue ningún patrón conocido. Más bien se trata de unos puntos aleatorios, ya que cada instructor, o, equivalentemente cada curso, tendrá un número de asistentes determinado, independiente del número asignado al instructor o al curso. La figura figure 4, en cambio, muestra la distribución de la variable *Size*. Esta variable tiene una distribución más cercana a la normal, aunque existen diferencias en función de la clase a la que pertenecen los datos.

La figura figure 5a muestra la frecuencia absoluta de las variables categóricas, separándolas según la clase de salida. Como puede observarse, existe una diferencia perceptible entre los asistentes nativos y no nativos, y también entre el tipo de semestre. Por un lado, hay menos asistentes nativos, pero suelen obtener mejores resultados. Lo mismo pasa con el semestre de verano: a pesar de haber menos asistentes, los resultados son mejores. Podría pensarse que esto viene dado por el número de alumnos, que tal vez sea menor en verano, y también es posible que sólo los más privilegiados puedan tener un asistente nativo. Sin embargo, un test de Kolmogorov–Smirnov bilateral nos indica

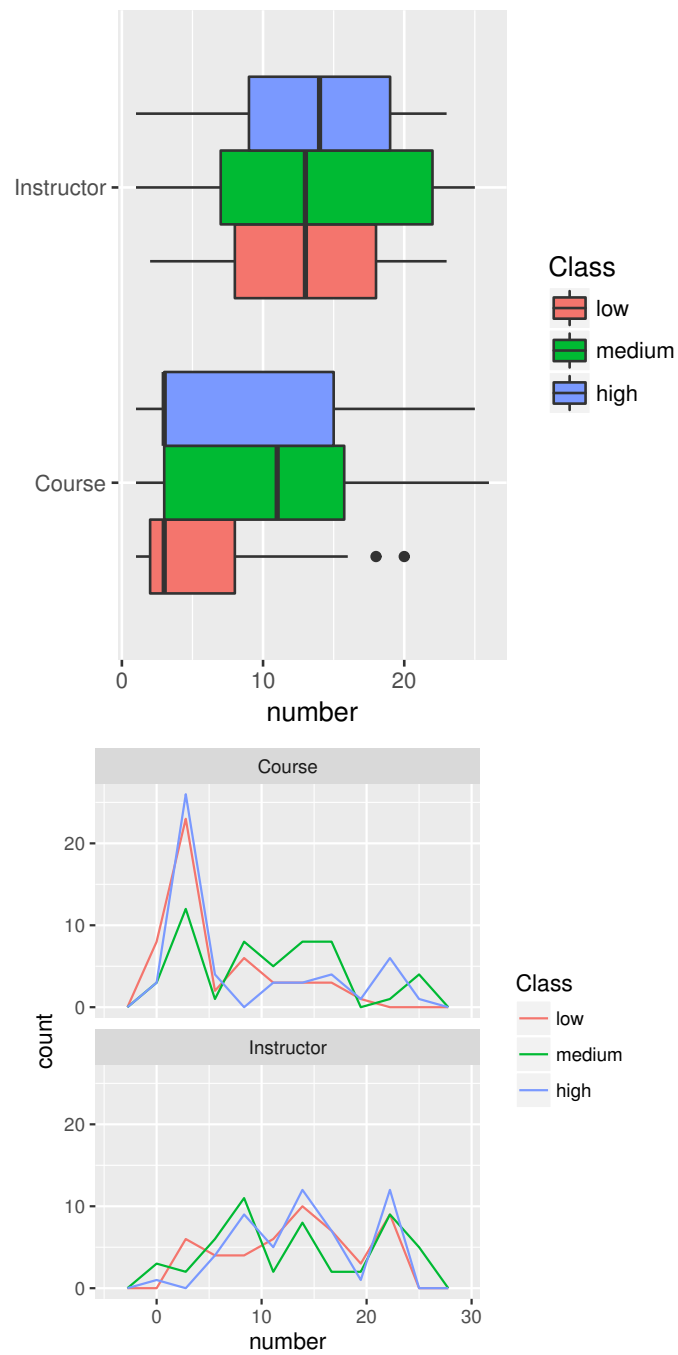


Figura 3: Diagrama de caja con bigotes de los valores numéricos de las variables `Course` e `Instructor`, junto a su distribución de frecuencias.

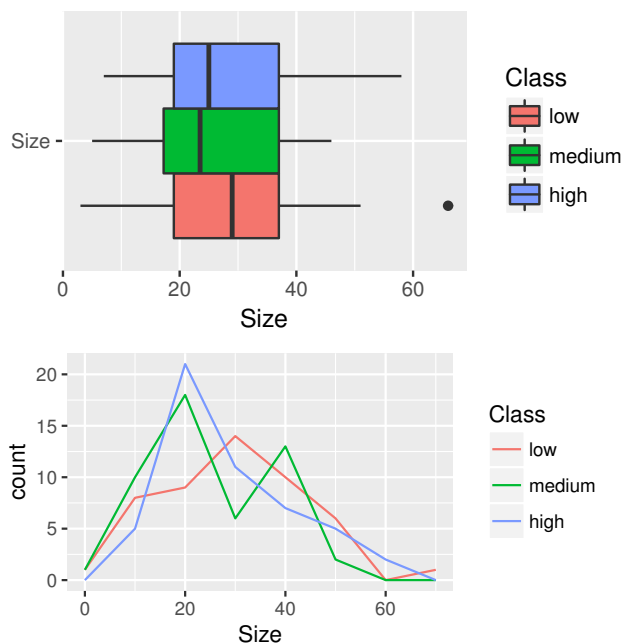


Figura 4: Diagrama de caja con bigotes de la variable numérica `Size` junto a su distribución de frecuencias.

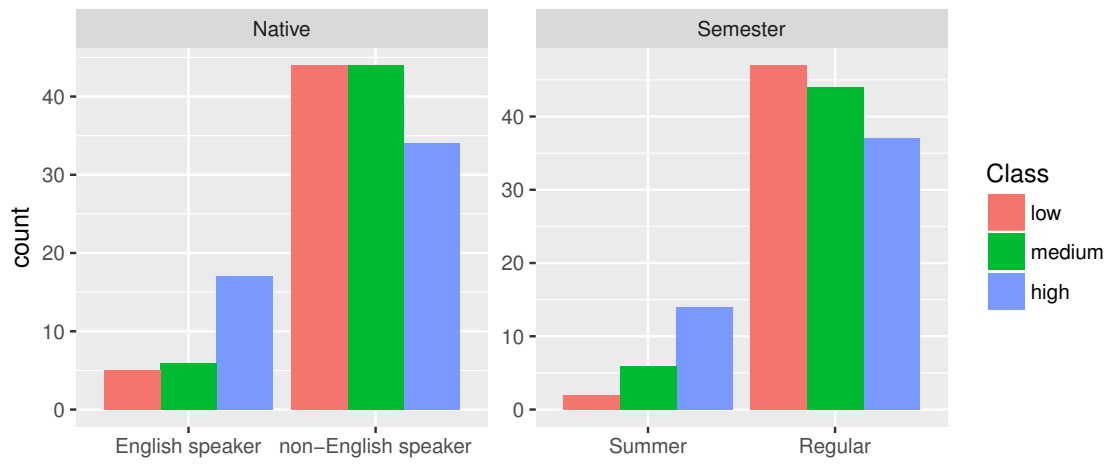
que no hay diferencia significativa en la distribución del tamaño del aula con asistentes nativos o no nativos. Sin embargo, el mismo test indica que sí existe una diferencia significativa en la distribución del tamaño del aula entre los semestres de verano e invierno.

Como se describía en el enunciado del problema, los tres grupos en los que se debe llevar a cabo la clasificación tienen una proporción parecida de registros, correspondientes a un tercio cada una. Esto se pone de manifiesto en la figura figure 5.

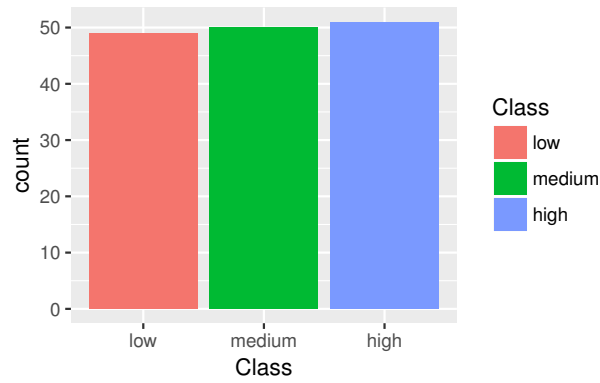
Volviendo a las variables categóricas comentadas anteriormente, debe decirse que el presente estudio ha sido realizado en 26 cursos, con 25 instructores diferentes. El número de asistentes que ha impartido clases en cada curso y con cada instructor ha sido variable. El valor numérico de cada curso y de cada instructor ha sido dado de manera arbitraria, de forma que los datos no presentan un orden. En la figura figure 6 se ha representado la frecuencia de las ocurrencias de `Course`, esto es, el número de asistentes que ha impartido clases en cada curso, ordenadas de mayor a menor. Como puede observarse, no hay un orden específico en el eje de abscisas, ni tampoco una distribución clara de los datos. Junto a la frecuencia absoluta, se ha representado la frecuencia relativa de cada clase, es decir, el número de asistentes con rendimiento alto, medio y bajo en proporción al número de asistentes en cada curso. En este caso, puede observarse que el rendimiento no depende del número de asistentes que haya tenido el curso, puesto que los valores no siguen ninguna tendencia en concreto. Este mismo análisis se ha realizado para la variable `Instructor`, obteniendo las mismas conclusiones (ver figura figure 7).

Un hecho que puede llamar la atención de este análisis es la existencia de cursos e instructores con resultados 100 % buenos (barritas azules), 100 % malos (barritas rojas), 100 % medios (barritas verdes) o mezcla (barritas de varios colores). Esto se debe, en parte, al número de asistentes que han dado clase en cada curso o con cada instructor. Cuando se tiene un número elevado de registros (muchos asistentes), los resultados de rendimiento suelen estar divididos. La tabla table 7 muestra un recuento de estos resultados, basándose en las observaciones de las figuras figure 6 y figure 7.

Como conclusión, hay que decir que se trata de un ejercicio de clasificación muy complicado. Esto se debe a que la mayor parte de las entradas son categóricas, mientras que los algoritmos de clasificación suelen esperar variables numéricas, y a la poca relación vista entre las distintas variables. Teniendo esto en cuenta, es posible predecir un rendimiento pobre en la clasificación.



(a) Variables de entrada.



(b) Variable de salida.

Figura 5: Conteo de variables categóricas.

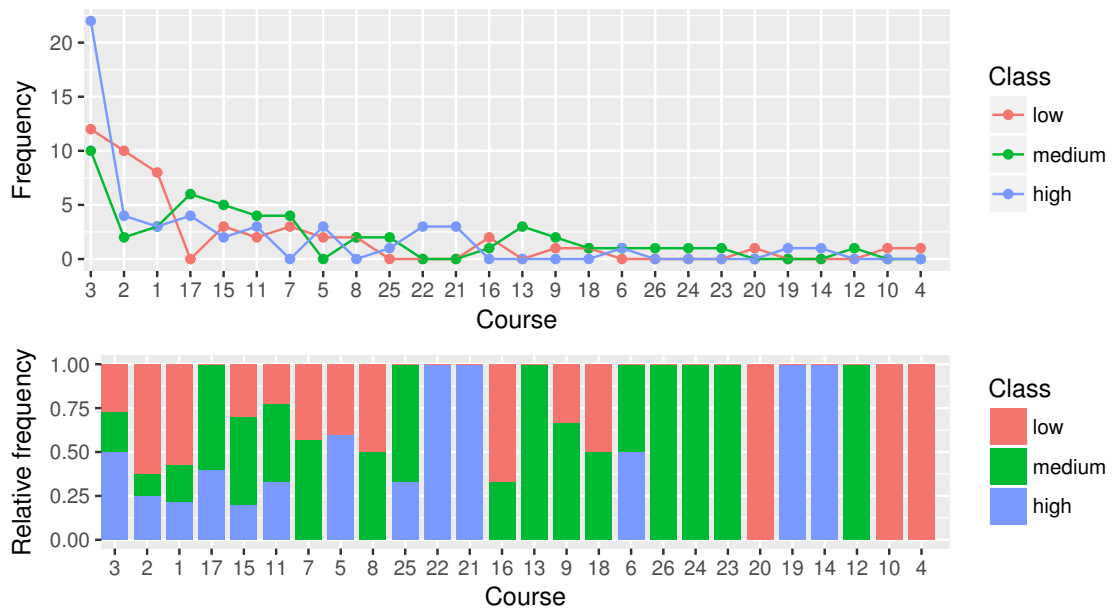


Figura 6: Frecuencia absoluta y relativa de las ocurrencias de Course en función de su clase. Los datos han sido ordenados para que las mayores frecuencias queden a la izquierda.

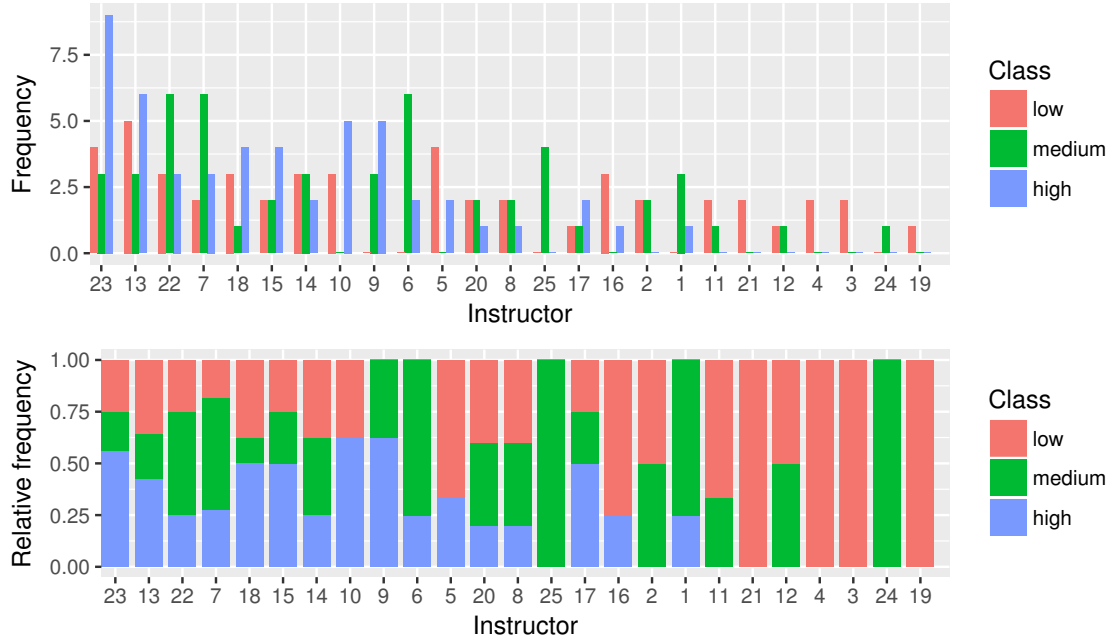


Figura 7: Frecuencia absoluta y relativa de las ocurrencias de `Instructor` en función de su clase. Los datos han sido ordenados para que las mayores frecuencias queden a la izquierda.

Class (performance)	Course	Instructor
low	3	4
medium-low	6	3
medium-low-high	5	10
medium	5	2
low-high	1	3
medium-high	2	3
high	4	0

Cuadro 7: Recuento del número de grupos e instructores en función del rendimiento.

Regressor	Summary of fit				Cross validation	
	Linear model ($\hat{Y}(X)$)	MSE	R^2	p -value	\overline{MSE}_{tra}	\overline{MSE}_{tst}
Sex	$10,71 + 0,42 X_F - 2,81 X_I$	8.39	0.1929	$< 2.2e-16$	8.39	8.40
Length	$2,10 + 14,94 X$	7.17	0.3098	$< 2.2e-16$	7.17	7.18
Diameter	$2,32 + 18,67 X$	6.96	0.3301	$< 2.2e-16$	6.96	6.97
Height	$3,94 + 42,96 X$	7.16	0.3107	$< 2.2e-16$	7.15	7.33
Whole_weight	$6,99 + 3,55 X$	7.36	0.2919	$< 2.2e-16$	7.36	7.37
Shucked_weight	$7,74 + 6,11 X$	8.55	0.1771	$< 2.2e-16$	8.55	8.57
Viscera_weight	$7,26 + 14,82 X$	7.76	0.2537	$< 2.2e-16$	7.76	7.77
Shell_weight	$6,46 + 14,53 X$	6.30	0.3937	$< 2.2e-16$	6.30	6.32

Cuadro 8: Regresiones lineales simples, con cada regresor.

3. Regresión sobre el conjunto de datos abalone

Tras haber visualizado los datos y haber explorado las posibles relaciones entre las distintas variables de entrada (ver sección 2.1), procedemos a aplicar la regresión lineal sobre el conjunto de datos. Para ello, se han explorado diversos modelos, que van desde la regresión lineal simple hasta modelos que contemplan las no linealidades e interacciones entre las variables. Posteriormente, se ha experimentado la regresión con el algoritmo k-NN. En este caso, se han evaluado los modelos encontrados con la regresión lineal.

En todos los casos, la bondad de las regresiones se ha medido usando validación cruzada con el método *k-fold cross validation*. De esta forma, ha sido posible medir el error cuadrático medio sobre el conjunto de test en cada regresión MSE_{tst} , realizando una media sobre los distintos conjuntos de test \overline{MSE}_{tst} . El error cuadrático medio sobre el conjunto de entrenamiento MSE_{tra} ha sido de utilidad para la evaluación del sobreajuste (*overfitting*). Para la ejecución de la validación cruzada se ha hecho uso de las particiones ya existentes.

3.1. Regresión lineal simple

Se ha realizado un modelo de regresión lineal sobre cada regresor. Los datos de los mismos pueden verse en la tabla 8. Como puede observarse, el regresor que ha dado mejor resultado ha sido Shell_weight, seguido por Diameter. Esta elección se ha basado en minimizar el error cuadrático medio sobre el conjunto de test (\overline{MSE}_{tst}), ya que consideraremos que el modelo tiene que tener la habilidad de generalizar (evitar sobreajuste). Llama la atención que los errores MSE medidos sobre los conjuntos de entrenamiento y el conjunto de test son muy parecidos. Esto se debe a que se está usando un modelo de regresión con una única variable, que difícilmente puede sobreajustar. En los apartados siguientes veremos un comportamiento bastante distinto, en el que \overline{MSE}_{tst} será claramente mayor a \overline{MSE}_{tra} .

Observando las gráficas de la figura 2, podemos intuir que existen no linealidades entre algunas de las variables de entrada y la de salida. En concreto, es claro que todos los pesos (variables _weight) siguen una curva parecida a una raíz o un logaritmo. Se han llevado a cabo diversas experiencias, aplicando operaciones sencillas a los datos de entrada. En particular, se ha probado la raíz cuadrada, cúbica y el logaritmo. Estas tres operaciones tienen un sentido físico sobre algunas de las variables. En concreto, el volumen ocupado (directamente proporcional al peso) es una función cúbica del radio, mientras que la superficie de contacto con el exterior es una función cuadrática. Esto es un factor limitante del crecimiento. Por otro lado, hay especies que crecen de forma exponencial, como las colonias de bacterias. De entre las operaciones probadas, el logaritmo ha dado el mejor resultado, al aumentar el índice de correlación lineal de Pearson y mejorar el ajuste lineal, como puede verse en la tabla 9 y la figura 9.

3.2. Regresión lineal múltiple

A continuación, se ha realizado la regresión lineal teniendo en cuenta todos los regresores, sin contemplar interacciones entre ellos. Se ha visto, mediante *backward selection* que los regresores que

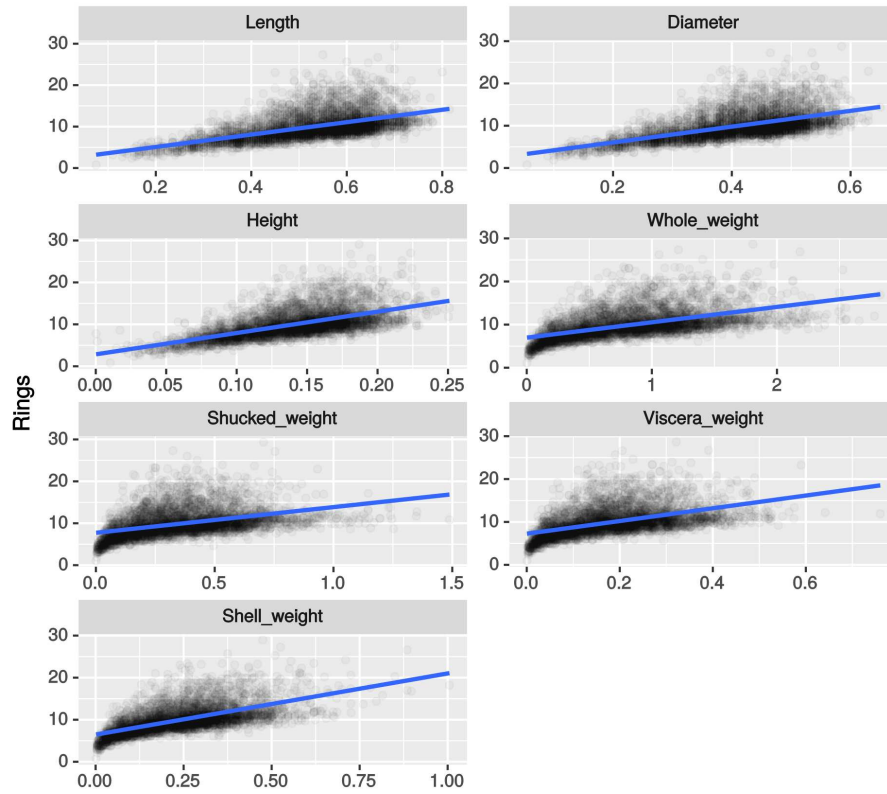


Figura 8: Regresiones lineales simples, con cada regresor.

Regressor	Correlation with Rings	Summary of fit			Cross validation	
		$\hat{Y}(X)$	R^2	p -value	\overline{MSE}_{tra}	\overline{MSE}_{tst}
Whole_weight	0.540	$6,99 + 3,55 X$	0.2919	$< 2.2e-16$	7.36	7.37
$\log(\text{Whole_weight})$	0.584	$10,93 + 2,25 X$	0.3411	$< 2.2e-16$	6.85	6.86
Shucked_weight	0.421	$7,74 + 6,11 X$	0.1771	$< 2.2e-16$	8.55	8.57
$\sqrt{\text{Shucked_weight}}$	0.477	$5,49 + 7,85 X$	0.2270	$< 2.2e-16$	8.03	8.04
$\sqrt[3]{\text{Shucked_weight}}$	0.492	$3,43 + 9,64 X$	0.2423	$< 2.2e-16$	7.87	7.89
$\log(\text{Shucked_weight})$	0.510	$12,40 + 1,91 X$	0.2604	$< 2.2e-16$	7.69	7.70
Viscera_weight	0.504	$7,26 + 14,82 X$	0.2537	$< 2.2e-16$	7.76	7.77
$\log(\text{Viscera_weight})$	0.566	$14,17 + 2,15 X$	0.3198	$< 2.2e-16$	7.07	7.08
Shell_weight	0.627	$6,46 + 14,53 X$	0.3937	$< 2.2e-16$	6.30	6.32
$\sqrt[3]{\text{Shell_weight}}$	0.649	$0,87 + 15,31 X$	0.4212	$< 2.2e-16$	6.02	6.03
$\log(\text{Shell_weight})$	0.634	$14,14 + 2,52 X$	0.4015	$< 2.2e-16$	6.22	6.23

Cuadro 9: Corrección de no linealidades en las distintas variables.

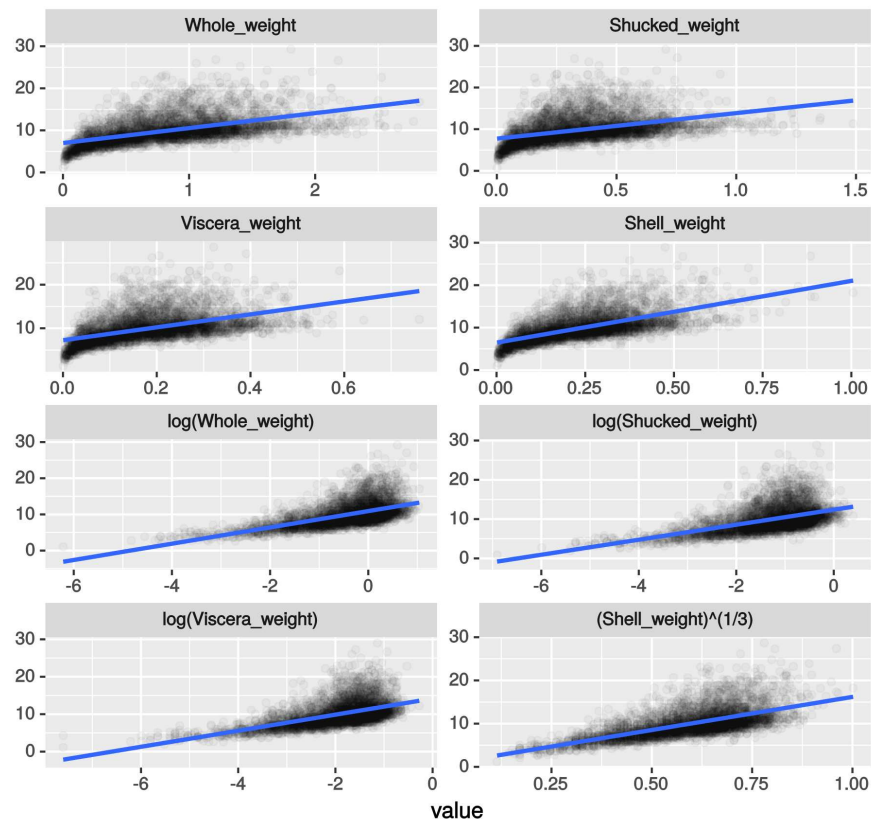


Figura 9: Regresión lineal simple en las variables en las que se han tenido en cuenta las no linealidades. Claramente, el modelo que aplica preprocesamiento tiene un mejor ajuste.

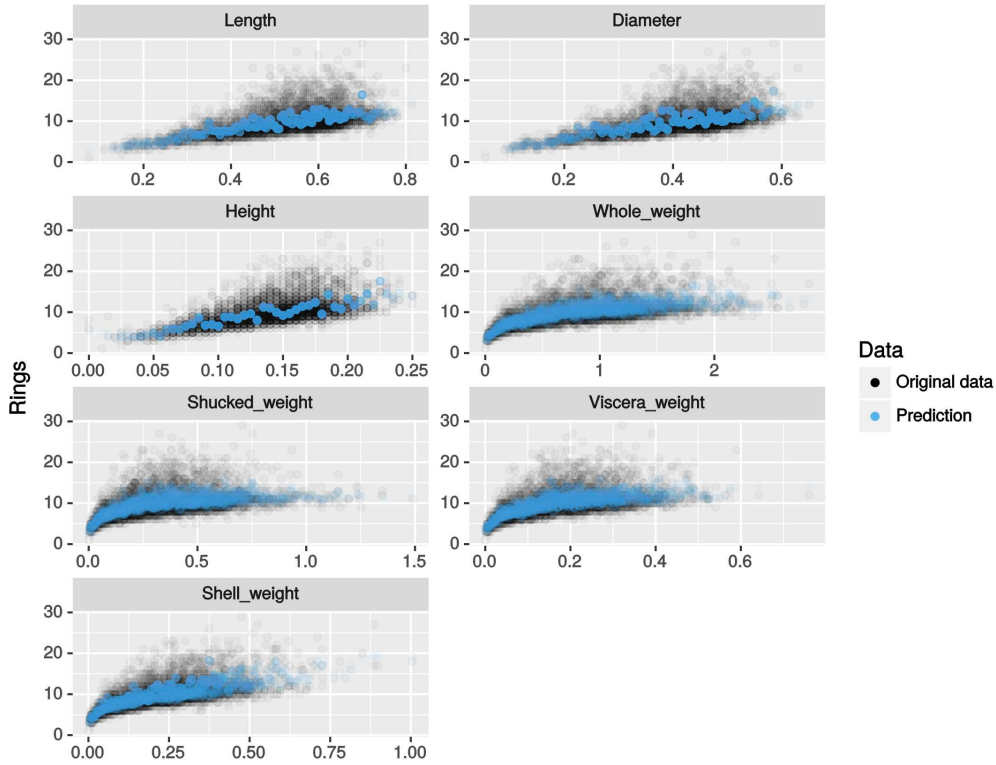


Figura 10: Regresiones con el algoritmo k-NN, con cada regresor de forma individual ($k = 7$).

producen mejor resultados son Diameter, Height y todos los `_weight`, pudiéndose descartar Sex y Length por presentar un p -value muy elevado. En la tabla 10 se puede contemplar los resultados de esta segunda aproximación, así como el de varias iteraciones en las que se han ido eliminando cada vez más regresores.

Así mismo, se ha realizado la regresión múltiple tratando de corregir las no linealidades vistas en el apartado anterior. En concreto, se han corregido las variables que miden los distintos pesos usando funciones logarítmicas y raíces cúbicas. Adicionalmente se han ensayado modelos con interacciones entre las variables. Al usar interacciones, se ha comprobado que es muy fácil conseguir modelos que produzcan sobreajuste, por lo que hay que seleccionar las interacciones con cuidado.

En este apartado hemos visto modelos que ajustaban muy bien el conjunto de *training*, pero en cambio fallaban al ser evaluados con el conjunto de test. Esto se debe a que los modelos que tienen en cuenta muchos parámetros tienen una gran varianza, perdiendo de vista las principales relaciones que hay entre las distintas variables. Esto hace que se sobreajusten al ruido presente en conjunto de datos, lo que se conoce como *overfitting*, y es el efecto contrario al producido por el sesgo que produce tener en cuenta pocas variables de entrada, o *underfitting*.

3.3. Regresión mediante el algoritmo k-NN

El algoritmo de los k vecinos más cercanos (k-Nearest Neighbours) puede usarse para regresión.

Comparando las figuras 8 y 10, podría parecer que k-NN con un único regresor ofrece resultados muy superiores a los del modelo lineal. No obstante, los resultados son muy dependientes del parámetro k . Para el ejemplo de las gráficas, con $k = 7$, el ajuste lineal tiene, en realidad, mejores resultados. En general, parece que k-NN produce más sobreajuste, al tener una diferencia mayor entre \overline{MSE}_{tra} y \overline{MSE}_{tst} . Esto puede verse claramente en la figura 11. En la misma, puede observarse que la diferencia del error en los sets de entrenamiento y de test es muy elevada, en especial con valores pequeños de k . Esta diferencia se acentúa en el caso de tener un modelo más complejo, en cuyo caso un valor bajo de k produce un sobreajuste muy importante.

Regressors	Summary of fit		Cross validation		Comments
	R^2	p -value	MSE_{tra}	MSE_{tst}	
Diameter + Height + Whole_weight + Shucked_weight + Viscera_weight + Shell_weight	0.528	< 2.2e-16	4.90	5.03	Modelo lineal eliminando las variables menos relevantes.
Length + Height + Whole_weight.log + Shucked_weight.log + Viscera_weight.log + Shell_weight.3	0.544	< 2.2e-16	4.74	4.78	Modelo lineal con preprocesamiento para eliminar las no linealidades.
poly(Length, 2) + poly(Height, 2) + poly(Whole_weight, 2) + poly(Shucked_weight, 2) + poly(Viscera_weight, 2) + poly(Shell_weight, 2)	0.565	< 2.2e-16	4.51	4.83	Modelo lineal usando polinomios de segundo grado para tratar de corregir las no linealidades. A pesar de que R^2 aumenta y \overline{MSE}_{tra} disminuye, el error en el conjunto de test crece, por lo que podemos sospechar que el modelo produce overfitting.
poly(Length, 3) + poly(Height, 3) + poly(Whole_weight, 3) + poly(Shucked_weight, 3) + poly(Viscera_weight, 3) + poly(Shell_weight, 3)	0.568	< 2.2e-16	4.48	5.68	Modelo lineal usando polinomios de grado 3 para tratar de corregir las no linealidades. El modelo produce overfitting de forma clara, lo cual es evidente si miramos la diferencia entre \overline{MSE}_{tra} y \overline{MSE}_{tst} .
Length + Height + Whole_weight.log * Shucked_weight.log * Viscera_weight * Shell_weight.3	0.574	< 2.2e-16	4.42	4.50	Modelo lineal teniendo en cuenta las interacciones entre algunas de las variables. El modelo produce mejores resultados que los anteriores, al reducir el \overline{MSE}_{tst} .

Cuadro 10: Regresiones lineales múltiples.

Regressor	k-NN ($k = 3$)		k-NN ($k = 7$)	
	MSE_{tra}	MSE_{tst}	MSE_{tra}	MSE_{tst}
Sex	9.39	9.35	8.81	8.81
Length	9.00	9.34	7.82	8.12
Diameter	9.49	9.82	7.67	7.97
Height	8.28	8.46	7.42	7.54
Whole_weight	3.96	9.34	5.17	7.79
Shucked_weight	5.70	10.56	5.97	8.88
Viscera_weight	6.75	9.87	5.97	8.22
Shell_weight	5.91	7.74	5.45	6.71
$\log(\text{Whole_weight})$	3.95	9.30	5.17	7.77
$\log(\text{Shucked_weight})$	5.72	10.57	5.96	8.91
$\log(\text{Viscera_weight})$	6.77	9.94	5.98	8.26
$\sqrt[3]{\text{Shell_weight}}$	5.90	7.65	5.42	6.64

Regressor	k-NN ($k = 15$)		Linear regression	
	MSE_{tra}	MSE_{tst}	MSE_{tra}	MSE_{tst}
Sex	9.06	9.10	8.39	8.40
Length	7.24	7.59	7.17	7.18
Diameter	7.08	7.36	6.96	6.97
Height	7.00	7.15	7.15	7.33
Whole_weight	5.98	7.29	7.36	7.37
Shucked_weight	6.72	8.20	8.55	8.57
Viscera_weight	6.31	7.56	7.76	7.77
Shell_weight	5.43	6.26	6.30	6.32
$\log(\text{Whole_weight})$	5.98	7.29	6.85	6.86
$\log(\text{Shucked_weight})$	6.72	8.21	7.69	7.70
$\log(\text{Viscera_weight})$	6.31	7.59	7.07	7.08
$\sqrt[3]{\text{Shell_weight}}$	5.41	6.23	6.02	6.03

Cuadro 11: Regresiones simples usando k-NN, con cada regresor. Comparativa con regresión lineal simple.

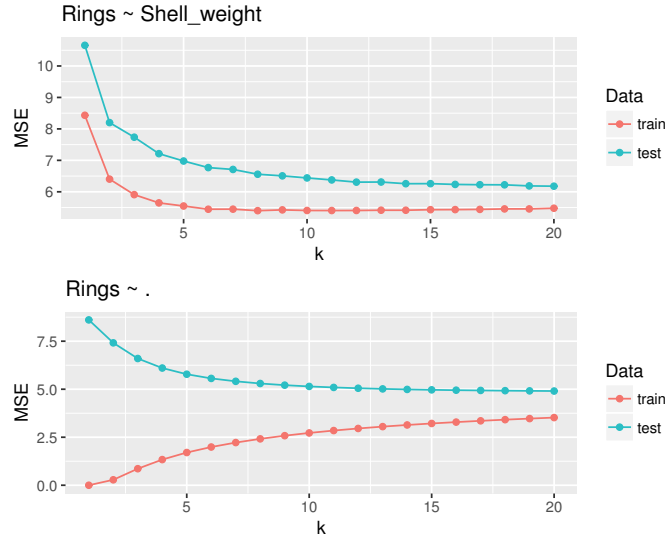


Figura 11: Comparación de diferentes modelos construidos con k-NN, cambiando el valor de k . Los errores representados son \overline{MSE}_{tra} y \overline{MSE}_{tst} .

Una de las mayores ventajas de k-NN es que se ajusta a las no linealidades sin necesidad de preprocesamiento, devolviendo resultados muy parecidos en ambos casos. Por este motivo, puede prescindirse de esta etapa al trabajar con modelos más complejos.

La figura 12 muestra una comparación del error cuadrático medio en modelos complejos. Al igual que ocurría en la regresión lineal, el modelo que tiene en cuenta las interacciones produce resultados muy buenos, siendo el que produce los mejores resultados con $k > 6$.

3.4. Comparación de los resultados

En este apartado se comparan los resultados obtenidos con el ajuste lineal, el algoritmo k-NN aplicado a la regresión. Más tarde, se comparan de forma genérica los modelos de ajuste lineal, k-NN y el modelo de regresión M5'. Para ello, se ha realizado un contraste de hipótesis sobre los resultados de la aplicación de estos métodos a los distintos conjuntos de datos usados en estas prácticas.

3.4.1. Comparación entre los modelos lineales y regresiones con k-NN para el caso estudiado

A la vista de las tablas 10 y 12, es posible decir que ambos métodos arrojan unos resultados parecidos, siempre que k sea lo suficientemente elevado. La figura 13 muestra los errores cometidos para 20 muestras de un conjunto de test, que no ha sido usado en el entrenamiento del modelo. En la misma puede apreciarse que ambos modelos cometen errores parecidos para la mayoría de los puntos. En cambio, hay puntos para los que los errores cometidos son considerablemente distintos. A pesar de estas diferencias, la media de los errores cuadráticos es muy parecida.

La figura 14 muestra una comparación de las distribuciones del error cuadrático para el mismo conjunto de test. Como puede observarse, las dos distribuciones son muy parecidas. Esto induce a pensar que ninguno de los dos métodos usados sea mejor que el otro. En efecto, un test de *Kolmogorov-Smirnov* no puede negar, con un p -value de 0.77, que las muestras hayan sido extraídas de la misma distribución. El test de *Mann-Whitney-Wilcoxon* arroja un resultado parecido, con un p -value de 0.61.

A la hora de elegir un modelo para explicar el conjunto de datos estudiado y realizar predicciones, debería usarse el modelo lineal que considera las interacciones ($\text{Rings} \sim \text{Length} + \text{Height} + \text{Whole_weight.log} * \text{Shucked_weight.log} * \text{Viscera_weight} * \text{Shell_weight.3}$). Para realizar esta elección no sólo se ha tenido en cuenta que el MSE medido ha sido ligeramente más bajo que el del correspondiente modelo k-NN, pues ya hemos visto que la diferencia no es significativa y ésta se debe, probablemente, a las particiones realizadas para la validación cruzada. En cambio, la elección

Regressors	k-NN ($k = 3$)		k-NN ($k = 7$)		k-NN ($k = 15$)	
	MSE_{tra}	MSE_{tst}	MSE_{tra}	MSE_{tst}	MSE_{tra}	MSE_{tst}
• (<i>sum of every variable</i>)	0.86	6.60	2.22	5.41	3.22	4.97
Diameter + Height + Whole_weight + Shucked_weight + Viscera_weight + Shell_weight	0.92	6.53	2.36	5.48	3.33	5.02
Length + Height + Whole_weight.log + Shucked_weight.log + Viscera_weight.log + Shell_weight.3	0.95	6.56	2.40	5.53	3.45	5.09
poly(Length, 2) + poly(Height, 2) + poly(Whole_weight, 2) + poly(Shucked_weight, 2) + poly(Viscera_weight, 2) + poly(Shell_weight, 2)	0.93	6.59	2.32	5.46	3.31	5.01
poly(Length, 3) + poly(Height, 3) + poly(Whole_weight, 3) + poly(Shucked_weight, 3) + poly(Viscera_weight, 3) + poly(Shell_weight, 3)	0.94	6.48	2.34	5.42	3.31	4.99
Length + Height + Whole_weight.log * Shucked_weight.log * Viscera_weight * Shell_weight.3	0.68	6.40	1.92	5.36	2.89	4.90

Cuadro 12: Regresiones kNN múltiples.

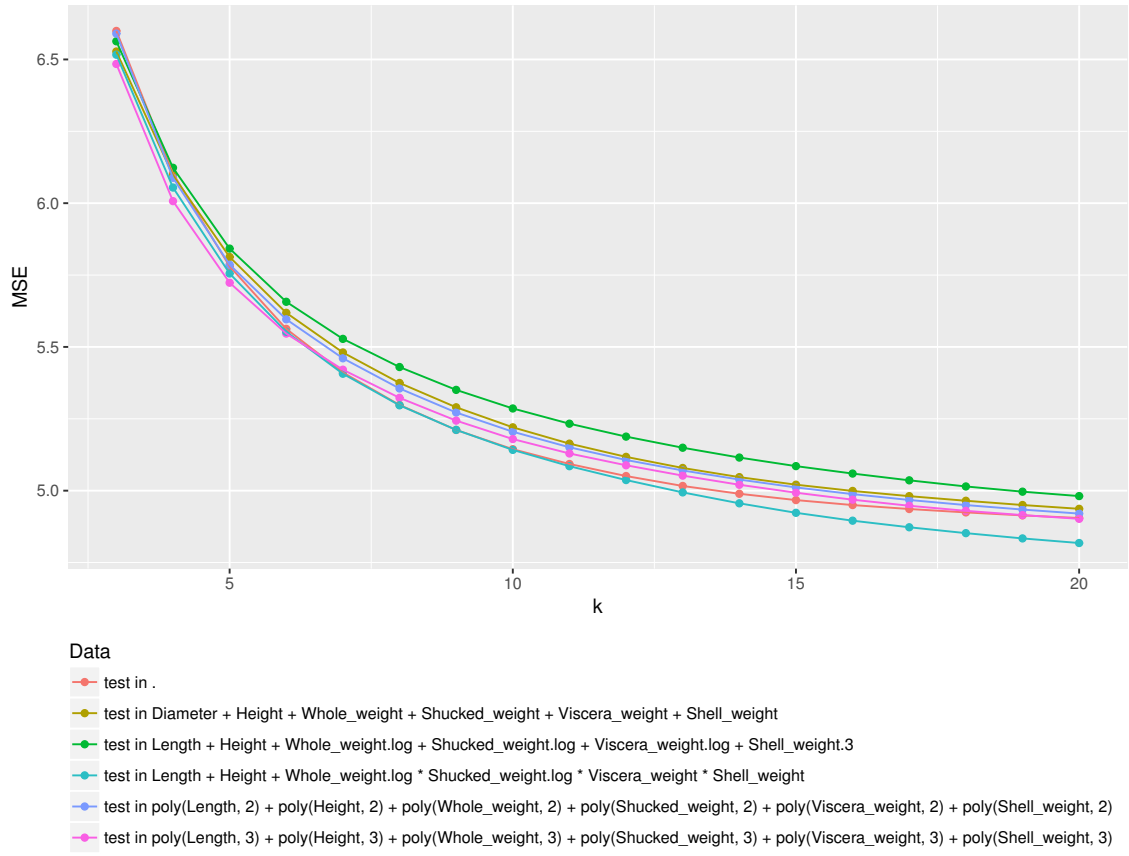


Figura 12: Comparación de diferentes modelos construidos con k-NN, cambiando el valor de k . El error representado es \overline{MSE}_{tst} .

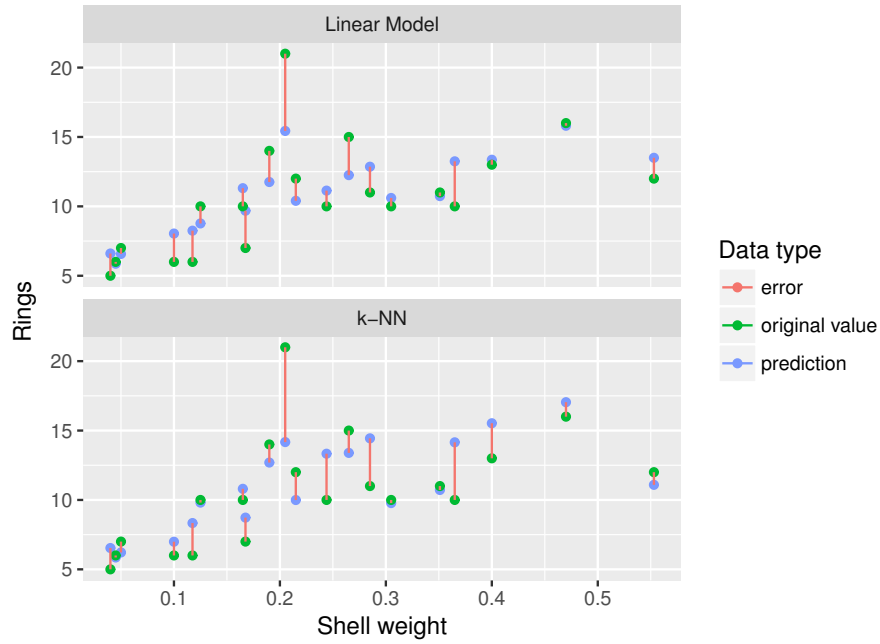


Figura 13: Comparación entre 20 datos reales sus predicciones con k-NN y LM para el modelo $\text{Rings} \sim \text{Length} + \text{Height} + \text{Whole_weight.log} * \text{Shucked_weight.log} * \text{Viscera_weight} * \text{Shell_weight}^3$ ($k = 7$). Aunque los errores tienen distintos valores para distintos puntos, la media es muy parecida.

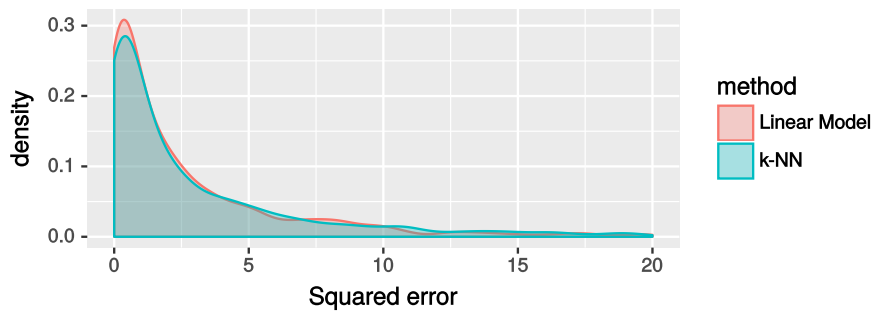


Figura 14: Comparación entre las distribuciones de error cuadrático con k-NN y LM para el modelo $\text{Rings} \sim \text{Length} + \text{Height} + \text{Whole_weight.log} * \text{Shucked_weight.log} * \text{Viscera_weight} * \text{Shell_weight}^3$ ($k = 7$).

R+	R-	p -value
78	93	0.7660

Cuadro 13: Resultado de aplicar el test de los rangos con signo de Wilcoxon a los resultados de MSE estudiados.

p -value	LM	k-NN	$1 - p$ -value	LM	k-NN
k-NN	0.67	—	k-NN	0.33	—
M5'	0.16	0.16	M5'	0.84	0.84

Cuadro 14: Resultados del test de los rangos con signo de Wilcoxon por parejas, aplicado a los resultados de *Accuracy* obtenidos por los diferentes algoritmos de regresión.

del modelo lineal se debe al hecho de que éste requiere un menor coste computacional, especialmente en la predicción, y que produce un modelo real que podría ser interpretado por un experto. También cabe destacar que dicho modelo produce menos sobreajuste que su equivalente k-NN, lo cual puede verse al comparar los errores cuadráticos medios en los conjunto de entrenamiento y de test.

3.4.2. Comparación genérica entre los modelos lineal, k-NN y M5'

Para comparar si existen diferencias significativas entre los modelos lineal, k-NN y M5', se han realizado diversos test estadísticos sobre los resultados de validación cruzada de modelos que usan todas las variables en distintos conjuntos de datos.

En primer lugar, se ha estudiado si existe diferencia significativa entre el modelo lineal y k-NN. Para ello, se ha realizado un test de los rangos con signo de Wilcoxon. El resultado, con un p -value de 0,77, es que no puede asegurarse que ninguno de los dos métodos produzca menor MSE.

Para comprobar si alguno de los tres métodos produce resultados significativamente superiores al resto se ha procedido a realizar un test de Friedman. En este caso, el p -value obtenido es de 0,0338, por lo que podemos negar la hipótesis nula que mantiene que no existe diferencia significativa entre el \overline{MSE}_{tst} medido con los métodos mencionados para el conjunto de datasets estudiado. Dicho esto, será necesario comparar los tres métodos entre sí, para lo que se hará uso de un test de los rangos con signo de Wilcoxon por parejas, cuyos resultado puede observarse en la tabla 14. Observando estos resultados, podemos decir que sólo existe un 33% de significancia de que LM produzca resultados distintos a k-NN. En cambio, tenemos un 84% de significancia de que M5' produzca resultados distintos tanto a LM como a k-NN.

4. Clasificación sobre el conjunto de datos tae

Tras haber visualizado los datos y haber explorado las posibles relaciones entre las distintas variables de entrada (ver sección 2.2), procedemos a aplicar distintos algoritmos de clasificación sobre el conjunto de datos. Para ello, se han utilizado en todo caso las particiones ya existentes, que serán también utilizadas para la validación cruzada.

4.1. Clasificación mediante el algoritmo k-NN

El algoritmo de los k vecinos más cercanos (k-NN) es un método no paramétrico que puede ser usado en clasificación. Tiene la principal ventaja de que suele conseguir muy buenas predicciones y la fase de entrenamiento es muy eficiente. Sin embargo, k-NN no produce un modelo para explicar el conjunto de datos, y la etapa de predicción es computacionalmente costosa.

Cabe destacar que, puesto que el principio básico de k-NN supone medir la distancia de cada elemento a sus k elementos más cercanos, es necesario definir una métrica y es conveniente que las variables de entrada estén normalizadas. Por este motivo, el primer paso en el procesamiento de los datos será normalizar la variable `Size`. Adicionalmente, podría procederse a normalizar el resto de variables, que, a pesar de ser categóricas, tienen un valor numérico. Esto supone prescindir de

Variable	Native	Instructor	Course	Semester	Size
Scale factor	x4	x4	x4	x1	<i>normalized</i>

Cuadro 15: Factores de normalización aplicados a las variables de entrada antes del algoritmo k-NN.



Figura 15: *Accuracy* obtenido con el algoritmo k-NN en función del parámetro k . Cada punto ha sido calculado haciendo la media de los resultados obtenidos para cada partición de validación cruzada (k -fold). La experiencia se ha repetido 9 veces para observar la variabilidad de los resultados producida por la componente aleatoria del algoritmo k-NN. La media se ha representado con una línea magenta más gruesa.

la transformación en *factors* de las variables de entrada. Aunque este paso no tiene un fundamento claro, se ha comprobado que mejora los resultados de la predicción. No obstante, tras probar con distintos factores de escalado en cada variable de entrada, se ha comprobado que la mayor mejora no se produce con una normalización¹, sino con un escalado como el de la tabla 15.

Con el fin de obtener los mejores resultados posibles, se ha realizado un estudio del rendimiento del algoritmo k-NN con distintos valores del parámetro k . Para ello, se ha calculado el *Accuracy* medio de las predicciones realizadas para cada partición de validación cruzada (k -fold) con distintos valores de k . Esta experiencia se ha repetido nueve veces, con el fin de observar la variabilidad de los resultados producida por la componente aleatoria del algoritmo k-NN. Por último, se ha calculado la media de los resultados de *Accuracy* observados, seleccionando la cifra más alta junto al valor de k que la produce. Los resultados óptimos, representados en la figura 15, se han conseguido para $k = 1$, llegando a un *Accuracy* medio de 0,621, con una desviación típica de 0,009.

4.2. Clasificación mediante análisis discriminante lineal (LDA)

LDA asume que los datos han sido extraídos de una distribución normal. Por esto, LDA no funcionará bien con datos que no sigan una distribución normal. Para comprobar si las variables de entrada siguen una distribución normal, se ha procedido a realizar un test de Shapiro-Wilk sobre las mismas. La tabla 16 muestra los resultados de este test, que deben interpretarse del siguiente modo: para cada una de las variables, debe rechazarse la hipótesis de que los datos provengan de una distribución normal. La figura 16 muestra los gráficos Q-Q de estas variables, pudiendo comprobarse que no se ajustan a la distribución normal. Tal vez las variables *Size* e *Instructor* tengan una forma más parecida a la recta esperada, cosa que también se pone de manifiesto en el p -value del test de Shapiro-Wilk, al ser éste un poco más elevado que el resto. Con estos resultados sobre la mesa, podemos asegurar que no se cumplen los requisitos establecidos por LDA, y que por tanto cabe esperar un rendimiento muy bajo del método.

Los resultados obtenidos con LDA pueden verse en la tabla 17, donde se comparan los resultados de los tres métodos estudiados.

¹ La normalización llevada a cabo se rige por la función $f(x) = \frac{x - \min(x)}{\max(x) - \min(x)}$

Variable	Statistic	<i>p</i> -value
Native	0.474	5.51e-21
Instructor	0.949	2.87e-05
Course	0.843	2.29e-11
Semester	0.421	6.89e-22
Size	0.971	2.71e-03

Cuadro 16: Test de normalidad Shapiro-Wilk aplicado a las variables de entrada del conjunto de datos *tae*.

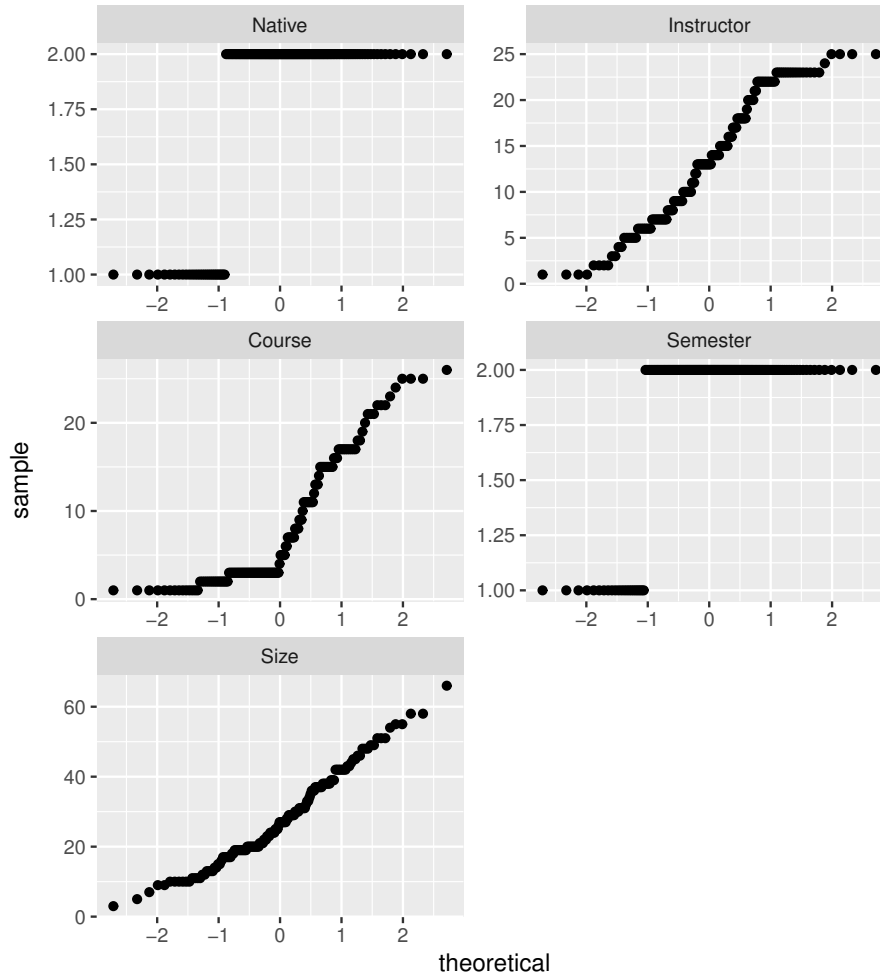


Figura 16: Gráficos Q-Q de las distintas variables de entrada.

k-fold	k-NN ($k = 1$)		LDA		QDA	
	Accuracy	Kappa	Accuracy	Kappa	Accuracy	Kappa
1	0.467	0.200	0.333	0.000	0.267	-0.100
2	0.714	0.569	0.500	0.222	0.214	-0.108
3	0.786	0.679	0.571	0.333	0.286	0.000
4	0.714	0.563	0.500	0.222	0.357	0.045
5	0.643	0.462	0.357	0.000	0.143	-0.244
6	0.603	0.408	0.357	0.000	0.286	-0.014
7	0.429	0.170	0.357	0.080	0.357	0.023
8	0.667	0.491	0.429	0.111	0.286	-0.029
9	0.571	0.349	0.429	0.188	0.286	-0.069
10	0.540	0.305	0.429	0.125	0.429	0.111
mean	0.613	0.420	0.426	0.128	0.291	-0.039

Cuadro 17: Comparación de los algoritmos de clasificación para cada uno de los k -folds.

p -value	k-NN	LDA
LDA	0.011	—
QDA	0.011	0.011

Cuadro 18: Resultados del test de los rangos con signo de Wilcoxon por parejas, aplicado a los resultados de Accuracy obtenidos por los diferentes algoritmos de clasificación.

4.3. Clasificación mediante análisis discriminante cuadrático (QDA)

QDA es un método muy parecido a LDA, que también requiere que los datos provengan de distribuciones normales. La principal ventaja de QDA sobre LDA es que el primero no requiere que los datos tengan la misma varianza ni covarianzas, puesto que éstas son estimadas de forma independiente. No obstante, el caso que nos ocupa no presenta ninguna distribución normal, por lo que no se esperan buenos resultados de este método.

Los resultados obtenidos con QDA pueden verse en la tabla 17, donde se comparan los resultados de los tres métodos estudiados.

4.4. Comparación de los algoritmos

Con el objetivo de mejorar los resultados de la clasificación, se ha probado a poner en una matriz de ocurrencias las variables Instructor y Course. A pesar de ello, las medidas de Accuracy con el algoritmo k-NN han sido muy similares a las obtenidas anteriormente. Es más, esta transformación no ha permitido el uso de los métodos LDA y QDA, ya que las matrices de ocurrencias tienen muchas columnas igual a cero, que son por tanto linealmente dependientes.

Los resultados obtenidos, tanto de Accuracy como Kappa son los mostrados en la tabla 17, cuyas medias están representadas gráficamente en la figura 17. Como puede observarse, k-NN es muy superior a los otros dos algoritmos en el caso estudiado. Para comprobar que efectivamente existen diferencias significativas entre los resultados de los distintos algoritmos, se han aplicado el test de Friedman y el test de los rangos con signo de Wilcoxon sobre la variable Accuracy. En el primero de los tests, se ha obtenido un p -value de $2.56e-05$, comprobándose que al menos uno de los métodos funciona mejor que el resto. El test de Wilcoxon revela que, en realidad, existen diferencias significativas entre los resultados de Accuracy de los tres algoritmos, al arrojar los p -value mostrados en la tabla 18.

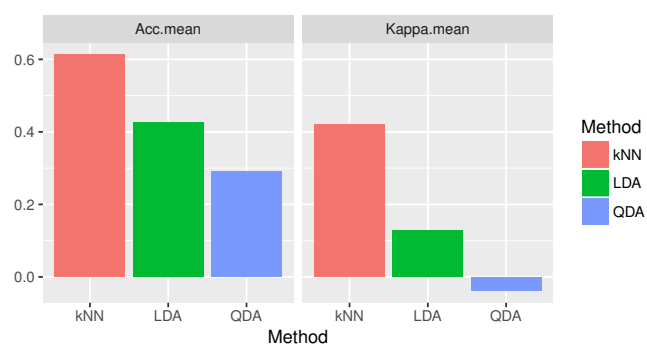


Figura 17: Comparación de los algoritmos de clasificación

Referencias

- [1] Inc. RStudio. Rstudio - free and open-source ide for r. <https://www.rstudio.com/>, 2017.
- [2] Jeffrey S Racine. Rstudio: A platform-independent ide for r and sweave. *Journal of Applied Econometrics*, 27(1):167–172, 2012.
- [3] Ross Ihaka and Robert Gentleman. R: a language for data analysis and graphics. *Journal of computational and graphical statistics*, 5(3):299–314, 1996.
- [4] Hadley Wickham. ggplot2 - a plotting system for r. <http://ggplot2.org/>, 2013.
- [5] Hadley Wickham. *ggplot2: Elegant Graphics for Data Analysis (Use R!)*. Springer, 2009.
- [6] Jesús Alcalá-Fdez, Alberto Fernández, Julián Luengo, Joaquín Derrac, Salvador García, Luciano Sánchez, and Francisco Herrera. Keel data-mining software tool: data set repository, integration of algorithms and experimental analysis framework. *Journal of Multiple-Valued Logic & Soft Computing*, 17, 2011.
- [7] Marine Resources Division and Sam Waugh. Abalone data set. <http://archive.ics.uci.edu/ml/datasets/Abalone>, 1994.
- [8] Keel-dataset - description for teaching assistant evaluation data set. <http://sci2s.ugr.es/keel/dataset.php?cod=188>, 2015.
- [9] Wei-Yin Loh and Tjen-Sien Lim. Teaching assistant evaluation data set. <http://archive.ics.uci.edu/ml/datasets/Teaching+Assistant+Evaluation>, 1997.
- [10] M. Lichman. UCI machine learning repository. <http://archive.ics.uci.edu/ml>, 2013.
- [11] Keel-dataset - description for abalone (regression version) data set. <http://sci2s.ugr.es/keel/dataset.php?cod=96>, 2015.
- [12] WJ Nash, TL Sellers, SR Talbot, AJ Cawthorn, and WB Ford. The population biology of abalone (haliotis species). *Blacklip Abalone (H. rubra) from the North Coast and Islands of Bass Strait. Sea Fisheries Division Technical Report*, 48, 1994.

A. Código fuente: Análisis de datos

A.1. Código fuente usado en la visualización de los datos abalone

```
1 #####
2 #      Introducción a la Ciencia de Datos      #
3 #      ANÁLISIS DE DATOS                      #
4 #                                              #
5 # (C) Cristian González Guerrero              #
6 #####
7
8 # Build the workspace (load variables and functions)
9 source("../regression/build-workspace.R")
10
11
12 # Get the data type and dimension
13 class(abalone)
14 dim(abalone)
15
16 # Get the data types of every field
17 lapply(abalone, class)
18
19 # Check if there are missing values
20 colSums(is.na(abalone))
21
22 # Get the statistical measurements
23 # (numeric variables) and save them into a.s.m
24 # (abalone.statistical.measurements matrix)
25 abalone2 = subset(abalone, select = -Sex)
26
27 a.s.m = matrix(
28   nrow = ncol(abalone2),
29   ncol = 8
30 )
31 colnames(a.s.m) = c(
32   "Min.",
33   "1st_Qu.",
34   "Median",
35   "Mean",
36   "3rd_Qu.",
37   "Max.",
38   "Var.",
39   "SD"
40 )
41 rownames(a.s.m) = names(abalone2)
42
43 a.s.m[,1] = sapply(abalone2, min)           # Minimum
44 a.s.m[,2] = sapply(abalone2, quantile, 0.25) # 1st Quart.
45 a.s.m[,3] = sapply(abalone2, median)        # Median
46 a.s.m[,4] = sapply(abalone2, mean)          # Mean
47 a.s.m[,5] = sapply(abalone2, quantile, 0.75) # 3rd Quart.
48 a.s.m[,6] = sapply(abalone2, max)           # Max
49 a.s.m[,7] = sapply(abalone2, var)           # Variance
50 a.s.m[,8] = sapply(abalone2, sd)            # Std. dev.
51
52 # Get the statistical measurements
53 # (categorical variables)
54 a.s.m2 = as.data.frame( table(abalone$Sex) )
55 names(a.s.m2)[1] = "Sex"
56 a.s.m2$relFreq = a.s.m2$Freq/nrow(abalone)
57
```

```

58 # Get summarized information
59 str(abalone)      # Structure
60 summary(abalone)  # Statistical measurements
61
62 # Get correlation matrix
63 cor.matrix = cor(abalone2)
64
65 # By taking a look to the file, it is clear that
66 # the output variable is Class, thus we have to
67 # classify the data by this category
68
69 # PLOTS
70
71 # Categorical variables (barplots)
72 ggplot(abalone, aes(x = Sex, fill = Sex)) +
73   geom_bar(stat = "count")
74
75 # Numeric variables Boxplots
76 meltData1 = melt::data.frame(data = abalone[,2:4])
77 meltData1$Measure = "length_(mm)"
78
79 meltData2 = melt::data.frame(data = abalone[,5:8])
80 meltData2$Measure = "weight_(grams)"
81
82 meltData = rbind(meltData1, meltData2)
83 ggplot(meltData, aes(factor(variable), value)) +
84   xlab("") + ylab("") + geom_boxplot() + coord_flip() +
85   facet_wrap(~ Measure, ncol=1, scales="free")
86
87 # Distribution of rings
88 ggplot(abalone, aes(Rings, fill = Sex, color = Sex)) +
89   geom_freqpoly(bins = 29)
90 ggplot(abalone, aes(Sex, Rings, fill = Sex)) + geom_boxplot()
91 ggplot(abalone, aes(Rings, fill = Sex, color = Sex)) +
92   stat_ecdf() + ylab("CDF")
93 ggplot(abalone, aes(Rings, fill = Sex, color = Sex)) +
94   stat_density(alpha = .8)
95
96 # Distribution of numeric variables against sex
97 # Remove outliers
98 abalone.resshaped = melt(abalone[abalone$Height<0.3,])
99 ggplot(abalone.resshaped, aes(value, colour = Sex)) +
100   geom_freqpoly(bins = 29) +
101   facet_wrap(~ variable, ncol = 2, scales="free")
102
103 # Histograms show that Male-Female distinction is not relevant,
104 # and does not characterize any of the above variables.
105 # Nonetheless, Infants' dimensions are usually lower,
106 # as is their age.
107 # Statistical test is needed to confirm this statement.
108
109
110 # Scatterplots
111 myData = melt::data.frame(
112   abalone[abalone$Height<0.3,],
113   id.vars=c("Sex", "Rings")
114 )
115 myData[myData$variable=="Height",] =
116   within(myData[myData$variable=="Height",], {
117     value = jitter(value, factor = 3)
118   })

```

```

119 )
120 ggplot(myData, aes(x = value, y = Rings)) +
121   geom_jitter(alpha = 0.03) +
122   facet_wrap( ~ variable, ncol = 2, scales = "free")
123
124
125 # Check whether the Rings distributions are the same
126 mean(abalone[abalone$Sex=="F",]$Rings)
127 summary(abalone[abalone$Sex=="F",]$Rings)
128 mean(abalone[abalone$Sex=="M",]$Rings)
129 summary(abalone[abalone$Sex=="M",]$Rings)
130
131 ks.test(
132   abalone[abalone$Sex=="M",]$Rings,
133   abalone[abalone$Sex=="F",]$Rings
134 )
135 # In the studied sample, there are more old females than
136 # old males. Actually, the distributions can be proved not
137 # to be the same.
138 # This can indicate that females usually live longer than males.
139 # This fact should be taken into account to test if there is
140 # any secondary sex characteristic, i.e. it does not suffice to
141 # test whether length distribution in females is the same than
142 # length distribution in males, it is necessary to test this
143 # taking age into account.
144 # From the minimum age of both female and male, it can be
145 # suspected that primary sex characteristics are easier
146 # detected on males.

```

End Of File 'abalone.R'

A.2. Código fuente usado en la visualización de los datos tae

```
1 #####
2 #      Introducción a la Ciencia de Datos      #
3 #      ANÁLISIS DE DATOS                      #
4 #                                              #
5 # (C) Cristian González Guerrero              #
6 #####
7
8 # Build the workspace
9 source("../classification/build-workspace.R")
10
11
12 # Get the data type and dimension
13 class(tae)
14 dim(tae)
15
16 # Get the data types of every field
17 lapply(tae, class)
18
19 # Check if there are missing values
20 colSums(is.na(tae))
21
22 # Get the statistical measurements (numeric variables)
23 attach(tae)
24 min      (Size)      # Minimum
25 quantile (Size, 0.25) # 1st Quartile
26 median   (Size)      # Median
27 mean     (Size)      # Mean
28 quantile (Size, 0.75) # 3rd Quartile
29 max      (Size)      # Max
30
31 var      (Size)      # Variance
32 sd       (Size)      # Standard deviation
33 detach(tae)
34
35 # Get summarized information
36 str(tae)      # Structure
37 summary(tae)  # Statistical measurements
38 summary(tae.original)
39
40 # By taking a look to the file, it is clear that
41 # the output variable is Class, thus we have to
42 # classify the data by this category
43
44
45 # Visualize numerical variables
46 myData = melt.data.frame(
47   data = subset(tae, select = c(Class, Size, Native, Semester)),
48   id.vars = c("Class", "Native", "Semester")
49 )
50
51 # Boxplots of numeric variables
52 ggplot(myData, aes(factor(variable), value, fill = Class)) +
53   xlab("") + ylab("Size") + geom_boxplot() + coord_flip()
54
55 # Distributions of Size
56 ggplot(myData, aes(value)) +
```

```

57   geom_freqpoly(bins = 8, aes(colour = Class)) +
58   xlab("Size") + xlim(c(0,70))
59
60 ggplot(myData, aes(value)) +
61   geom_freqpoly(bins = 8, aes(colour = Semester)) +
62   xlab("Size") + xlim(c(0,70))
63 ggplot(myData, aes(value)) +
64   geom_freqpoly(bins = 8, aes(colour = Native)) +
65   xlab("Size") + xlim(c(0,70))
66 ks.test(
67   subset(tae, subset = Semester=="Summer")$Size,
68   subset(tae, subset = Semester=="Regular")$Size
69 )
70 ks.test(
71   subset(tae, subset = Native=="English_speaker")$Size,
72   subset(tae, subset = Native=="non-English_speaker")$Size
73 )
74
75 # Visualize pseudo-numeric variables
76 myData = melt.data.frame(
77   data = subset(
78     tae.original,
79     select = c(Class, Course, Instructor)
80   ),
81   id.vars = "Class"
82 )
83 ggplot(myData, aes(factor(variable), value, fill = Class)) +
84   xlab("") + ylab("number") + geom_boxplot() + coord_flip()
85 ggplot(myData, aes(value, colour = Class)) +
86   geom_freqpoly(bins = 10) + xlab("number") +
87   facet_wrap(~ variable, ncol = 1)
88
89 #####
90 # Barplots of factors #
91 #####
92
93 ## Native & Semester
94 myData = melt.data.frame(
95   data = subset(tae, select = c(-Size, -Instructor, -Course)),
96   id.vars = "Class"
97 )
98
99
100 ggplot(myData, aes(value, fill = Class)) +
101   geom_bar(stat = "count", position = "dodge") +
102   xlab("") +
103   facet_wrap(~ variable, ncol = 2, scales = "free")
104
105 ## Class
106 ggplot(tae, aes(Class, fill = Class)) +
107   geom_bar(stat = "count")
108
109
110 ## Course
111 freq = ave(rep(1, times=nrow(tae)), tae$Course, FUN=sum)
112 myData = tae[order(freq, tae$Course, decreasing = TRUE), ]
113 myData = melt.data.frame(
114   data = subset(myData, select = c(Class, Course)),
115   id.vars = "Class"
116 )
117

```

```

118
119 ggplot(myData, aes(value, fill = Class)) +
120   geom_bar(stat = "count", position = "dodge") +
121   xlab("Course") +
122   scale_x_discrete(
123     limits = as.factor(
124       order(
125         count(tae, "Course")$freq,
126         decreasing = T
127       )
128     )
129   )
130
131
132 myData = as.data.frame(table(tae$Class, tae$Course))
133 names(myData) = c("Class", "Course", "freq")
134 freq = ave(myData$freq, myData$Course, FUN = sum)
135 myData = myData[order(freq, myData$Course, decreasing = TRUE), ]
136 ggplot(myData, aes(x = Course, y = freq, color = Class)) +
137   geom_point() +
138   geom_line(aes(group = Class)) +
139   scale_x_discrete(limits = myData$Course) +
140   ylab("Frequency")
141
142 myData$rel.freq =
143   myData$freq/freq[order(freq, myData$Course, decreasing = T)]
144 ggplot(myData, aes(x = Course, y = rel.freq, fill = Class)) +
145   geom_col(width = 2.5) +
146   scale_x_discrete(limits = myData$Course) +
147   ylab("Relative frequency")
148
149
150 ## Instructor
151 myData = as.data.frame(table(tae$Class, tae$Instructor))
152 names(myData) = c("Class", "Instructor", "freq")
153 freq = ave(myData$freq, myData$Instructor, FUN = sum)
154 myData =
155   myData[order(freq, myData$Instructor, decreasing = T), ]
156 myData[myData$freq==0,]$freq = 0.03
157 ggplot(myData, aes(x = Instructor, y = freq)) +
158   geom_col(width = 2, position = "dodge", aes(fill = Class)) +
159   scale_x_discrete(limits = myData$Instructor) +
160   ylab("Frequency")
161
162 myData[myData$freq==0.03,]$freq = 0
163 myData$rel.freq =
164   myData$freq/freq[order(freq, myData$Instructor, decreasing = T)]
165 ggplot(myData, aes(x = Instructor, y = rel.freq, fill = Class)) +
166   geom_col(width = 2.5) +
167   scale_x_discrete(limits = myData$Instructor) +
168   ylab("Relative frequency")
169
170 ggplot(tae, aes(as.integer(Instructor), color = Class)) +
171   geom_density()
172 ggplot(tae, aes(as.integer(Course), color = Class)) +
173   geom_density()

```

End Of File 'tae.R'

B. Código fuente: Regresión

B.1. Código fuente de las funciones usadas en los ejercicios de regresión

```
1 #####
2 #      Introducción a la Ciencia de Datos      #
3 #      FUNCIONES Y VARIABLES PARA REGRESIÓN    #
4 #                                              #
5 # (C) Cristian González Guerrero              #
6 #####
7
8 # Load required libraries
9 library(utils)
10 library(stats)
11 library(foreign)
12 library(ggplot2)
13 library(reshape)
14 library(kknn)
15
16 # ABALONE
17
18 # Load the dataset and provide it with the
19 # structure from Keel
20 abalone = read.csv(
21   "datasets/abalone/abalone.dat",
22   comment.char="@ "
23 )
24 names(abalone) = c(
25   "Sex",
26   "Length",
27   "Diameter",
28   "Height",
29   "Whole_weight",
30   "Shucked_weight",
31   "Viscera_weight",
32   "Shell_weight",
33   "Rings"
34 )
35 abalone$Sex = factor(
36   abalone$Sex,
37   levels = c(1, 2, 3),
38   labels = c("M", "F", "I")
39 )
40
41 # Load training data and test data
42 abalone.tra = list()
43 abalone.tst = list()
44 for (i in 1:5) {
45   for (j in 1:2) {
46     filename = paste(
47       "datasets/abalone/abalone-5-",
48       as.character(i),
49       ifelse(j==1, "tra", "tst"),
50       ".dat",
51       sep = " "
52     )
53     x = read.csv(
54       filename,
```



```

55     comment.char="@ "
56   )
57   names(x) = names(abalone)
58   x$Sex = factor(
59     x$Sex,
60     levels = c(1, 2, 3),
61     labels = c("M", "F", "I")
62   )
63   if(j==1) {
64     abalone.tra[[i]] = x
65   } else {
66     abalone.tst[[i]] = x
67   }
68 }
69 }
70
71
72 # FUNCTIONS
73
74 # Run k-fold cross validation on LM fit
75 run_lm_fold = function(
76   i,
77   tra = abalone.tra,
78   tst = abalone.tst,
79   model = Rings~.,
80   tt = "test"
81 ) {
82   x_tra = tra[[i]]
83   x_tst = tst[[i]]
84
85   if (tt != "test") {
86     test = x_tra
87   } else {
88     test = x_tst
89   }
90
91   # Perform LM fit
92   lm.fit = lm(model, x_tra)
93   output.var = as.character(model[2])
94
95   # Get MSE Error
96   yprime = predict(lm.fit, test)
97   sum(abs(test[, output.var] - yprime)^2) / length(yprime)
98 }
99
100
101 run_knn_fold = function(
102   i,
103   tra = abalone.tra,
104   tst = abalone.tst,
105   model = Rings~.,
106   tt = "test",
107   ...
108 ) {
109   x_tra = tra[[i]]
110   x_tst = tst[[i]]
111
112   if (tt != "test") {
113     test = x_tra
114   } else {
115     test = x_tst

```

```

116 }
117
118 # Perform knn fit
119 knn.fit = kknn(model, x_tra, test, ...)
120 output.var = as.character(model[2])
121
122 # Get MSE Error
123 yprime = knn.fit$fitted.values
124 sum(abs(test[, output.var] - yprime)^2) / length(yprime)
125 }
126
127 # Add non linearities to abalone training/test set
128 add.non.linearities = function(df) {
129   # Remove new fields
130   df = df[, 1:9]
131
132   df$Whole_weight.2 = (df$Whole_weight)^(1/2)
133   df$Whole_weight.3 = (df$Whole_weight)^(1/3)
134   df$Whole_weight.log = log(df$Whole_weight)
135   df$Shucked_weight.2 = (df$Shucked_weight)^(1/2)
136   df$Shucked_weight.3 = (df$Shucked_weight)^(1/3)
137   df$Shucked_weight.log = log(df$Shucked_weight)
138   df$Viscera_weight.2 = (df$Viscera_weight)^(1/2)
139   df$Viscera_weight.3 = (df$Viscera_weight)^(1/3)
140   df$Viscera_weight.log = log(df$Viscera_weight)
141   df$Shell_weight.2 = (df$Shell_weight)^(1/2)
142   df$Shell_weight.3 = (df$Shell_weight)^(1/3)
143   df$Shell_weight.log = log(df$Shell_weight)
144
145   return(df)
146 }
147
148 rm(i, j, filename, x)

```

End Of File 'build-workspace.R'

B.2. Código fuente usado en el ejercicio de regresión lineal

```
1 #####
2 #      Introducción a la Ciencia de Datos      #
3 #      REGRESIÓN LINEAR                        #
4 #                                              #
5 # (C) Cristian González Guerrero              #
6 #####
7
8 # Build the workspace
9 source("../regression/build-workspace.R")
10
11
12 # Get simple linear regression models
13 simple.linear.fit = list()
14 for (i in 1:(length(abalone)-1)) {
15   simple.linear.fit[[i]] = lm(abalone$Rings~abalone[,i])
16   names(simple.linear.fit)[i] = names(abalone)[i]
17 }
18
19 summary(simple.linear.fit$Sex)
20 summary(simple.linear.fit$Length)
21 summary(simple.linear.fit$Diameter)
22 summary(simple.linear.fit$Height)
23 summary(simple.linear.fit$Whole_weight)
24 summary(simple.linear.fit$Shucked_weight)
25 summary(simple.linear.fit$Viscera_weight)
26 summary(simple.linear.fit$Shell_weight)
27
28
29 MSE = matrix(nrow = 8, ncol = 2)
30 colnames(MSE) = c("train", "test")
31 rownames(MSE) = names(abalone)[1:8]
32 for (tt in colnames(MSE)) {
33   for (i in rownames(MSE)) {
34     MSE[i, tt] = mean(sapply(
35       1:5,
36       run_lm_fold,
37       model = Rings~eval(parse(text = i)),
38       tt = tt
39     ))
40   }
41 }
42
43 # Using the shell weight seems to be the best model,
44 # since it produces less error.
45
46 # Scatterplots
47 myData = melt.data.frame(
48   abalone[abalone$Height<0.3,],
49   id.vars=c("Sex", "Rings")
50 )
51 myData[myData$variable=="Height",] =
52   within(myData[myData$variable=="Height",], {
53     value = jitter(value, factor = 3)
54   })
55
56 ggplot(myData, aes(x = value, y = Rings)) +
```

```

57 geom_jitter(alpha = 0.03) +
58 geom_smooth(method = lm) +
59 facet_wrap( ~ variable, ncol = 2, scales = "free") +
60 xlab("")
61
62
63 # Analysis of preprocessing
64 abalone = add.non.linearities(abalone)
65 abalone.tra = lapply(abalone.tra, add.non.linearities)
66 abalone.tst = lapply(abalone.tst, add.non.linearities)
67
68
69 ## Calculate correlations
70 Y.vars = grepl("Rings", names(abalone))
71 X.vars = grepl("_weight", names(abalone))
72 cor.Shucked_weight = cor(
73   subset(abalone, select = X.vars),
74   subset(abalone, select = Y.vars)
75 )
76
77 ## New models, based on new variables
78 simple.linear.fit = list()
79 for (i in 1:(length(abalone))) {
80   simple.linear.fit[[i]] = lm(abalone$Rings~abalone[,i])
81   names(simple.linear.fit)[i] = names(abalone)[i]
82 }
83
84 summary(simple.linear.fit$Whole_weight.log)
85 summary(simple.linear.fit$Shucked_weight.2)
86 summary(simple.linear.fit$Shucked_weight.3)
87 summary(simple.linear.fit$Shucked_weight.log)
88 summary(simple.linear.fit$Viscera_weight.log)
89 summary(simple.linear.fit$Shell_weight.3)
90 summary(simple.linear.fit$Shell_weight.log)
91
92 MSE = matrix(nrow = ncol(abalone), ncol = 2)
93 colnames(MSE) = c("train", "test")
94 rownames(MSE) = names(abalone)
95 for (tt in colnames(MSE)) {
96   for (i in rownames(MSE)) {
97     MSE[i, tt] = mean(sapply(
98       1:5,
99       run_lm_fold,
100       model = Rings~eval(parse(text = i)),
101       tt = tt
102     ))
103   }
104 }
105
106
107 # Plot
108 myData = melt.data.frame(
109   abalone[abalone$Height<0.3,],
110   id.vars=c("Sex", names(abalone)[Y.vars])
111 )
112
113 selected.fields = c(
114   "Whole_weight",
115   "log(Whole_weight)",
116   "Shucked_weight",
117   "log(Shucked_weight)",

```

```

118   "Viscera_weight",
119   "log(Viscera_weight)",
120   "Shell_weight",
121   "(Shell_weight)^(1/3)"
122 )
123 names(selected.fields) = c(
124   "Whole_weight",
125   "Whole_weight.log",
126   "Shucked_weight",
127   "Shucked_weight.log",
128   "Viscera_weight",
129   "Viscera_weight.log",
130   "Shell_weight",
131   "Shell_weight.3"
132 )
133
134 ggplot(
135   subset(myData, variable %in% names(selected.fields)),
136   aes(x = value, y = Rings)
137 ) +
138   geom_jitter(alpha = 0.03) + geom_smooth(method = lm) +
139   facet_wrap(
140     ~ variable,
141     ncol = 2,
142     scales = "free",
143     labeller = as_labeller(selected.fields)
144   ) +
145   xlab("")
146
147 # Linear regression using multiple variables
148 ## Only linear models
149 abalone = abalone[,1:9]
150
151 myFit = lm(Rings~., abalone)
152 summary(myFit)
153
154 myModel = Rings ~
155   Diameter +
156   Height +
157   Whole_weight +
158   Shucked_weight +
159   Viscera_weight +
160   Shell_weight
161 myFit = lm(myModel, abalone)
162 summary(myFit)
163
164
165 MSE = matrix(nrow = 1, ncol = 2)
166 colnames(MSE) = c("train", "test")
167 for (tt in colnames(MSE)) {
168   MSE[1, tt] = mean(sapply(
169     1:5,
170     run_lm_fold,
171     model = myModel,
172     tt = tt
173   ))
174 }
175
176 # Add non-linearities
177 abalone = add.non.linearities(abalone)
178

```

```

179 myModel = Rings ~
180   Length +
181   Height +
182   Whole_weight.log +
183   Shucked_weight.log +
184   Viscera_weight.log +
185   Shell_weight.3
186 myFit = lm(Rings~., abalone)
187 summary(myFit)
188
189 myFit = lm(myModel, abalone)
190 summary(myFit)
191
192 MSE = matrix(nrow = 1, ncol = 2)
193 colnames(MSE) = c("train", "test")
194 for (tt in colnames(MSE)) {
195   MSE[1, tt] = mean(sapply(
196     1:5,
197     run_lm_fold,
198     model = myModel,
199     tt = tt
200   ))
201 }
202
203 myModel = Rings ~
204   poly(Length, 2) +
205   poly(Height, 2) +
206   poly(Whole_weight, 2) +
207   poly(Shucked_weight, 2) +
208   poly(Viscera_weight, 2) +
209   poly(Shell_weight, 2)
210 myFit = lm(myModel, abalone)
211 summary(myFit)
212
213 MSE = matrix(nrow = 1, ncol = 2)
214 colnames(MSE) = c("train", "test")
215 for (tt in colnames(MSE)) {
216   MSE[1, tt] = mean(sapply(
217     1:5,
218     run_lm_fold,
219     model = myModel,
220     tt = tt
221   ))
222 }
223
224 # Add interactions
225 myModel = Rings ~
226   Length +
227   Height +
228   Whole_weight.log *
229   Shucked_weight.log *
230   Viscera_weight *
231   Shell_weight.3
232
233 myFit = lm(myModel, abalone)
234 summary(myFit)
235
236 MSE = matrix(nrow = 1, ncol = 2)
237 colnames(MSE) = c("train", "test")
238 for (tt in colnames(MSE)) {
239   MSE[1, tt] =

```

```
240     mean(sapply(1:5, run_lm_fold, model = myModel, tt = tt))
241 }
```

End Of File 'linear-regression.R'

B.3. Código fuente usado en el ejercicio de regresión k-NN

```
1 #####
2 #      Introducción a la Ciencia de Datos      #
3 #      REGRESIÓN k-NN                          #
4 #                                              #
5 # (C) Cristian González Guerrero              #
6 #####
7
8 # Build the workspace
9 source("../regression/build-workspace.R")
10
11 abalone = abalone[abalone$Height<0.3,]
12
13 # Get simple knn models
14 simple.knn.fit = list()
15 for (i in 1:(length(abalone)-1)) {
16   simple.knn.fit[[i]] =
17     knn(abalone$Rings~abalone[,i], abalone, abalone)
18   names(simple.knn.fit)[i] = names(abalone)[i]
19 }
20
21
22 ## Scatterplots
23 myData1 = melt::data.frame(
24   abalone,
25   id.vars=c("Sex", "Rings")
26 )
27
28 myData2 = data.frame()
29 for (variable in names(simple.knn.fit)[2:8]) {
30   Prediction = simple.knn.fit[[variable]]$fitted.values
31   myData2 = rbind(myData2, data.frame(Prediction))
32 }
33
34 myData = cbind(myData1, myData2)
35
36 Data = "Original_data"
37 b = "Prediction"
38 colorPalette <- c("#000000", "#56B4E9")
39
40 ggplot(myData) +
41   geom_point(
42     aes(x = value, y = Rings, color = Data),
43     alpha = 0.03
44   ) +
45   geom_point(
46     aes(x = value, y = Prediction, color = b),
47     alpha = 0.03
48   ) +
49   facet_wrap(~ variable, ncol = 2, scales = "free") +
50   xlab("") +
51   scale_colour_manual(values=colorPalette) +
52   guides(colour = guide_legend(override.aes = list(alpha = 1)))
53
54 ## MSE
55 MSE = matrix(nrow = 8, ncol = 2)
56 colnames(MSE) = c("train", "test")
```



```

57 rownames(MSE) = names(abalone)[1:8]
58 for (tt in colnames(MSE)) {
59   for (i in rownames(MSE)) {
60     MSE[i, tt] = mean(sapply(
61       1:5,
62       run_knn_fold,
63       model = Rings~eval(parse(text = i)),
64       tt = tt,
65       k = 7 # Modify k
66     ))
67   }
68 }
69
70 ## Comparative k
71 k = 1:20
72 MSE.Shell_weight = data.frame(k)
73 for (tt in c("train", "test")) {
74   for (i in k) {
75     MSE.Shell_weight[i,tt] = mean(sapply(
76       1:5,
77       run_knn_fold,
78       model = Rings~Shell_weight,
79       tt = tt,
80       k = i
81     ))
82   }
83 }
84
85 Data = "test"
86 b = "train"
87 ggplot(MSE.Shell_weight) +
88   geom_point(aes(x = k, y = test, color = Data)) +
89   geom_line(aes(x = k, y = test, color = Data)) +
90   geom_point(aes(x = k, y = train, color = b)) +
91   geom_line(aes(x = k, y = train, color = b)) +
92   ylab("MSE") +
93   scale_colour_manual(values=colorPalette)
94
95
96
97 # Non-linearities
98 abalone = add.non.linearities(abalone)
99 abalone.tra = lapply(abalone.tra, add.non.linearities)
100 abalone.tst = lapply(abalone.tst, add.non.linearities)
101
102 MSE = matrix(nrow = ncol(abalone), ncol = 2)
103 colnames(MSE) = c("train", "test")
104 rownames(MSE) = names(abalone)
105 for (tt in colnames(MSE)) {
106   for (i in rownames(MSE)) {
107     MSE[i, tt] = mean(sapply(
108       1:5,
109       run_knn_fold,
110       model = Rings~eval(parse(text = i)),
111       tt = tt,
112       k = 7 # Modify k
113     ))
114   }
115 }
116
117

```

```

118 # Multiple variables
119 abalone = abalone[,1:9]
120 abalone.tra = lapply(abalone.tra, function(df){return(df[,1:9])})
121 abalone.tst = lapply(abalone.tst, function(df){return(df[,1:9])})
122 MSE.record = data.frame()
123
124 knn.MSE.data.frame = function(
125   model = Rigns~.,
126   k = 1:20, i = 1:5,
127   tt = c("train", "test")
128 ) {
129   MSE = data.frame(k)
130   tt_ = tt
131   i_ = i
132   for (tt in tt_) {
133     for (i in k) {
134       MSE[i, tt] = mean(sapply(
135         i_,
136         run_knn_fold,
137         model = model,
138         tt = tt,
139         k = i # Modify k
140       ))
141     }
142   }
143   MSE$model = as.character(myModel)[3]
144
145   return(MSE)
146 }
147
148 plot.MSE.data.frame = function(MSE, model = "") {
149   model = as.character(myModel)
150   if(length(model) == 3) {
151     model = paste(model[2], model[1], model[3])
152   }
153   myData = melt.data.frame(
154     MSE,
155     id.vars = c("k", "model"),
156     variable_name = "Data"
157   )
158
159   ggplot(myData, aes(x = k, y = value, color = Data)) +
160     geom_point() + geom_line() + ylab("MSE") +
161     ggtitle(model)
162 }
163
164 # Simple regression
165 myModel = Rings~Shell_weight
166
167 MSE = knn.MSE.data.frame(myModel)
168 MSE.record = rbind(MSE.record, MSE)
169 plot.MSE.data.frame(MSE, myModel)
170
171 # Every variable
172 myModel = Rings~.
173
174 MSE = knn.MSE.data.frame(myModel)
175 MSE.record = rbind(MSE.record, MSE)
176 plot.MSE.data.frame(MSE, myModel)
177
178

```

```

179 ## No interactions
180 myModel = Rings ~
181   Diameter +
182   Height +
183   Whole_weight +
184   Shucked_weight +
185   Viscera_weight +
186   Shell_weight
187
188 MSE = knn.MSE.data.frame(myModel)
189 MSE.record = rbind(MSE.record, MSE)
190 plot.MSE.data.frame(MSE)
191
192
193 ## Non-linearities
194 abalone = add.non.linearities(abalone)
195 abalone.tra = lapply(abalone.tra, add.non.linearities)
196 abalone.tst = lapply(abalone.tst, add.non.linearities)
197
198 myModel = Rings ~
199   Length +
200   Height +
201   Whole_weight.log +
202   Shucked_weight.log +
203   Viscera_weight.log +
204   Shell_weight.3
205
206 MSE = knn.MSE.data.frame(myModel)
207 MSE.record = rbind(MSE.record, MSE)
208 plot.MSE.data.frame(MSE)
209
210
211 ## Non-linearities (poly)
212 myModel = Rings ~
213   poly(Length, 2) +
214   poly(Height, 2) +
215   poly(Whole_weight, 2) +
216   poly(Shucked_weight, 2) +
217   poly(Viscera_weight, 2) +
218   poly(Shell_weight, 2)
219
220 MSE = knn.MSE.data.frame(myModel)
221 MSE.record = rbind(MSE.record, MSE)
222 plot.MSE.data.frame(MSE)
223
224 #--
225 myModel = Rings ~
226   poly(Length, 3) +
227   poly(Height, 3) +
228   poly(Whole_weight, 3) +
229   poly(Shucked_weight, 3) +
230   poly(Viscera_weight, 3) +
231   poly(Shell_weight, 3)
232
233 MSE = knn.MSE.data.frame(myModel)
234 MSE.record = rbind(MSE.record, MSE)
235 plot.MSE.data.frame(MSE)
236
237
238 ## Interactions
239 myModel = Rings ~

```

```

240 Length +
241 Height +
242 Whole_weight.log *
243 Shucked_weight.log *
244 Viscera_weight *
245 Shell_weight.3
246
247 MSE = knn.MSE.data.frame(myModel)
248 MSE.record = rbind(MSE.record, MSE)
249 plot.MSE.data.frame(MSE)
250
251 ## Compare k
252 myData = subset(MSE.record, select = -train)
253 myData = myData[myData$k>2,]
254 myData = melt.data.frame(
255   myData,
256   id.vars = c("k", "model"),
257   variable_name = "Data"
258 )
259 myData$Data = paste(myData$Data, "in", myData$model)
260
261 ggplot(myData, aes(x = k, y = value, color = Data)) +
262   geom_point() + geom_line() + ylab("MSE") +
263   theme(legend.position = "bottom", legend.direction = "vertical")

```

End Of File 'knn-regression.R'

B.4. Código fuente de la comparación entre los distintos algoritmos de regresión

```
1 #####
2 #      Introducción a la Ciencia de Datos      #
3 #      COMPARATIVA DE ALGORITMOS                #
4 #                                                #
5 # (C) Cristian González Guerrero                #
6 #####
7
8 # Build the workspace
9 source("../regression/build-workspace.R")
10
11 abalone.tra = lapply(abalone.tra, add.non.linearities)
12 abalone.tst = lapply(abalone.tst, add.non.linearities)
13
14
15 # Comparison of data points
16
17 myModel = Rings ~
18   Length +
19   Height +
20   Whole_weight.log *
21   Shucked_weight.log *
22   Viscera_weight *
23   Shell_weight.3
24
25 myTrain = abalone.tra[[1]]
26 myTest  = abalone.tst[[1]][sample(1:nrow(abalone.tst[[1]]), 20), ]
27
28 ## Linear regression
29 # Plot some points of the last model
30 myFit = lm(myModel, myTrain)
31 myResult = predict(myFit, myTest)
32
33 myData1 = myTest
34 myData1$Rings = myResult
35 myData1 = cbind(
36   myData1,
37   data.type = "prediction",
38   method = "Linear_Model"
39 )
40 myData1 = rbind(
41   myData1,
42   cbind(
43     myTest,
44     data.type = "original_value",
45     method = "Linear_Model"
46   )
47 )
48 myData1 = melt.data.frame(
49   data = myData1,
50   id.vars = c("Shell_weight", "data.type", "method"),
51   measure.vars = "Rings"
52 )
53
54 ## k-NN
55 # Plot some points of the last model
56 myResult = kknm(myModel, myTrain, myTest, k = 7)
```

```

57
58 myData2 = myTest
59 myData2$Rings = myResult$fitted.values
60 myData2 = cbind(
61   myData2,
62   data.type = "prediction",
63   method = "k-NN"
64 )
65 myData2 = rbind(
66   myData2,
67   cbind(
68     myTest,
69     data.type = "original_value",
70     method = "k-NN"
71   )
72 )
73 myData2 = melt.data.frame(
74   data = myData2,
75   id.vars = c("Shell_weight", "data.type", "method"),
76   measure.vars = "Rings"
77 )
78
79
80 ## Plot
81 myData = rbind(myData1, myData2)
82
83 ggplot(
84   myData,
85   aes(x = Shell_weight, y = value, color = data.type)
86 ) +
87   geom_point() +
88   ylab("Rings") +
89   xlab("Shell_weight") +
90   geom_line(aes(group = Shell_weight, color = "error")) +
91   labs(color = "Data_type") +
92   facet_wrap( ~ method, ncol = 1)
93
94
95 # Error distribution
96 myTest = abalone.tst[[1]]
97 lmErrors =
98   myTest$Rings - predict(myFit, myTest)
99 knnErrors =
100   myTest$Rings - knn(myModel, myTrain, myTest)$fitted.values
101
102 ### Tests
103 wilcox.test(lmErrors^2, knnErrors^2, paired = FALSE)
104 ks.test(lmErrors^2, knnErrors^2)
105
106 ### Plot
107 myData1 = data.frame(method = "Linear_Model", Err = lmErrors)
108 myData2 = data.frame(method = "k-NN", Err = knnErrors)
109 myData = melt.data.frame(
110   data = rbind(myData1, myData2),
111   id.vars = "method"
112 )
113 ggplot(myData, aes(x=value^2)) +
114   geom_density(aes(color = method, fill = method), alpha = 0.3) +
115   xlab("Squared_error") + xlim(0,20)
116
117

```

```

118 # General algorithm comparison (using MSE from every database)
119 ## Read data
120 testResults = read.csv("../regression/regr_test_alumnos.csv")
121 testTable = testResults[, 2:ncol(testResults)]
122 rownames(testTable) = testResults[,1]
123
124 trainResults = read.csv("../regression/regr_train_alumnos.csv")
125 trainTable = trainResults[, 2:ncol(testResults)]
126 rownames(trainTable) = trainResults[,1]
127
128 ## Comparison between LM and kNN
129 ## (kNN is assumed to provide better results)
130 difs = (testTable[,1] - testTable[,2])/testTable[,1]
131
132 wilc_1_2 = cbind(
133     ifelse(difs<0, abs(difs)+0.1, 0.1),
134     ifelse(difs>0, abs(difs)+0.1, 0.1)
135 )
136 colnames(wilc_1_2) = colnames(testTable)[1:2]
137
138 ### Wilcoxon tests
139 LMvsKNNtst = wilcox.test(
140     wilc_1_2[,1],
141     wilc_1_2[,2],
142     alternative = "two.sided",
143     paired = TRUE
144 )
145
146 Rplus = LMvsKNNtst$statistic
147 pvalue = LMvsKNNtst$p.value
148
149 LMvsKNNtst = wilcox.test(
150     wilc_1_2[,2],
151     wilc_1_2[,1],
152     alternative = "two.sided",
153     paired = TRUE
154 )
155
156 Rminus = LMvsKNNtst$statistic
157
158 WilcoxonTestOutput = cbind(Rplus, Rminus, pvalue)
159
160
161 ## Comparison amongst LM, kNN and M5'
162 friedman.test(as.matrix(testTable))
163
164 #### p-value < 0.05 => significative difference exist
165
166 groups = rep(1:ncol(testTable), each = nrow(testTable))
167 pairwise.wilcox.test(
168     as.matrix(testTable),
169     groups,
170     p.adjust.method = "holm",
171     paired = TRUE
172 )
173
174
175
176 #### M5' seems to work better (confidence: 84%)
177 #### LM and kNN seems not to present any difference

```

End Of File 'comparison.R'

C. Código fuente: Clasificación

C.1. Código fuente de de las funciones usadas en el ejercicio de clasificación

```
1 #####
2 #      Introducción a la Ciencia de Datos      #
3 #      FUNCIONES Y VARIABLES PARA CLASIFICACIÓN  #
4 #                                                                 #
5 # (C) Cristian González Guerrero                #
6 #####
7
8 # Load required libraries
9 library(utils)
10 library(stats)
11 library(foreign)
12 library(ggplot2)
13 library(reshape)
14 library(GGally)
15 library(class)
16 library(caret)
17 library(MASS)
18 library(plyr)
19 library(ISLR)
20 library(klaR)
21
22 #TAE
23
24 # Load the dataset and provide it with the
25 # structure from Keel
26 tae = read.csv(
27   "datasets/tae/tae.dat",
28   comment.char="@")
29 )
30 names(tae) = c(
31   "Native",
32   "Instructor",
33   "Course",
34   "Semester",
35   "Size",
36   "Class"
37 )
38
39
40 tae$Class      = factor(
41   tae$Class,
42   levels = c(1,2,3),
43   labels = c("low", "medium", "high")
44 )
45 tae.original = tae
46 tae$Native    = factor(
47   tae$Native,
48   levels = c(1,2),
49   labels = c(
50     "English_speaker",
51     "non-English_speaker"
52   )
53 )
54 tae$Semester  = factor(
```



```

55   tae$Semester,
56   levels = c(1,2),
57   labels = c("Summer", "Regular")
58 )
59 tae$Instructor = factor(tae$Instructor)
60 tae$Course      = factor(tae$Course)
61
62
63 # Load traning data and test data
64 tae.tra = list()
65 tae.tra.cl = list()
66 tae.tra.dt = list()
67 tae.tst = list()
68 tae.tst.cl = list()
69 tae.tst.dt = list()
70 for (i in 1:10) {
71   for (j in 1:2) {
72     filename = paste(
73       "datasets/tae/tae-10-",
74       as.character(i),
75       ifelse(j==1, "tra", "tst"),
76       ".dat",
77       sep = ""
78     )
79     x = read.csv(
80       filename,
81       comment.char="@ "
82     )
83     names(x) = names(tae)
84     # Everything except the Class will be
85     # treated as numeric for classification
86     x$Class = factor(
87       x$Class,
88       levels = c(1,2,3),
89       labels = c("low", "medium", "high")
90     )
91
92     if(j==1) {
93       tae.tra[[i]] = x
94       tae.tra.dt[[i]] = subset(x, select = -Class)
95       tae.tra.cl[[i]] = x$Class
96     } else {
97       tae.tst[[i]] = x
98       tae.tst.dt[[i]] = subset(x, select = -Class)
99       tae.tst.cl[[i]] = x$Class
100    }
101  }
102 }
103
104
105 # FUNCTIONS
106 normalize = function(x) {
107   return((x-min(x)) / (max(x)-min(x)))
108 }
109
110 normalize.tae = function(df) {
111   df$Instructor = df$Instructor*4
112   df$Course      = df$Course*4
113   df$Native      = df$Native*4
114   df$Semester    = df$Semester
115   df$Size        = normalize(df$Size)

```

```

116
117     return(df)
118 }
119
120 # Accuracy function
121 accuracy = function(ConfusionMatrix) {
122     return(sum(diag(ConfusionMatrix))/sum(ConfusionMatrix))
123
124     meanAccuracy      = 0
125     for (k in 1:ncol(ConfusionMatrix)) {
126         TP = ConfusionMatrix[k,k]
127         FP = sum(ConfusionMatrix[k,]) - TP
128         FN = sum(ConfusionMatrix[,k]) - TP
129         TN = sum(ConfusionMatrix[-k,-k])
130         Accuracy = (TP+TN)/(TP+TN+FP+FN)
131
132         meanAccuracy = meanAccuracy + Accuracy/ncol(ConfusionMatrix)
133     }
134     return(meanAccuracy)
135 }
136
137 set.seed(1)
138
139 rm(i, j, filename, x)

```

End Of File 'build-workspace.R'

C.2. Código fuente usado en el ejercicio de clasificación

```
1 #####
2 #      Introducción a la Ciencia de Datos      #
3 #                  CLASIFICACIÓN                  #
4 #                  #                              #
5 # (C) Cristian González Guerrero                #
6 #####
7
8 # Build the workspace
9 source("../classification/build-workspace.R")
10
11 # Some info about the classes
12 round(prop.table(table(tae$Class)) * 100, digits = 1)
13
14 # Normalize numeric data
15 tae$Size = normalize(tae$Size)
16 tae.tra.dt = lapply(tae.tra.dt, normalize.tae)
17 tae.tst.dt = lapply(tae.tst.dt, normalize.tae)
18
19
20
21 #####
22 # k Nearest Neighbours (k-NN) #
23 #####
24
25 # Apply k-NN algorithm to the first set
26 tae_test_pred = knn(
27   train = tae.tra.dt[[1]], # Dataframe
28   test  = tae.tst.dt[[1]], # Dataframe
29   cl     = tae.tra.cl[[1]], # Vector (factor)
30   k      = 21
31 )
32
33 # See the results on a contingency table (confusion matrix)
34 table(tae_test_pred, tae.tst.cl[[1]])
35
36 # Apply k-NN algorithm to every dataset and check the
37 # overall accuracy for several values of k. Repeat this process
38 # m times, to get the random part working.
39 k.max = 35
40 Acc = data.frame(k = 1:k.max)
41 for (m in 1:9) {
42   for (k in 1:k.max) {
43     acc = rep(0, 10)
44     for (i in 1:10) {
45       tae_test_pred = knn(
46         train = tae.tra.dt[[i]], # Dataframe
47         test  = tae.tst.dt[[i]], # Dataframe
48         cl     = tae.tra.cl[[i]], # Vector (factor)
49         k      = k
50       )
51       acc[i] = accuracy(
52         table(tae_test_pred, tae.tst.cl[[i]])
53       )
54     }
55     Acc[k, paste("Acc", m, sep = "")] = mean(acc)
56   }
}
```

```

57 }
58
59 Acc$Acc_mean = rowMeans(Acc[, -1])
60
61 myData = melt::data.frame(
62   data = Acc,
63   id.vars = "k"
64 )
65 myData$lineSize = 1
66 myData[myData$variable=="Acc_mean",]$lineSize = 2
67
68 max(Acc$Acc_mean)
69 which.max(Acc$Acc_mean)
70
71 ggplot(
72   myData,
73   aes(x = k, y = value, color = variable, size = lineSize)
74 ) +
75   geom_line() + scale_size(range = c(0.5, 2), guide="none") +
76   ylab("Accuracy")
77 ## It seems that k=1 performs the best, with Acc = 0.61
78 # This plot could also be done with training data
79 # in order to test for overfitting
80
81
82 # Fill in the table in the document
83 Acc.knn = matrix(ncol = 10, nrow = 10)
84 rownames(Acc.knn) = paste("K", 1:10, sep = "")
85 colnames(Acc.knn) =
86   c(paste("Acc", 1:9, sep = ""), "Acc_mean")
87 Kappa.knn = matrix(ncol = 10, nrow = 10)
88 rownames(Kappa.knn) = paste("K", 1:10, sep = "")
89 colnames(Kappa.knn) =
90   c(paste("Kappa", 1:9, sep = ""), "Kappa_mean")
91 for (m in 1:9) {
92   myAcc = rep(0, 10)
93   myKappa = rep(0, 10)
94   for (i in 1:10) {
95     tae_test_pred = knn(
96       train = tae.tra.dt[[i]], # Dataframe
97       test = tae.tst.dt[[i]], # Dataframe
98       cl = tae.tra.cl[[i]], # Vector (factor)
99       k = 1
100     )
101     cm = confusionMatrix(
102       table(tae_test_pred, tae.tst.cl[[i]])
103     )
104     myAcc[i] = cm$overall[1]
105     myKappa[i] = cm$overall[2]
106   }
107   Acc.knn[,m] = myAcc
108   Kappa.knn[,m] = myKappa
109 }
110 Acc.knn[, 10] = rowMeans(Acc.knn[, 1:9])
111 Kappa.knn[, 10] = rowMeans(Kappa.knn[, 1:9])
112
113 #####
114 # Linear Discriminant Analysis (LDA) #
115 #####
116
117

```

```

118 # Check for normality of (pseudo-)numeric variables
119 sapply(tae.original[1:5], shapiro.test)
120
121 myData = melt.data.frame(
122   data = tae.original,
123   measure.vars = c(
124     "Native",
125     "Instructor",
126     "Course",
127     "Semester",
128     "Size"
129   )
130 )
131 ggplot(myData) +
132   stat_qq(aes(sample = value)) +
133   facet_wrap(~variable, ncol = 2, scales = "free")
134
135 # LDA fit
136 lda.fit = lda(
137   Class ~ .,
138   data = tae.tra[[1]]
139 )
140
141 lda.pred = predict(lda.fit, tae.tst.dt[[1]])
142
143 table(lda.pred$class, tae.tst.cl[[1]])
144
145 # Do it for every dataset in k-fold CV
146 acc = rep(0, 10)
147 for (i in 1:10) {
148   lda.fit = lda(
149     Class ~ .,
150     data = tae.tra[[i]]
151   )
152   lda.pred = predict(lda.fit, tae.tst.dt[[i]])
153
154   acc[i] = accuracy(
155     table(lda.pred$class, tae.tst.cl[[i]])
156   )
157 }
158 Acc.LDA = data.frame(fold = 1:10, Acc = acc)
159
160 mean(Acc.LDA$Acc)
161
162 # LDA performs poorly if every variable is considered.
163 # However, there are not real numeric variables, only
164 # ordered factors, with no normal distribution.
165 # Other variables should be computed from available
166 # information (and pray for them to have a quasi-normal
167 # distribution). Otherwise, LDA will not perform correctly.
168 # Actually, mean(acc) equals 0.348 if the model considers
169 # solely the "numeric" variables, i.e. Course+Instructor+Size,
170 # meaning a loss of less than 5% in accuracy.
171
172
173 # Fill in the table in the document
174 Acc.lda = matrix(ncol = 10, nrow = 10)
175 rownames(Acc.lda) = paste("K", 1:10, sep = "")
176 colnames(Acc.lda) =
177   c(paste("Acc", 1:9, sep = ""), "Acc_mean")
178 Kappa.lda = matrix(ncol = 10, nrow = 10)

```

```

179 rownames(Kappa.lda) = paste("K", 1:10, sep = "")
180 colnames(Kappa.lda) =
181   c(paste("Kappa", 1:9, sep = ""), "Kappa_mean")
182 for (m in 1:9) {
183   myAcc = rep(0, 10)
184   myKappa = rep(0, 10)
185   for (i in 1:10) {
186     lda.fit = lda(
187       Class ~ .,
188       data = tae.tra[[i]]
189     )
190     lda.pred = predict(lda.fit, tae.tst.dt[[i]])
191
192     cm = confusionMatrix(
193       table(lda.pred$class, tae.tst.cl[[i]])
194     )
195     myAcc[i] = cm$overall[1]
196     myKappa[i] = cm$overall[2]
197   }
198   Acc.lda[,m] = myAcc
199   Kappa.lda[,m] = myKappa
200 }
201 Acc.lda[,10] = rowMeans(Acc.lda[,1:9])
202 Kappa.lda[,10] = rowMeans(Kappa.lda[,1:9])
203
204
205 ## Some plots
206 partimat(
207   Class~Size+Instructor,
208   data = tae.tra[[1]],
209   method = "lda"
210 )
211 partimat(
212   Class~Size+Course,
213   data = tae.tra[[1]],
214   method = "lda"
215 )
216 partimat(
217   Class~Instructor+Course,
218   data = tae.tra[[1]],
219   method = "lda"
220 )
221
222
223
224 #####
225 # Quadratic Discriminant Analysis (QDA) #
226 #####
227 qda.fit = qda(
228   Class ~ .,
229   data = tae.tra[[1]]
230 )
231
232 qda.pred = predict(qda.fit, tae.tst.dt[[1]])
233
234 table(qda.pred$class, tae.tst.cl[[1]])
235
236 # Do it for every dataset in k-fold CV
237 acc = rep(0, 10)
238 for (i in 1:10) {
239   qda.fit = qda(

```

```

240     Class ~ .,
241     data = tae.tra[[i]]
242 )
243 qda.pred = predict(qda.fit, tae.tst.dt[[i]])
244
245 acc[i] = accuracy(
246     table(qda.pred$class, tae.tst.cl[[i]])
247 )
248 }
249 Acc.QDA = data.frame(fold = 1:10, Acc = acc)
250
251 mean(Acc.QDA$Acc)
252
253 # QDA performs poorly if every variable is considered.
254
255
256 # Fill in the table in the document
257 Acc.qda = matrix(ncol = 10, nrow = 10)
258 rownames(Acc.qda) = paste("K", 1:10, sep = "")
259 colnames(Acc.qda) =
260     c(paste("Acc", 1:9, sep = ""), "Acc_mean")
261 Kappa.qda = matrix(ncol = 10, nrow = 10)
262 rownames(Kappa.qda) = paste("K", 1:10, sep = "")
263 colnames(Kappa.qda) =
264     c(paste("Kappa", 1:9, sep = ""), "Kappa_mean")
265 for (m in 1:9) {
266     myAcc = rep(0, 10)
267     myKappa = rep(0, 10)
268     for (i in 1:10) {
269         qda.fit = qda(
270             Class ~ .,
271             data = tae.tra[[i]]
272         )
273         qda.pred = predict(qda.fit, tae.tst.dt[[i]])
274
275         cm = confusionMatrix(
276             table(qda.pred$class, tae.tst.cl[[i]])
277         )
278         myAcc[i] = cm$overall[1]
279         myKappa[i] = cm$overall[2]
280     }
281     Acc.qda[,m] = myAcc
282     Kappa.qda[,m] = myKappa
283 }
284 Acc.qda[,10] = rowMeans(Acc.qda[,1:9])
285 Kappa.qda[,10] = rowMeans(Kappa.qda[,1:9])
286
287
288 ## Some plots
289 partimat(
290     Class~Size+Instructor,
291     data = tae.tra[[1]],
292     method = "qda"
293 )
294 partimat(
295     Class~Size+Course,
296     data = tae.tra[[1]],
297     method = "qda"
298 )
299 partimat(
300     Class~Instructor+Course,

```

```

301 data = tae.tra[[1]],
302 method = "qda"
303 )
304
305
306
307 #####
308 # Algorithm comparison #
309 #####
310
311 # Acc and Kappa for every k-fold
312 Acc = cbind(Acc.knn[,10], Acc.lda[,10], Acc.qda[,10])
313 Acc = rbind(Acc, mean = colMeans(Acc))
314 colnames(Acc) = c("kNN", "LDA", "QDA")
315 Kappa = cbind(Kappa.knn[,10], Kappa.lda[,10], Kappa.qda[,10])
316 Kappa = rbind(Kappa, mean = colMeans(Kappa))
317 colnames(Kappa) = c("kNN", "LDA", "QDA")
318
319 # Check whether there is a method outperforming the others
320 friedman.test(as.matrix(Acc))
321
322 ## Low p-value clearly shows that there's a method
323 ## outperforming the others
324
325 groups = rep(1:ncol(Acc), each = nrow(Acc))
326 pairwise.wilcox.test(
327   Acc,
328   groups,
329   p.adjust.method = "holm",
330   paired = TRUE
331 )
332
333 # There are significant differences amongst the three methods,
334 # a plot will clarify it
335
336 myAcc = as.data.frame(cbind(Acc, k.fold = 1:11))[1:10,]
337 myData = melt.data.frame(
338   data = myAcc,
339   id.vars = c("k.fold"),
340   variable_name = "Method"
341 )
342 names(myData)[3] = "Acc"
343 ggplot(myData, aes(x = k.fold, y = Acc, fill = Method)) +
344   geom_col(position = "dodge")
345
346 myData = data.frame(Measure = c("Acc.mean", "Kappa.mean"))
347 myData = cbind(myData, rbind(t(Acc[11,]), t(Kappa[11,])))
348 myData = melt.data.frame(
349   data = as.data.frame(myData),
350   id.vars = "Measure",
351   variable_name = "Method"
352 )
353 ggplot(myData, aes(x = Method, y = value, fill = Method)) +
354   geom_col(position = "dodge") + ylab("") +
355   facet_wrap(~Measure, nrow = 1)

```

End Of File 'classification.R'

