

Tarea-Examen 3: Aprendizaje Supervisado

Curso Avanzado de Estadística. Profa. Guillermina Eslava Gómez.

Aldo Sayeg Pasos Trejo. César Cossio Guerrero.

Posgrado en Ciencias Matemáticas. Universidad Nacional Autónoma de México.

16 de abril de 2020

1. Problema 1

Buscamos realizar un estudio de clasificación sobre la base de datos **Pima**, que consiste en 532 observaciones de 7 variables usadas para describir el estado de salud de las mujeres de la tribu Pima, cuya tasa de enfermedad de Diabetes Mellitus es inusualmente alta. La descripción [1] de las variables se muestra en la tabla 13.

Variable	Descripción
“npreg”	número de embarazos
“glu”	Concentración de glucosa en el plasma para una prueba oral
“bp”	Presión arterial diastólica (mmHg)
“skin”	grueso de la piel del triceps (mm)
“bmi”	Índice de masa corporal (kg/m^2)
“ped”	Función pedigree para diabetes
“age”	Edad (años)
“type”	Indicador de diabetes (“yes” o “no”)

Tabla 1: Descripción de las variables

Antes de realizar cualquier estudio sobre el sistema, analizamos la relación entre sus variables analizando su correlación en la figura 1, su separación por clase en una gráfica a pares como en la figura 2, por componentes principales como en 3 y sus distribuciones individuales en las figuras 4,5,6,7,8,9 y 10.

De ese análisis, podemos concluir rápidamente que “bmi” presenta una alta correlación con “skin”. Esto se puede entender claramente ya que el índice de masa corporal debe de tener un efecto claro

sobre la piel en el tríceps. De igual manera, la correlación entre “age” y “npreg” es totalmente obvia debido a que el número de embarazos se relaciona directamente con la edad.

En cuanto a las diferencias entre clases, nuevamente no se pueden observar muchas diferencias sustanciales y las que se muestran como la de la figura 10 se entiende nuevamente de manera intuitiva ya que la incidencia de Diabetes debe de ser poca en mujeres jóvenes. La diferencia entre la distribución de “glu” en la figura 5 también se puede explicar de manera trivial al igual que la de la figura 8 al ser la obesidad un factor de riesgo para la diabetes.

Esta análisis nos permite especular que las variables con mayor significancia estadística pueden ser “bmi”, “glu”, “age” y que podemos usar una sola variable para representar variables altamente correlacionadas. En principio, nos interesa usar 4 modelos para clasificar

- a) Análisis de discriminante lineal
- b) Naive Bayes
- c) Regresión logística
- d) Support Vector Machines (SVM)

Yendo un poco más lejos, consideramos realizar un análisis incluyendo también el Análisis de discriminante cuadrático y K-Neighbors. La figura 11 muestra las tasas de clasificación global y local para esos modelos para distintos tamaños del conjunto de entrenamiento de los datos. La separación de los datos se realizó sin tomar en cuenta el mantener la proporcionalidad entre las clases. Notemos que la tasa aparente corresponde a la fracción $s = 1.0$. Por limitaciones del poder de cómputo disponible, se realizaron $n = 10$ iteraciones para las tasas no aparentes.

Podemos observar que, en cuanto a la tasa de error global aparente, el modelo K-Neighbors presenta el menor valor.

Sin embargo, pareciera intuitivo pensar que las tasas no aparentes deben de mejorar si ahora forzamos a que la separación entre conjunto de entrenamiento y de prueba mantenga la proporcionalidad existen en las clases. La tabla 26 y la gráfica 12 muestran dichos resultados.

En efecto, se muestra que las tasas no aparentes mejoran aunque no significativamente. Si pensamos en la capacidad predictiva del modelo y en el contexto del diagnóstico médico, es claro que la tasa de aceptación que más nos interesa son los falsos negativos pues nos interesa diagnosticar con mayor precisión. En ese sentido, los modelos entrenados manteniendo la proporcionalidad presentan tasas muy elevadas de error para fracciones pequeñas.

Antes de proceder modelo por modelo, buscamos para cada variable numérica y positiva del modelo, cual es el valor λ que maximiza la verosimilitud de una transformación Box-Cox para esa variable. Dichos valores, junto con sus intervalos de confianza, se pueden consultar en la tabla 27

Habiendo realizado primero este análisis, ahora nos disponemos a presentar los 4 modelos relevantes.

1.1. Análisis de discriminante lineal

Para discriminante lineal, no hay muchos parámetros que podamos ajustar para posibles modelos. Nos enfocamos más bien en buscar transformaciones o interacciones entre las variables.

Después de mostrar que las transformaciones Box-Cox no presentan una mejoría sustancial para ninguno de los modelos, y que esta es considerada la mejor clase de transformaciones dado que las figuras 1, 2 no nos muestra que otra transformación podría mejorar la separación con una variable, nos limitamos a no realizar transformaciones sobre las variables.

Así, solo añadimos posibles interacciones. Entre interacciones a pares, encontramos que la interacción que presentaba menor tasa global aparente de error era la interacción entre “ped” y “age”.

Añadiendo esta interacción al modelo, después consideramos interacciones terciarias que lo hagan óptimo. Sin embargo, todas las tasas de interacciones muestran tasas similares (mayores 1 0.20), por lo que no consideramos que representen una mejoría sobre nuestro modelo y, así, nos quedamos con este modelo en interacción a pares como el mejor clasificador.

1.2. Naive Bayes

Nuevamente, no existen muchos parámetros modificables en el método de Naive Bayes y, por los argumentos anteriormente mostrados, recurrimos a otra interacción binaria, la interacción entre “ped” e “bp”.

1.3. Regresión Logística

Existen mejores estimadores para determinar los valores de la regresión logística. La tabla de los p -values para una regresión con todas las variables la podemos presentar en la tabla 28.

Después de revisar la significancia estadística de la regresión, nos podemos quedar solamente con las variables “glu”, “bmi”, “ped” y “age”. Además, considerando que la edad presenta una gran influencia, podemos incluir un término cuadrático para la edad “age²”. Así, este modelo es el que consideramos óptimo y representativo de este método.

1.4. Support Vector Machines

Por ser el método que presenta más parámetros para el clasificador, primero realizamos una corrida para kernels lineales, de funciones de base radial de grado 1,3 y 5 y polinomiales de los mismos grados. Los resultados se pueden consultar en la figura 13. Aunque existen también los Kernels sigmoides [2], al probar con estos encontramos tasas de clasificación muy bajas, por lo que los excluimos del análisis.

El kernel lineal se desempeña muy bien también al ser comparado por los otros. Tanto clasificadores de base radial como los polinomiales muestran muy buen desempeño aun para grados bajos. Sin embargo, el clasificador que consideramos como el óptimo es el que minimiza la tasa global aparente, es decir, **el de kernel polinomial y grado 5**.

Por las limitaciones del poder computacional, resultó imposible analizar distintos valores del parámetro γ . Aunque ahora podríamos analizar interacciones, por la estructura de los parámetros del sis-

Método	Variables en el modelo	Notas extra
Análisis de Discriminante Lineal	“npreg”, “glu”, “bp”, “skin”, “bmi”, “ped”, “age”, “type”, “ped*age”	-
Naive Bayes	“npreg”, “glu”, “bp”, “skin”, “bmi”, “ped”, “age”, “type”, “ped*bp”	-
Regresión logística	“glu”, “bmi”, “ped”, “age”, “age ² ”	-
Support Vector Machines	“npreg”, “glu”, “bp”, “skin”, “bmi”, “ped”, “age”, “type”	Kernel Polynomial de grado 5

Tabla 2: Variables de los modelos

tema, sabemos que estas ya se están contemplando al aumentar el grado del kernel. Así, dejamos este modelo como el definitivo para este clasificador y presentamos sus tasas de error de manera gráfica individual como tabular en la figura 14.

1.5. Conclusiones

Para resumir, en la tabla 2 mostramos las variables utilizadas por el modelo correspondiente a cada método: Comparamos los cuatro modelos seleccionados para cada método. La tabla 30 y la figura 15 muestran dicha comparación. De Inmediatamente podemos concluir que ningún método alcanza una tasa de error global menor al 20 % . Todos los modelos logran mantener a la tasa de error global, tanto aparente como no aparente, en valores menores a 25 %. Sintetizamos dichos resultados en las tablas 4 y 6.

	LDA	Naive Bayes	Logistic	SVM, kernel = poly, deg = 5
Mean global error, entrenamiento: 1.0	0.205	0.205	0.207	0.205
Mean Yes error, entrenamiento: 1.0	0.412	0.412	0.390	0.475
Mean No error, entrenamiento: 1.0	0.101	0.101	0.115	0.070

Tabla 4: Tasas de error aparente locales y globales para los cuatro métodos con conjuntos divididos manteniendo proporcionalidad

	LDA	Naive Bayes	Logistic	SVM, kernel = poly, deg = 5
Mean global error, entrenamiento: 0.7	0.216	0.210	0.207	0.223
Mean Yes error, entrenamiento: 0.7	0.412	0.413	0.381	0.496
Mean No error, entrenamiento: 0.7	0.106	0.104	0.117	0.071

Tabla 6: Tasas de error aparente locales y globales para los cuatro métodos con conjuntos divididos manteniendo proporcionalidad

En general, el modelo que menos falsos negativos da es la SVM. Sin embargo, su tasa de falsos positivos siempre es muy alta. En general, consideramos que el mejor método resulta ser el modelo de regresión logística pues sus tasas de error locales son bastante estables independientemente del tamaño de la muestra y, aunque la global no aparente no tiene tanta estabilidad, la global aparente presente un valor bajo y no distinto a los otros valores.

2. Problema 2

Tenemos una base de datos, consistente en 236 observaciones que representa a pacientes de una clínica danesa de cardiología [3], con 14 variables categóricas cuya descripción muestra la tabla 7

Variable	Descripción
“Sex”	Sexo (“male” or “female”)
“AngPec”	Indicador de si padece angina de pecho (“Atypical”, “None” or “Typical”)
“AMI”	Indicador de problemas cardiacos (“yes” or “no”)
“QWave”	Indicador de presencia de una onda en el electrocardiograma (ECG) (“yes” or “no”)
“QWavecode”	Indicador de si es representativo la presencia de la onda Q (“Usable” or “Nonusable”)
“STcode”	Indicador de si es representativa la presencia de una onda T en el ECG (“Usable” or “Nonusable”)
“STchange”	Indicador de la presencia de una onda T en el ECG (“yes” or “no”)
“SuffHeartF”	Indicador de si ha presentado infartos cardiacos (“yes” o “no”)
“Hypertrophi”	Indicador de si ha habido hipertrofia arterial (“yes” o “no”)
“Hyperchol”	Indicador de si hay colesterol alto (“yes” o “no”)
“Smoker”	Indicador de si es fumador (“yes” o “no”)
“Inherit”	Indicador de si hay predisposiciones hereditarias (“yes” o “no”)
“Heartfail”	Indicador de si ha tenido paros cardiacos (“yes” o “no”)
“CAD”	Indicador de presencia de la enfermedad de las arterias coronarias (“yes” o “no”)

Tabla 7: Descripción de las variables

Nos interesa clasificar los valores de “CAD” correctamente. Para visualizar las variables no podemos graficar ni realizar un cambio a componentes principales de estas variables al ser categóricas, podemos ver un histograma de sus distintos valores entre las clases que las componen, como muestran las figuras 16,18,17,19,20,21,22,23,24,25,26,27 y28.

De la valoración cualitativa, vemos que las variables que muestran una mayor diferencia en sus frecuencias para cada categoría entre sus dos grupos son “AngPec”, “AMI”, “QWave” y “Hyperchol”. Es intuitivo pensar que el alto colesterol, la angina de pecho y la presencia de malestares cardiacos en general puede contribuir a la enfermedad CAD, por lo que las diferencias entre los grupos nos parecen razonables. Sin embargo, sin mayor conocimiento sobre las variables, no podemos realizar más a severaciones.

Procedemos a construir modelos de clasificación sobre con los siguientes métodos

- a) Naive Bayes
- b) Regresión logística
- c) Support Vector machines

Nuevamente, de manera previa, realizamos un barrido con muchos métodos de clasificación para obtener un panorama general de su utilidad. La tabla 32 y la figura 29 muestran los resultados de dichas cosas. La tasa global de error aparente es menor al 20 % para todos, lo que indica alta precisión en todos los métodos. También observamos que las tasas locales para cada clase tienen valores cercanos a la local, lo que indica la estabilidad entre clases de los métodos. Procedemos a analizar método tras método.

2.1. Naive Bayes

Nuevamente, no nos es posible refinar este método de muchas maneras, por lo que buscamos la interacción binaria que minimiza el error global aparente. Encontramos que es de la variable “Sex” con “AMI”. Mantenemos nuestro modelo con ese clasificador.

2.2. Regresión logística

Realizando de manera inmediata una regresión logística sobre los datos, podemos obtener los p -values para estimar la significancia estadística de cada variable. La tabla 33 muestra dichos resultados. Como podemos observar, las únicas que presentan un p -value que podríamos considerar como significativo son las variables “AngPec”, “AMI”, “STcode”, “STchange” y “Hyperchol”, por lo que son las únicas que consideraremos para el clasificador.

2.3. Support Vector machines

Realizamos el mismo análisis que en el ejercicio anterior, primero realizamos una corrida para kernels lineales, de funciones de base radial de grado 1,3 y 5 y polinomiales de los mismos grados. Los resultados se pueden consultar en la figura . Notamos que el clasificador polinomial de grado 5 es el que presenta la menor tasa de error global aparente, por lo que la consideramos la más adecuada para nuestros datos.

Las limitaciones computacionales nos impiden explorar valores de la variable γ distintos.

Método	Variables en el modelo	Notas extra
Naive Bayes	“Sex”, “AngPec”, “AMI”, “QWave”, “QWavecode”, “STcode”, “STchange”, “SuffHeartF”, “Hypertrop-hi”, “Hyperchol”, “Smoker”, “Inherit”, “Heartfail”, “CAD”, “Sex*AMI”	-
Regresión logística	“AngPec”, “AMI”, “STcode”, “STchange”, “Hyperchol”	-
Support Vector Machines	“Sex”, “AngPec”, “AMI”, “QWave”, “QWavecode”, “STcode”, “STchange”, “SuffHeartF”, “Hypertrop-hi”, “Hyperchol”, “Smoker”, “Inherit”, “Heartfail”, “CAD”	Kernel Polynomial de grado 5

Tabla 8: Variables de los modelos

2.4. Conclusiones

En la tabla 8 mostramos las variables utilizadas por el modelo correspondiente a cada método: La gráfica 31 y la tabla 35 muestran los resultados para los tres modelos presentados. Sin ninguna duda, el modelo de la SVM presenta la menor tasa global de error aparente, así como estabilidad directa sobre la diferencia entre las tasas de errores locales no aparentes. En general, sus tasas locales y globales no aparentes también son las menores entre todos los valores presentes. Presentamos un resumen de dichas tablas en las tablas 10 y 12

	Naive Bayes	Logistic	SVM, kernel = poly, deg = 5
Mean global error, entrenamiento: 1.0	0.144	0.169	0.038
Mean Yes error, entrenamiento: 1.0	0.150	0.243	0.037
Mean No error, entrenamiento: 1.0	0.140	0.109	0.039

Tabla 10: Tasas de error aparente locales y globales para los cuatro métodos con conjuntos divididos manteniendo proporcionalidad

	Naive Bayes	Logistic	SVM, kernel = poly, deg = 5
Mean global error, en- trenamiento: 0.7	0.170	0.201	0.200
Mean Yes error, entre- namiento: 0.7	0.171	0.196	0.104
Mean No error, entrena- miento: 0.7	0.154	0.153	0.060

Tabla 12: Tasas de error aparente locales y globales para los cuatro métodos con conjuntos divididos manteniendo proporcionalidad

Cabe destacar que dicho modelo, el de SVM, no utiliza ninguna interacción ni ninguna transformación de variables, por lo que lo consideramos óptimo en el sentido de que maneja los datos de manera directa

3. Problema 3

Se nos presenta una base de datos de 145 observaciones de 5 variables numéricas y una variable categórica, las cuales describimos a continuación

Variable	Descripción
“Weight”	Peso Kg()
“Fglucose”	Concentración de glucosa en el plasma en ayunoas
“Glucoseint”	Area bajo una curva de glucosa para una medición constante en 3 horas
“InsulineResist”	Area bajo una curva de insulina para una medición constante en 3 horas
“InsulinResp”	Respuesta de insulina a la glucosa
“Class”	Clase de paciente (“Overt”, “Chemical” y “Normal”)

Tabla 13: Descripción de las variables

Nuevamente, primero analizamos la relación entre sus variables analizando su correlación en la figura 32, su separación por clase en una gráfica a pares como en la figura 33 y sobre las componentes principales en la figura 34.

El análisis exploratorio nos muestra que existe una alta correlación entre “Glucoseint” y “Fglucose”, lo cual se puede explicar de manera inmediata. Si realizamos un ajuste con un clasificador multinomial tomando como referencia la clase 3 y un clasificador que compara cada clase con las demás (one vs others), se obtienen las tasas mostradas en la figura 39 y en la tabla 37. Usaremos el clasificador multinomial para el análisis.

Para seleccionar modelos predictivos y descriptivos, primero se revisó el artículo en cuestión [3-3] que trabajó con los datos. El modelo de referencia es el que el artículo propone y consta de las siguientes tres variables “GlucoseInt”, “InsulinResp” y “InsulineResist”. Cabe mencionar que los nombres utilizados en el artículo original para denotar a las variables son distintos de los presentados en la base de datos.

Por otra parte, seleccionamos la categoría de “Chemical” como referencia, esto debido a que en el artículo parecen estar interesados en diferenciar esta con la categoría Over. Esta elección, desde luego, obedece más a qué pregunta queramos contestar. Sin embargo, en nuestro caso lo utilizamos porque parece que nos ayudaría a identificar con mayor facilidad los grupos.

3.1. Modelo descriptivo

Se realizaron diversas funciones para encontrar la combinación óptima dentro de un conjunto de modelos que consideraba términos simples, interacciones y términos cuadráticos. Para lograr esto se implementó una rutina que utiliza una búsqueda tabú con una función de optimización el AIC, BIC y Deviance de cada modelo. En dichas rutinas se consideró la que todas las variables involucradas fueran significativas respecto de la estadística F . El resultado de dicho análisis arrojó los modelos presentados en la tabla 14. En la tabla se muestran modelos solo de tres términos, y esto se debe

Variables en el modelo	AIC	BIC	Deviance
“Fglucose” , “Fglucose ² ” , “GlucoseInt ² ”	34.9	58.7	18.9
“Fglucose” , “GlucoseInt” , “Fglucose ² ”	30.7	54.5	14.7
“Weight” , “Weight*InsulinResp” , “InsulinResp ² ”	230.2	254.0	214.2

Tabla 14: Resumen de los posibles modelos descriptivos con los parámetros de su ajuste

principalmente a dos motivos: por una parte, como indica el análisis del artículo, el número debería ser próximo a 3 términos; y por el otro lado, casi todos los modelos que optimizaron cada medida tenían ese número de términos (como se puede apreciar en la tabla 14). El mejor modelo, en términos descriptivos es el segundo, que tomaremos como el modelo descriptivo.

3.2. Modelo predictivo

Para calcular un modelo predictivo, para lo cual se utilizó la búsqueda tabú pero ahora la función objetivo que se buscó minimizar fue el mínimo del máximo de las tasas de error por grupos en errores no aparentes para una fracción de entrenamiento de 0.7. La tabla 15 resume tres posibles modelos, mientras que la tabla 17 muestra sus tasas de error no aparentes.

Nombre del modelo	Variables en el modelo
Modelo 1	“Fglucose” , “GlucoseInt*GlucoseInt”
Modelo 2	“InsulinResp”, “Fglucose*InsulinResp” , “GlucoseInt*InsulinResp”
Modelo 3	“Fglucose” , “Weight*InsulineResist” , “Fglucose*GlucoseInt” , “GlucoseInt*GlucoseInt”

Tabla 15: Resumen de los posibles modelos predictivos con sus respectivas tasas de error no aparente utilizando K fold cross-validation

	Modelo 1	Modelo 2	Modelo 3
Mean global error, entrenamiento: 0.7	0.140	0.030	0.230
Mean 1 error, entrenamiento: 0.7	0.121	0.009	0.002
Mean 2 error, entrenamiento: 0.7	0.474	0.062	0.894
Mean 3 error, entrenamiento: 0.7	0.005	0.010	0.009

Tabla 17: Tasas de error no aparente locales y globales para los tres modelos con conjuntos divididos manteniendo proporcionalidad

El modelo que minimiza dichos valores, muy por encima de todos los otros contrincantes, es el modelo 2, por lo que lo seleccionamos como el modelo predictivo.

3.3. Conclusiones

La tabla 18 muestra las variables de los dos modelos trabajados

Modelo	Variables	Notas
Descriptivo	“Fglucose” , “GlucoseInt” , “Fglucose ² ”	-
Predictivo	“InsulinResp”, “Fglucose*InsulinResp” , “GlucoseInt*InsulinResp”	-

Tabla 18: Modelos descriptivo y predictivo

A continuación, mostramos las tasas de error aparente y no aparente para el modelo predictivo y descriptivo en la tablas

	Descriptivo	Predictivo
Mean global error, entrenamiento: 1.0	0.021	0.021
Mean 1 error, entrenamiento: 1.0	0.000	0.000
Mean 2 error, entrenamiento: 1.0	0.056	0.056
Mean 3 error, entrenamiento: 1.0	0.013	0.013

Tabla 20: Tasas de error aparente locales y globales para el modelo descriptivo y predictivo

	Descriptivo	Predictivo
Mean global error, entrenamiento: 0.7	0.023	0.023
Mean 1 error, entrenamiento: 0.7	0.000	0.002
Mean 2 error, entrenamiento: 0.7	0.056	0.057
Mean 3 error, entrenamiento: 0.7	0.022	0.011

Tabla 22: Tasas de error no aparente locales y globales para el modelo descriptivo y predictivo con conjuntos divididos manteniendo proporcionalidad

Es claro que el modelo descriptivo y el modelo predictivo difieren. También, en el último no aparecen los términos individuales de las interacciones; no obstante, al agregar dichos términos individuales al modelo no eran significativos según la estadística F . Ambos modelos tienen un alto poder predictivo, mientras que la interpretación del descriptivo es mucho más sencilla.

Anexo 1: Tablas relevantes**Problema 1**

	LDA	QDA	Naive Bayes	Logistic	KNC, k=5	SVM
Mean global error, entrenamiento: 0.1	0.241	0.278	0.250	0.241	0.258	0.243
std, entrenamiento: 0.1	0.017	0.027	0.020	0.017	0.027	0.036
Mean Yes error, entrenamiento: 0.1	0.406	0.475	0.352	0.426	0.515	0.537
std, entrenamiento: 0.1	0.089	0.089	0.088	0.098	0.127	0.180
Mean No error, entrenamiento: 0.1	0.144	0.157	0.190	0.139	0.118	0.094
std, entrenamiento: 0.1	0.042	0.054	0.057	0.038	0.046	0.049
Mean global error, entrenamiento: 0.2	0.229	0.266	0.247	0.234	0.250	0.240
std, entrenamiento: 0.2	0.014	0.017	0.021	0.014	0.015	0.034
Mean Yes error, entrenamiento: 0.2	0.390	0.423	0.342	0.405	0.429	0.584
std, entrenamiento: 0.2	0.044	0.099	0.038	0.049	0.048	0.151
Mean No error, entrenamiento: 0.2	0.136	0.158	0.187	0.135	0.143	0.070
std, entrenamiento: 0.2	0.021	0.044	0.038	0.024	0.028	0.038
Mean global error, entrenamiento: 0.3	0.219	0.251	0.234	0.228	0.251	0.232
std, entrenamiento: 0.3	0.014	0.014	0.014	0.013	0.014	0.019
Mean Yes error, entrenamiento: 0.3	0.429	0.394	0.357	0.412	0.408	0.523
std, entrenamiento: 0.3	0.041	0.034	0.034	0.049	0.040	0.074
Mean No error, entrenamiento: 0.3	0.110	0.145	0.164	0.123	0.143	0.077
std, entrenamiento: 0.3	0.019	0.014	0.025	0.019	0.022	0.033
Mean global error, entrenamiento: 0.4	0.233	0.242	0.241	0.227	0.265	0.224
std, entrenamiento: 0.4	0.020	0.014	0.013	0.018	0.014	0.019
Mean Yes error, entrenamiento: 0.4	0.412	0.394	0.347	0.406	0.459	0.521
std, entrenamiento: 0.4	0.032	0.047	0.026	0.022	0.054	0.046
Mean No error, entrenamiento: 0.4	0.125	0.140	0.175	0.122	0.119	0.075
std, entrenamiento: 0.4	0.015	0.024	0.017	0.014	0.019	0.018
Mean global error, entrenamiento: 0.5	0.215	0.240	0.247	0.227	0.250	0.220
std, entrenamiento: 0.5	0.016	0.020	0.018	0.015	0.018	0.017

Continued on next page

	LDA	QDA	Naive Bayes	Logistic	KNC, k=5	SVM
Mean Yes error, entrenamiento: 0.5	0.407	0.406	0.351	0.420	0.396	0.493
std, entrenamiento: 0.5	0.028	0.029	0.025	0.024	0.034	0.017
Mean No error, entrenamiento: 0.5	0.117	0.133	0.174	0.112	0.128	0.084
std, entrenamiento: 0.5	0.013	0.015	0.015	0.013	0.016	0.010
Mean global error, entrenamiento: 0.6	0.217	0.248	0.234	0.228	0.263	0.236
std, entrenamiento: 0.6	0.020	0.031	0.016	0.019	0.020	0.036
Mean Yes error, entrenamiento: 0.6	0.425	0.408	0.346	0.414	0.398	0.482
std, entrenamiento: 0.6	0.021	0.016	0.020	0.023	0.037	0.033
Mean No error, entrenamiento: 0.6	0.108	0.134	0.174	0.119	0.123	0.093
std, entrenamiento: 0.6	0.012	0.008	0.015	0.013	0.017	0.021
Mean global error, entrenamiento: 0.7	0.204	0.252	0.228	0.230	0.250	0.221
std, entrenamiento: 0.7	0.021	0.029	0.030	0.034	0.032	0.028
Mean Yes error, entrenamiento: 0.7	0.414	0.384	0.339	0.422	0.337	0.477
std, entrenamiento: 0.7	0.016	0.019	0.009	0.030	0.032	0.023
Mean No error, entrenamiento: 0.7	0.116	0.139	0.181	0.111	0.140	0.094
std, entrenamiento: 0.7	0.011	0.011	0.007	0.011	0.008	0.015
Mean global error, entrenamiento: 0.8	0.213	0.246	0.238	0.254	0.266	0.221
std, entrenamiento: 0.8	0.036	0.033	0.053	0.035	0.044	0.038
Mean Yes error, entrenamiento: 0.8	0.423	0.390	0.337	0.421	0.329	0.496
std, entrenamiento: 0.8	0.018	0.012	0.016	0.021	0.022	0.021
Mean No error, entrenamiento: 0.8	0.109	0.139	0.180	0.109	0.121	0.082
std, entrenamiento: 0.8	0.009	0.005	0.006	0.007	0.012	0.014
Mean global error, entrenamiento: 0.9	0.209	0.263	0.259	0.217	0.252	0.231
std, entrenamiento: 0.9	0.046	0.062	0.048	0.031	0.054	0.048
Mean Yes error, entrenamiento: 0.9	0.417	0.397	0.348	0.426	0.317	0.482
std, entrenamiento: 0.9	0.011	0.009	0.006	0.013	0.011	0.009
Mean No error, entrenamiento: 0.9	0.114	0.141	0.179	0.111	0.114	0.090
std, entrenamiento: 0.9	0.006	0.005	0.004	0.006	0.015	0.009

Continued on next page

	LDA	QDA	Naive Bayes	Logistic	KNC, k=5	SVM
Mean global error, entrenamiento: 1.0	0.212	0.224	0.239	0.212	0.171	0.216
std, entrenamiento: 1.0	0.000	0.000	0.000	0.000	0.000	0.000
Mean Yes error, entrenamiento: 1.0	0.424	0.390	0.345	0.424	0.294	0.480
std, entrenamiento: 1.0	0.000	0.000	0.000	0.000	0.000	0.000
Mean No error, entrenamiento: 1.0	0.107	0.141	0.186	0.107	0.110	0.085
std, entrenamiento: 1.0	0.000	0.000	0.000	0.000	0.000	0.000

Tabla 24: Tasas de error locales y globales para varios clasificadores con conjuntos de entrenamiento divididos sin mantener proporcionalidad

	LDA	QDA	Naive Bayes	Logistic	KNC, k=5	SVM
Mean global error, entrenamiento: 0.1	0.247	0.287	0.252	0.235	0.258	0.247
std, entrenamiento: 0.1	0.025	0.024	0.014	0.015	0.020	0.032
Mean Yes error, entrenamiento: 0.1	0.455	0.461	0.377	0.429	0.449	0.600
std, entrenamiento: 0.1	0.055	0.104	0.046	0.082	0.058	0.164
Mean No error, entrenamiento: 0.1	0.130	0.175	0.183	0.130	0.149	0.067
std, entrenamiento: 0.1	0.033	0.051	0.032	0.041	0.029	0.037
Mean global error, entrenamiento: 0.2	0.229	0.266	0.251	0.233	0.252	0.230
std, entrenamiento: 0.2	0.011	0.021	0.017	0.013	0.018	0.007
Mean Yes error, entrenamiento: 0.2	0.428	0.437	0.363	0.426	0.425	0.521
std, entrenamiento: 0.2	0.028	0.073	0.045	0.047	0.044	0.065
Mean No error, entrenamiento: 0.2	0.120	0.147	0.175	0.126	0.148	0.081
std, entrenamiento: 0.2	0.015	0.031	0.032	0.025	0.030	0.026
Mean global error, entrenamiento: 0.3	0.227	0.242	0.241	0.231	0.265	0.229
std, entrenamiento: 0.3	0.020	0.010	0.013	0.012	0.014	0.011
Mean Yes error, entrenamiento: 0.3	0.382	0.407	0.350	0.410	0.420	0.528
std, entrenamiento: 0.3	0.025	0.042	0.026	0.038	0.033	0.036
Mean No error, entrenamiento: 0.3	0.135	0.138	0.175	0.123	0.154	0.073

Continued on next page

	LDA	QDA	Naive Bayes	Logistic	KNC, k=5	SVM
std, entrenamiento: 0.3	0.020	0.014	0.019	0.015	0.022	0.013
Mean global error, entrenamiento: 0.4	0.226	0.251	0.237	0.221	0.258	0.227
std, entrenamiento: 0.4	0.016	0.011	0.013	0.014	0.019	0.012
Mean Yes error, entrenamiento: 0.4	0.413	0.408	0.346	0.408	0.402	0.511
std, entrenamiento: 0.4	0.030	0.035	0.024	0.028	0.037	0.052
Mean No error, entrenamiento: 0.4	0.115	0.140	0.175	0.120	0.145	0.078
std, entrenamiento: 0.4	0.011	0.014	0.022	0.021	0.021	0.022
Mean global error, entrenamiento: 0.5	0.212	0.238	0.241	0.225	0.253	0.223
std, entrenamiento: 0.5	0.013	0.017	0.020	0.012	0.015	0.012
Mean Yes error, entrenamiento: 0.5	0.410	0.402	0.338	0.405	0.383	0.510
std, entrenamiento: 0.5	0.022	0.025	0.010	0.016	0.036	0.026
Mean No error, entrenamiento: 0.5	0.114	0.136	0.171	0.121	0.134	0.078
std, entrenamiento: 0.5	0.013	0.012	0.011	0.009	0.018	0.015
Mean global error, entrenamiento: 0.6	0.225	0.246	0.230	0.220	0.251	0.225
std, entrenamiento: 0.6	0.025	0.023	0.022	0.024	0.024	0.024
Mean Yes error, entrenamiento: 0.6	0.412	0.382	0.354	0.412	0.379	0.481
std, entrenamiento: 0.6	0.021	0.024	0.019	0.017	0.026	0.033
Mean No error, entrenamiento: 0.6	0.116	0.142	0.167	0.118	0.122	0.093
std, entrenamiento: 0.6	0.009	0.010	0.012	0.010	0.011	0.019
Mean global error, entrenamiento: 0.7	0.222	0.244	0.236	0.237	0.238	0.225
std, entrenamiento: 0.7	0.034	0.028	0.023	0.020	0.020	0.026
Mean Yes error, entrenamiento: 0.7	0.420	0.385	0.341	0.406	0.351	0.477
std, entrenamiento: 0.7	0.016	0.014	0.015	0.015	0.022	0.018
Mean No error, entrenamiento: 0.7	0.113	0.143	0.182	0.120	0.121	0.095
std, entrenamiento: 0.7	0.009	0.006	0.009	0.009	0.017	0.013
Mean global error, entrenamiento: 0.8	0.232	0.234	0.245	0.219	0.250	0.233
std, entrenamiento: 0.8	0.031	0.034	0.034	0.022	0.032	0.031
Mean Yes error, entrenamiento: 0.8	0.415	0.393	0.339	0.431	0.356	0.485
std, entrenamiento: 0.8	0.008	0.022	0.007	0.010	0.033	0.005

Continued on next page

	LDA	QDA	Naive Bayes	Logistic	KNC, k=5	SVM
Mean No error, entrenamiento: 0.8	0.116	0.136	0.179	0.110	0.116	0.089
std, entrenamiento: 0.8	0.008	0.005	0.009	0.008	0.014	0.006
Mean global error, entrenamiento: 0.9	0.215	0.239	0.228	0.237	0.259	0.198
std, entrenamiento: 0.9	0.051	0.046	0.054	0.053	0.057	0.039
Mean Yes error, entrenamiento: 0.9	0.419	0.390	0.339	0.422	0.315	0.488
std, entrenamiento: 0.9	0.008	0.009	0.007	0.011	0.013	0.008
Mean No error, entrenamiento: 0.9	0.110	0.139	0.181	0.114	0.110	0.083
std, entrenamiento: 0.9	0.006	0.003	0.004	0.006	0.006	0.005
Mean global error, entrenamiento: 1.0	0.212	0.224	0.239	0.212	0.171	0.216
std, entrenamiento: 1.0	0.000	0.000	0.000	0.000	0.000	0.000
Mean Yes error, entrenamiento: 1.0	0.424	0.390	0.345	0.424	0.294	0.480
std, entrenamiento: 1.0	0.000	0.000	0.000	0.000	0.000	0.000
Mean No error, entrenamiento: 1.0	0.107	0.141	0.186	0.107	0.110	0.085
std, entrenamiento: 1.0	0.000	0.000	0.000	0.000	0.000	0.000

Tabla 26: Tasas de error locales y globales para varios clasificadores con conjuntos de entrenamiento divididos sin mantener proporcionalidad

	lambda	Lower confidence interval, alpha = 0.05	Upper confidence interval, alpha = 0.05
glu	-0.1487	-0.4522	0.1580
bp	1.0763	0.7520	1.4125
skin	0.5490	0.3709	0.7263
bmi	0.1249	-0.2044	0.4510
ped	-0.0090	-0.1269	0.1084
age	-1.4435	-1.7747	-1.1202

Tabla 27: Valores de λ que maximizan la verosimilitud de una transformación Box-Cox

Coefficient	Estimate	Std. Error	z value	Pr($p > z $)
(Intercept)	-0.97191	0.12103	-8.030	9.72e-16
glu	1.06690	0.13040	8.182	2.80e-16
bp	-0.10413	0.12553	-0.830	0.406802
skin	0.09456	0.15746	0.601	0.548172
bmi	0.53147	0.15900	3.342	0.000830
ped	0.42914	0.12360	3.472	0.000516
age	0.54557	0.12267	4.448	8.68e-06

Tabla 28: Coeficientes y p -values para la regresión logística del modelo completo

	LDA	Naive Bayes	Logistic	SVM, kernel = poly, deg = 5
Mean global error, entrenamiento: 0.1	0.231	0.249	0.219	0.253
std, entrenamiento: 0.1	0.018	0.031	0.015	0.023
Mean Yes error, entrenamiento: 0.1	0.389	0.434	0.392	0.529
std, entrenamiento: 0.1	0.071	0.047	0.075	0.123
Mean No error, entrenamiento: 0.1	0.143	0.145	0.128	0.098
std, entrenamiento: 0.1	0.029	0.052	0.035	0.039
Mean global error, entrenamiento: 0.2	0.224	0.236	0.216	0.245
std, entrenamiento: 0.2	0.010	0.015	0.011	0.012
Mean Yes error, entrenamiento: 0.2	0.392	0.423	0.392	0.533
std, entrenamiento: 0.2	0.028	0.042	0.044	0.083
Mean No error, entrenamiento: 0.2	0.128	0.128	0.121	0.085
std, entrenamiento: 0.2	0.027	0.021	0.021	0.029
Mean global error, entrenamiento: 0.3	0.213	0.232	0.207	0.237
std, entrenamiento: 0.3	0.008	0.010	0.010	0.018
Mean Yes error, entrenamiento: 0.3	0.410	0.407	0.395	0.519
std, entrenamiento: 0.3	0.028	0.023	0.027	0.042
Mean No error, entrenamiento: 0.3	0.110	0.121	0.111	0.077
std, entrenamiento: 0.3	0.014	0.018	0.019	0.018
Mean global error, entrenamiento: 0.4	0.216	0.219	0.211	0.243
std, entrenamiento: 0.4	0.009	0.012	0.021	0.014

Continued on next page

	LDA	Naive Bayes	Logistic	SVM, kernel = poly, deg = 5
Mean Yes error, entrenamiento: 0.4	0.404	0.401	0.387	0.472
std, entrenamiento: 0.4	0.025	0.018	0.021	0.031
Mean No error, entrenamiento: 0.4	0.115	0.119	0.113	0.087
std, entrenamiento: 0.4	0.016	0.011	0.012	0.013
Mean global error, entrenamiento: 0.5	0.215	0.202	0.213	0.229
std, entrenamiento: 0.5	0.020	0.020	0.014	0.010
Mean Yes error, entrenamiento: 0.5	0.413	0.406	0.393	0.490
std, entrenamiento: 0.5	0.023	0.019	0.022	0.024
Mean No error, entrenamiento: 0.5	0.103	0.110	0.113	0.075
std, entrenamiento: 0.5	0.012	0.012	0.011	0.009
Mean global error, entrenamiento: 0.6	0.214	0.223	0.196	0.224
std, entrenamiento: 0.6	0.023	0.032	0.024	0.019
Mean Yes error, entrenamiento: 0.6	0.407	0.415	0.385	0.484
std, entrenamiento: 0.6	0.015	0.016	0.015	0.020
Mean No error, entrenamiento: 0.6	0.111	0.109	0.111	0.077
std, entrenamiento: 0.6	0.008	0.016	0.006	0.010
Mean global error, entrenamiento: 0.7	0.216	0.210	0.207	0.223
std, entrenamiento: 0.7	0.018	0.021	0.023	0.032
Mean Yes error, entrenamiento: 0.7	0.412	0.413	0.381	0.496
std, entrenamiento: 0.7	0.018	0.015	0.009	0.024
Mean No error, entrenamiento: 0.7	0.106	0.104	0.117	0.071
std, entrenamiento: 0.7	0.014	0.010	0.008	0.009
Mean global error, entrenamiento: 0.8	0.190	0.218	0.207	0.228
std, entrenamiento: 0.8	0.034	0.025	0.035	0.043
Mean Yes error, entrenamiento: 0.8	0.410	0.415	0.388	0.484
std, entrenamiento: 0.8	0.013	0.014	0.010	0.016
Mean No error, entrenamiento: 0.8	0.100	0.107	0.114	0.072

Continued on next page

	LDA	Naive Bayes	Logistic	SVM, kernel = poly, deg = 5
std, entrenamiento: 0.8	0.008	0.006	0.008	0.004
Mean global error, entrenamiento: 0.9	0.224	0.226	0.209	0.222
std, entrenamiento: 0.9	0.044	0.055	0.036	0.043
Mean Yes error, entrenamiento: 0.9	0.418	0.406	0.383	0.481
std, entrenamiento: 0.9	0.007	0.011	0.009	0.009
Mean No error, entrenamiento: 0.9	0.104	0.106	0.114	0.071
std, entrenamiento: 0.9	0.008	0.006	0.005	0.004
Mean global error, entrenamiento: 1.0	0.205	0.205	0.207	0.205
std, entrenamiento: 1.0	0.000	0.000	0.000	0.000
Mean Yes error, entrenamiento: 1.0	0.412	0.412	0.390	0.475
std, entrenamiento: 1.0	0.000	0.000	0.000	0.000
Mean No error, entrenamiento: 1.0	0.101	0.101	0.115	0.070
std, entrenamiento: 1.0	0.000	0.000	0.000	0.000

Tabla 30: Tasas de error locales y globales para varios clasificadores con conjuntos de entrenamiento divididos sin mantener proporcionalidad

Problema 2

	LDA	QDA	Naive Bayes	Logistic	KNC, k=5	SVM
Mean global error, entrenamiento: 0.1	0.258	0.378	0.295	0.214	0.187	0.215
std, entrenamiento: 0.1	0.054	0.079	0.081	0.033	0.032	0.040
Mean Yes error, entrenamiento: 0.1	0.293	0.307	0.207	0.262	0.214	0.246
std, entrenamiento: 0.1	0.045	0.273	0.080	0.062	0.049	0.092
Mean No error, entrenamiento: 0.1	0.193	0.381	0.360	0.149	0.150	0.155
std, entrenamiento: 0.1	0.081	0.232	0.193	0.053	0.050	0.068
Mean global error, entrenamiento: 0.2	0.207	0.320	0.252	0.176	0.171	0.180
std, entrenamiento: 0.2	0.028	0.111	0.103	0.033	0.016	0.029

Continued on next page

	LDA	QDA	Naive Bayes	Logistic	KNC, k=5	SVM
Mean Yes error, entrenamiento: 0.2	0.229	0.252	0.172	0.193	0.174	0.177
std, entrenamiento: 0.2	0.068	0.250	0.047	0.041	0.037	0.051
Mean No error, entrenamiento: 0.2	0.150	0.335	0.294	0.126	0.154	0.136
std, entrenamiento: 0.2	0.035	0.276	0.187	0.038	0.026	0.062
Mean global error, entrenamiento: 0.3	0.165	0.245	0.260	0.169	0.181	0.149
std, entrenamiento: 0.3	0.030	0.090	0.122	0.020	0.022	0.011
Mean Yes error, entrenamiento: 0.3	0.182	0.179	0.199	0.187	0.192	0.134
std, entrenamiento: 0.3	0.027	0.087	0.120	0.029	0.047	0.028
Mean No error, entrenamiento: 0.3	0.113	0.211	0.284	0.117	0.130	0.116
std, entrenamiento: 0.3	0.024	0.224	0.245	0.036	0.023	0.028
Mean global error, entrenamiento: 0.4	0.168	0.224	0.263	0.162	0.159	0.150
std, entrenamiento: 0.4	0.033	0.086	0.117	0.025	0.019	0.015
Mean Yes error, entrenamiento: 0.4	0.189	0.139	0.179	0.177	0.167	0.121
std, entrenamiento: 0.4	0.034	0.028	0.070	0.033	0.029	0.023
Mean No error, entrenamiento: 0.4	0.108	0.218	0.291	0.105	0.120	0.116
std, entrenamiento: 0.4	0.028	0.179	0.260	0.037	0.037	0.022
Mean global error, entrenamiento: 0.5	0.160	0.214	0.196	0.167	0.152	0.138
std, entrenamiento: 0.5	0.017	0.053	0.036	0.028	0.026	0.035
Mean Yes error, entrenamiento: 0.5	0.153	0.166	0.195	0.178	0.158	0.134
std, entrenamiento: 0.5	0.030	0.028	0.037	0.034	0.032	0.035
Mean No error, entrenamiento: 0.5	0.116	0.146	0.167	0.109	0.110	0.090
std, entrenamiento: 0.5	0.031	0.111	0.017	0.026	0.029	0.022
Mean global error, entrenamiento: 0.6	0.158	0.173	0.182	0.162	0.155	0.156
std, entrenamiento: 0.6	0.032	0.034	0.029	0.039	0.019	0.033
Mean Yes error, entrenamiento: 0.6	0.175	0.118	0.164	0.164	0.166	0.130
std, entrenamiento: 0.6	0.016	0.032	0.017	0.024	0.027	0.032
Mean No error, entrenamiento: 0.6	0.101	0.125	0.169	0.109	0.091	0.093
std, entrenamiento: 0.6	0.027	0.021	0.020	0.023	0.025	0.025

Continued on next page

	LDA	QDA	Naive Bayes	Logistic	KNC, k=5	SVM
Mean global error, entrenamiento: 0.7	0.155	0.158	0.168	0.168	0.149	0.151
std, entrenamiento: 0.7	0.027	0.040	0.040	0.033	0.042	0.028
Mean Yes error, entrenamiento: 0.7	0.159	0.094	0.166	0.154	0.166	0.106
std, entrenamiento: 0.7	0.017	0.025	0.017	0.018	0.039	0.020
Mean No error, entrenamiento: 0.7	0.095	0.121	0.164	0.116	0.086	0.100
std, entrenamiento: 0.7	0.021	0.017	0.016	0.022	0.020	0.021
Mean global error, entrenamiento: 0.8	0.127	0.196	0.204	0.144	0.173	0.156
std, entrenamiento: 0.8	0.040	0.045	0.053	0.064	0.032	0.038
Mean Yes error, entrenamiento: 0.8	0.152	0.114	0.165	0.160	0.174	0.129
std, entrenamiento: 0.8	0.009	0.016	0.015	0.018	0.019	0.025
Mean No error, entrenamiento: 0.8	0.096	0.120	0.163	0.091	0.089	0.074
std, entrenamiento: 0.8	0.013	0.014	0.021	0.016	0.007	0.016
Mean global error, entrenamiento: 0.9	0.158	0.204	0.196	0.171	0.113	0.125
std, entrenamiento: 0.9	0.069	0.057	0.046	0.088	0.079	0.059
Mean Yes error, entrenamiento: 0.9	0.150	0.114	0.167	0.150	0.161	0.111
std, entrenamiento: 0.9	0.012	0.013	0.011	0.014	0.020	0.023
Mean No error, entrenamiento: 0.9	0.098	0.119	0.155	0.119	0.078	0.078
std, entrenamiento: 0.9	0.014	0.016	0.011	0.013	0.014	0.020
Mean global error, entrenamiento: 1.0	0.114	0.097	0.153	0.127	0.123	0.089
std, entrenamiento: 1.0	0.000	0.000	0.000	0.000	0.000	0.000
Mean Yes error, entrenamiento: 1.0	0.140	0.093	0.150	0.150	0.178	0.112
std, entrenamiento: 1.0	0.000	0.000	0.000	0.000	0.000	0.000
Mean No error, entrenamiento: 1.0	0.093	0.101	0.155	0.109	0.078	0.070
std, entrenamiento: 1.0	0.000	0.000	0.000	0.000	0.000	0.000

Tabla 32: Tasas de error locales y globales para varios clasificadores con conjuntos de entrenamiento divididos sin mantener proporcionalidad

Coefficient	Estimate	Std. Error	z value	Pr(> z)
(Intercept)	-1.9907	1.3717	-1.451	0.146690
SexMale	0.4269	0.5242	0.814	0.415394
AngPecNone	-0.9167	0.7861	-1.166	0.243542
AngPecTypical	2.0258	0.6704	3.022	0.002514
AMINotCertain	-1.8638	0.5475	-3.404	0.000663
QWaveYes	1.7011	0.5516	3.084	0.002043
QWavecodeUsable	0.8020	1.0821	0.741	0.458595
STcodeUsable	-1.9136	0.7652	-2.501	0.012394
STchangeYes	2.4963	0.6071	4.112	3.92e-05
SuffHeartFYes	0.3750	0.5581	0.672	0.501609
HypertrophYes	-0.7961	0.5992	-1.329	0.183919
HypercholYes	1.2769	0.4521	2.825	0.004733
SmokerYes	0.3492	0.5504	0.634	0.525844
InheritYes	0.5909	0.4846	1.219	0.222732
HeartfailYes	-0.8047	0.6662	-1.208	0.227045

Tabla 33: Coeficientes y p -values para la regresión logística del modelo completo

	Naive Bayes	Logistic	SVM, kernel = poly, deg = 5
Mean global error, en- trenamiento: 0.1	0.281	0.204	0.200
std, entrenamiento: 0.1	0.095	0.041	0.028
Mean Yes error, entre- namiento: 0.1	0.285	0.277	0.272
std, entrenamiento: 0.1	0.173	0.125	0.062
Mean No error, entrena- miento: 0.1	0.253	0.135	0.105
std, entrenamiento: 0.1	0.217	0.047	0.037
Mean global error, en- trenamiento: 0.2	0.332	0.192	0.201
std, entrenamiento: 0.2	0.104	0.032	0.031
Mean Yes error, entre- namiento: 0.2	0.179	0.201	0.254
std, entrenamiento: 0.2	0.123	0.050	0.052
Mean No error, entrena- miento: 0.2	0.429	0.160	0.086
std, entrenamiento: 0.2	0.255	0.048	0.026
Mean global error, en- trenamiento: 0.3	0.211	0.192	0.195
std, entrenamiento: 0.3	0.078	0.033	0.027

Continued on next page

	Naive Bayes	Logistic	SVM, kernel = poly, deg = 5
Mean Yes error, entrenamiento: 0.3	0.246	0.229	0.198
std, entrenamiento: 0.3	0.199	0.061	0.061
Mean No error, entrenamiento: 0.3	0.153	0.159	0.093
std, entrenamiento: 0.3	0.035	0.042	0.037
Mean global error, entrenamiento: 0.4	0.292	0.182	0.189
std, entrenamiento: 0.4	0.131	0.039	0.025
Mean Yes error, entrenamiento: 0.4	0.137	0.164	0.173
std, entrenamiento: 0.4	0.078	0.044	0.040
Mean No error, entrenamiento: 0.4	0.378	0.179	0.083
std, entrenamiento: 0.4	0.291	0.019	0.021
Mean global error, entrenamiento: 0.5	0.217	0.190	0.171
std, entrenamiento: 0.5	0.078	0.033	0.039
Mean Yes error, entrenamiento: 0.5	0.174	0.176	0.121
std, entrenamiento: 0.5	0.051	0.038	0.030
Mean No error, entrenamiento: 0.5	0.210	0.158	0.078
std, entrenamiento: 0.5	0.185	0.030	0.024
Mean global error, entrenamiento: 0.6	0.165	0.169	0.183
std, entrenamiento: 0.6	0.035	0.039	0.029
Mean Yes error, entrenamiento: 0.6	0.153	0.199	0.116
std, entrenamiento: 0.6	0.017	0.045	0.031
Mean No error, entrenamiento: 0.6	0.156	0.157	0.070
std, entrenamiento: 0.6	0.021	0.025	0.013
Mean global error, entrenamiento: 0.7	0.170	0.201	0.200
std, entrenamiento: 0.7	0.042	0.037	0.039
Mean Yes error, entrenamiento: 0.7	0.171	0.196	0.104
std, entrenamiento: 0.7	0.016	0.030	0.031
Mean No error, entrenamiento: 0.7	0.154	0.153	0.060

Continued on next page

	Naive Bayes	Logistic	SVM, kernel = poly, deg = 5
std, entrenamiento: 0.7	0.022	0.031	0.027
Mean global error, en- trenamiento: 0.8	0.169	0.188	0.204
std, entrenamiento: 0.8	0.065	0.055	0.057
Mean Yes error, entre- namiento: 0.8	0.169	0.168	0.082
std, entrenamiento: 0.8	0.015	0.044	0.026
Mean No error, entrena- miento: 0.8	0.145	0.155	0.054
std, entrenamiento: 0.8	0.018	0.024	0.010
Mean global error, en- trenamiento: 0.9	0.142	0.175	0.162
std, entrenamiento: 0.9	0.068	0.064	0.044
Mean Yes error, entre- namiento: 0.9	0.160	0.197	0.058
std, entrenamiento: 0.9	0.012	0.045	0.017
Mean No error, entrena- miento: 0.9	0.142	0.136	0.039
std, entrenamiento: 0.9	0.005	0.029	0.010
Mean global error, en- trenamiento: 1.0	0.144	0.169	0.038
std, entrenamiento: 1.0	0.000	0.000	0.000
Mean Yes error, entre- namiento: 1.0	0.150	0.243	0.037
std, entrenamiento: 1.0	0.000	0.000	0.000
Mean No error, entrena- miento: 1.0	0.140	0.109	0.039
std, entrenamiento: 1.0	0.000	0.000	0.000

Tabla 35: Tasas de error locales y globales para varios clasificadores con conjuntos de entrenamiento divididos sin mantener proporcionalidad

Problema 3

	One- versus- others	multinomial
Mean global error, en- trenamiento: 0.1	0.160	0.112
Continued on next page		

	One- versus- others	multinomial
std, entrenamiento: 0.1	0.019	0.029
Mean 1 error, entrenamiento: 0.1	0.141	0.212
std, entrenamiento: 0.1	0.052	0.065
Mean 2 error, entrenamiento: 0.1	0.370	0.093
std, entrenamiento: 0.1	0.092	0.112
Mean 3 error, entrenamiento: 0.1	0.039	0.057
std, entrenamiento: 0.1	0.032	0.022
Mean global error, entrenamiento: 0.2	0.135	0.129
std, entrenamiento: 0.2	0.011	0.012
Mean 1 error, entrenamiento: 0.2	0.071	0.101
std, entrenamiento: 0.2	0.038	0.062
Mean 2 error, entrenamiento: 0.2	0.287	0.278
std, entrenamiento: 0.2	0.047	0.079
Mean 3 error, entrenamiento: 0.2	0.039	0.022
std, entrenamiento: 0.2	0.019	0.006
Mean global error, entrenamiento: 0.3	0.046	0.056
std, entrenamiento: 0.3	0.009	0.024
Mean 1 error, entrenamiento: 0.3	0.071	0.071
std, entrenamiento: 0.3	0.038	0.014
Mean 2 error, entrenamiento: 0.3	0.037	0.056
std, entrenamiento: 0.3	0.026	0.039
Mean 3 error, entrenamiento: 0.3	0.013	0.018
std, entrenamiento: 0.3	0.011	0.016
Mean global error, entrenamiento: 0.4	0.065	0.042
std, entrenamiento: 0.4	0.005	0.024
Mean 1 error, entrenamiento: 0.4	0.071	0.030
std, entrenamiento: 0.4	0.029	0.043
Mean 2 error, entrenamiento: 0.4	0.028	0.028

Continued on next page

	One- versus- others	multinomial
std, entrenamiento: 0.4	0.023	0.023
Mean 3 error, entrenamiento: 0.4	0.031	0.022
std, entrenamiento: 0.4	0.012	0.022
Mean global error, entrenamiento: 0.5	0.055	0.050
std, entrenamiento: 0.5	0.019	0.036
Mean 1 error, entrenamiento: 0.5	0.081	0.061
std, entrenamiento: 0.5	0.029	0.049
Mean 2 error, entrenamiento: 0.5	0.019	0.019
std, entrenamiento: 0.5	0.013	0.026
Mean 3 error, entrenamiento: 0.5	0.009	0.013
std, entrenamiento: 0.5	0.012	0.011
Mean global error, entrenamiento: 0.6	0.092	0.057
std, entrenamiento: 0.6	0.043	0.035
Mean 1 error, entrenamiento: 0.6	0.061	0.040
std, entrenamiento: 0.6	0.049	0.014
Mean 2 error, entrenamiento: 0.6	0.056	0.009
std, entrenamiento: 0.6	0.060	0.013
Mean 3 error, entrenamiento: 0.6	0.018	0.022
std, entrenamiento: 0.6	0.016	0.016
Mean global error, entrenamiento: 0.7	0.076	0.008
std, entrenamiento: 0.7	0.060	0.011
Mean 1 error, entrenamiento: 0.7	0.020	0.000
std, entrenamiento: 0.7	0.014	0.000
Mean 2 error, entrenamiento: 0.7	0.065	0.000
std, entrenamiento: 0.7	0.035	0.000
Mean 3 error, entrenamiento: 0.7	0.022	0.004
std, entrenamiento: 0.7	0.022	0.006
Mean global error, entrenamiento: 0.8	0.057	0.046

Continued on next page

	One- versus- others	multinomial
std, entrenamiento: 0.8	0.016	0.016
Mean 1 error, entrenamiento: 0.8	0.020	0.000
std, entrenamiento: 0.8	0.029	0.000
Mean 2 error, entrenamiento: 0.8	0.046	0.000
std, entrenamiento: 0.8	0.035	0.000
Mean 3 error, entrenamiento: 0.8	0.018	0.018
std, entrenamiento: 0.8	0.012	0.006
Mean global error, entrenamiento: 0.9	0.089	0.022
std, entrenamiento: 0.9	0.031	0.031
Mean 1 error, entrenamiento: 0.9	0.010	0.000
std, entrenamiento: 0.9	0.014	0.000
Mean 2 error, entrenamiento: 0.9	0.102	0.000
std, entrenamiento: 0.9	0.073	0.000
Mean 3 error, entrenamiento: 0.9	0.018	0.004
std, entrenamiento: 0.9	0.012	0.006
Mean global error, entrenamiento: 1.0	0.034	0.000
std, entrenamiento: 1.0	0.000	0.000
Mean 1 error, entrenamiento: 1.0	0.000	0.000
std, entrenamiento: 1.0	0.000	0.000
Mean 2 error, entrenamiento: 1.0	0.083	0.000
std, entrenamiento: 1.0	0.000	0.000
Mean 3 error, entrenamiento: 1.0	0.026	0.000
std, entrenamiento: 1.0	0.000	0.000

Tabla 37: Tasas de error locales y globales para varios clasificadores con conjuntos de entrenamiento divididos sin mantener proporcionalidad

Anexo 2: Figuras relevantes

Problema 1

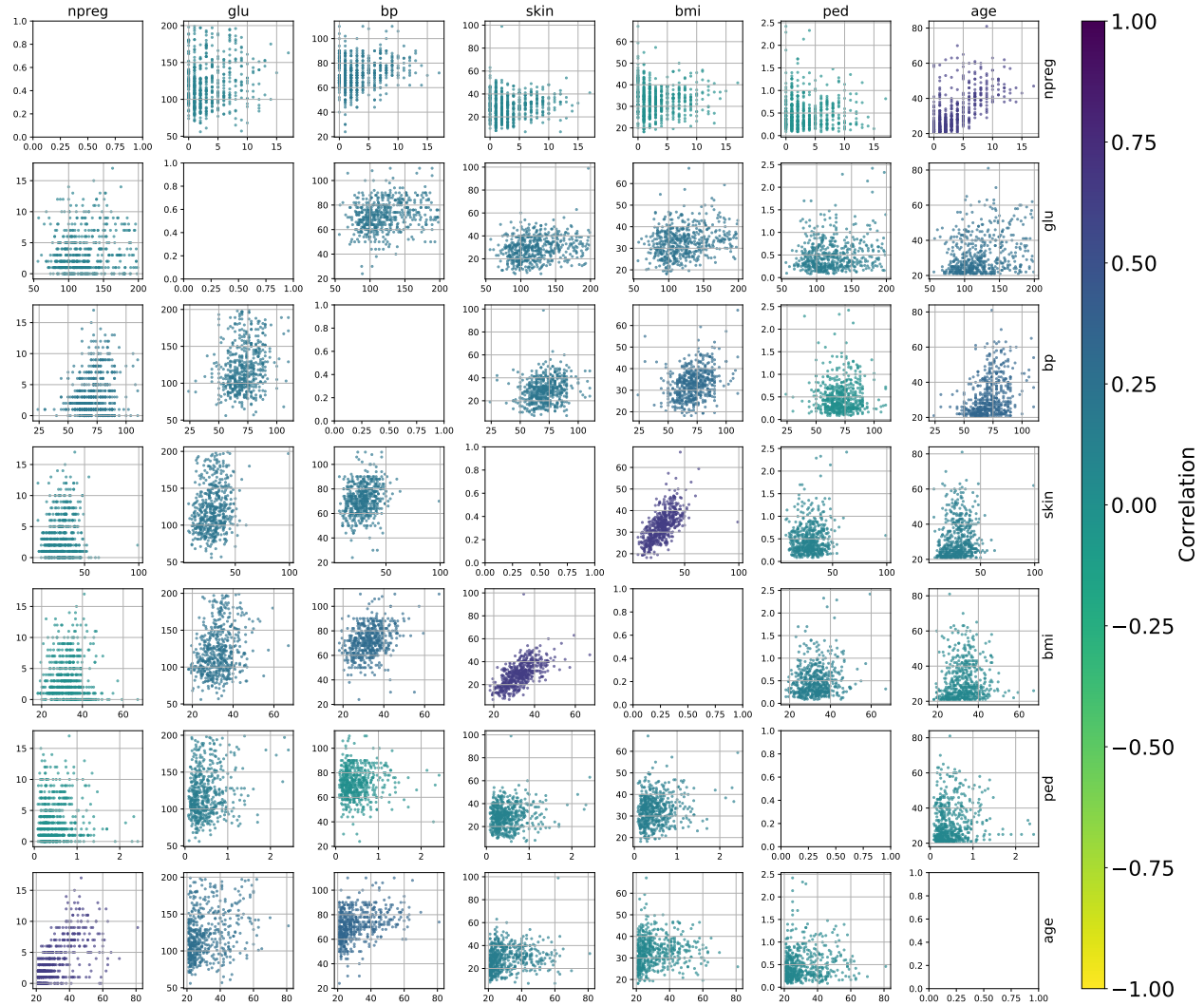


Figura 1: Gráfica de correlación para las variables numericas del problema 1

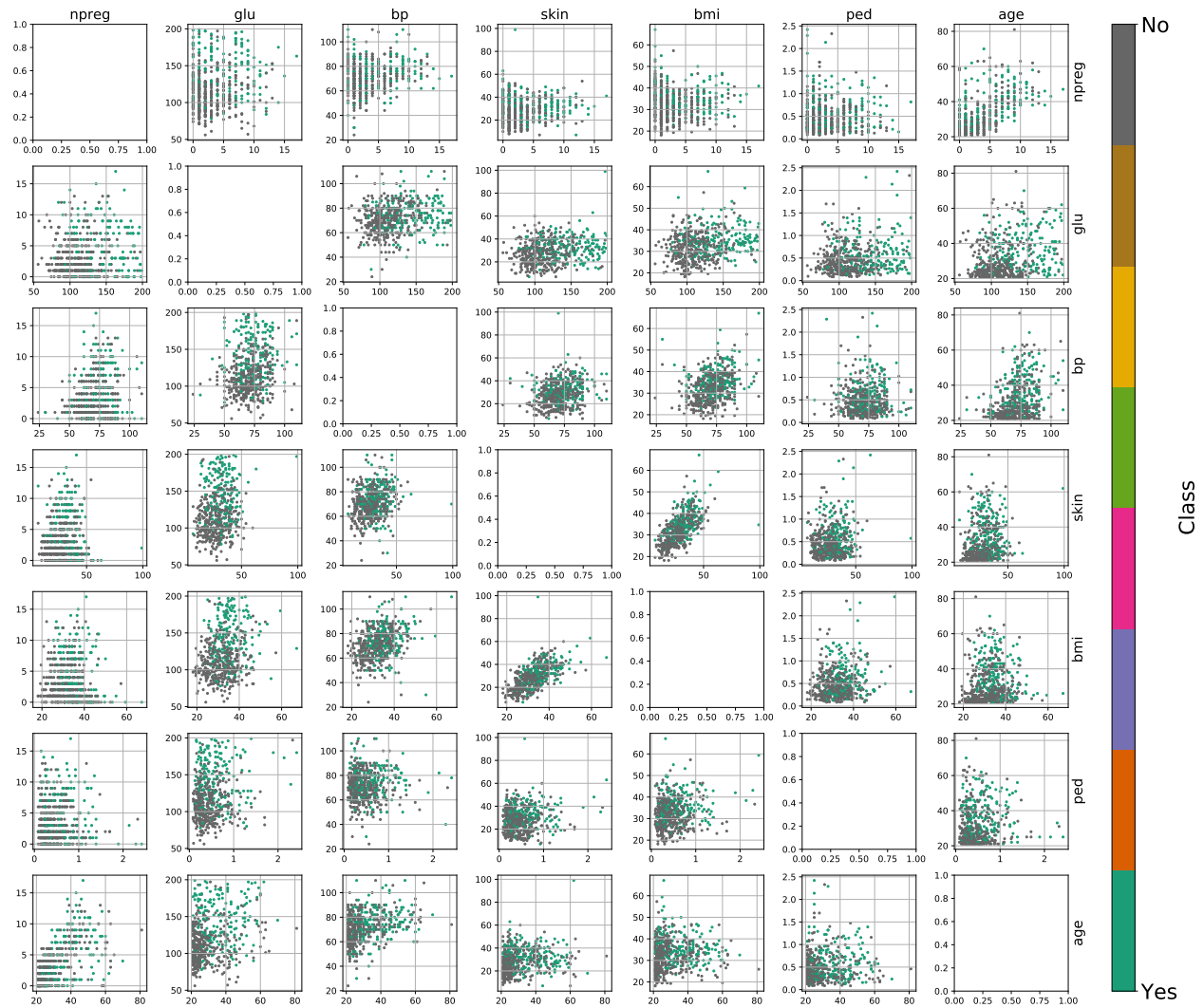


Figura 2: Gráficas de una variable contra otra divididas en clases para las variables numéricas del problema 1

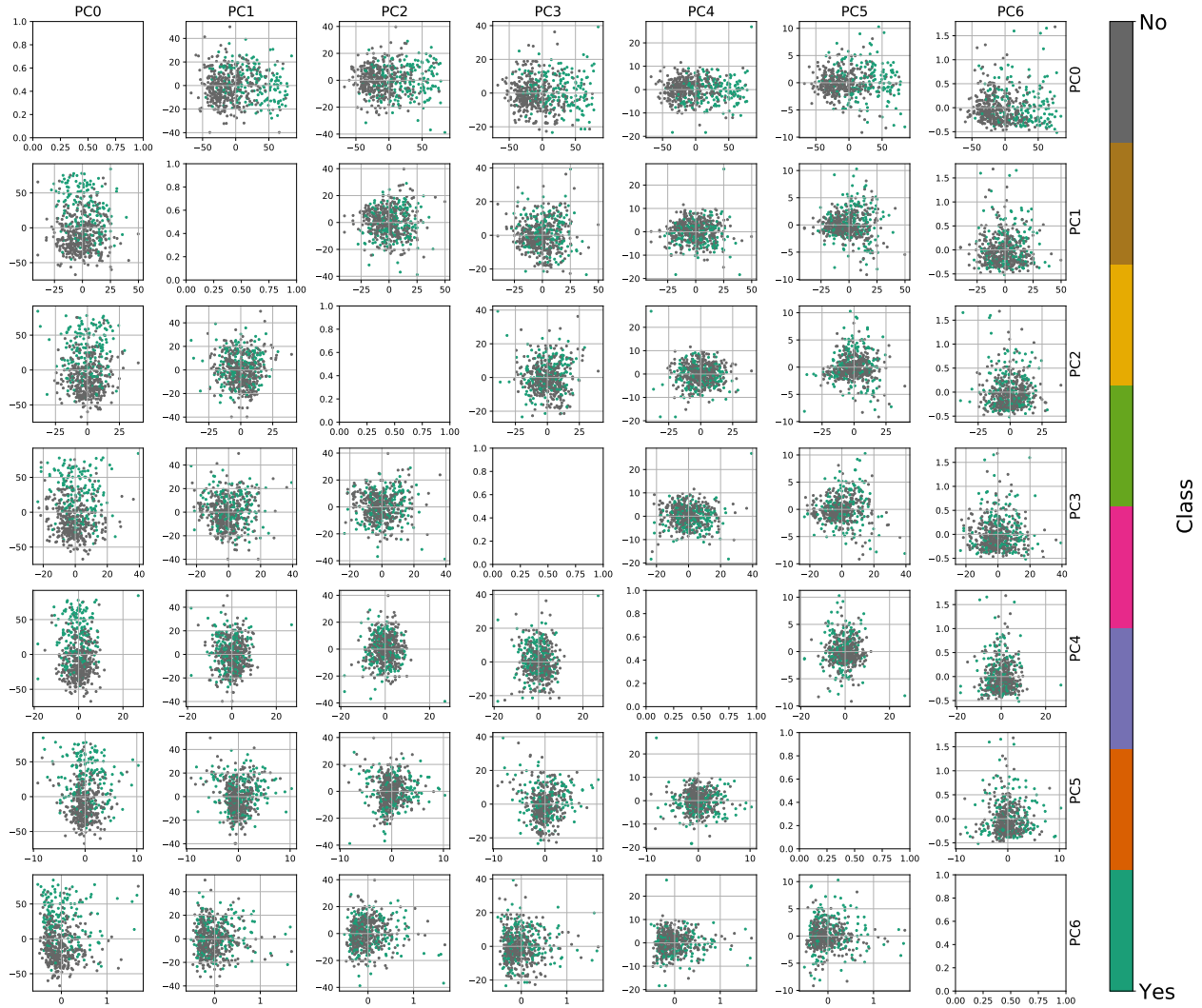


Figura 3: Gráficas de una variable contra otra divididas en clases para las componentes principales del problema 1

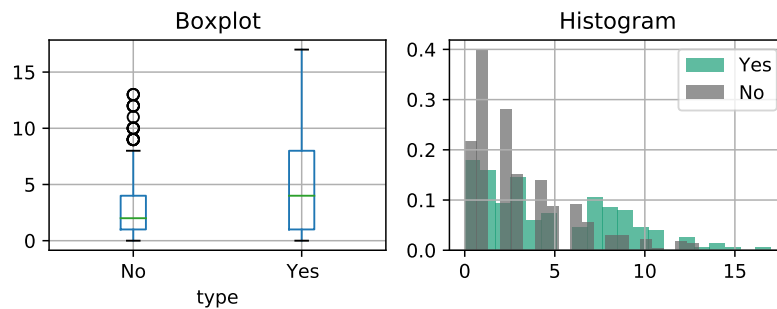


Figura 4: Distribución de la variable “npreg”

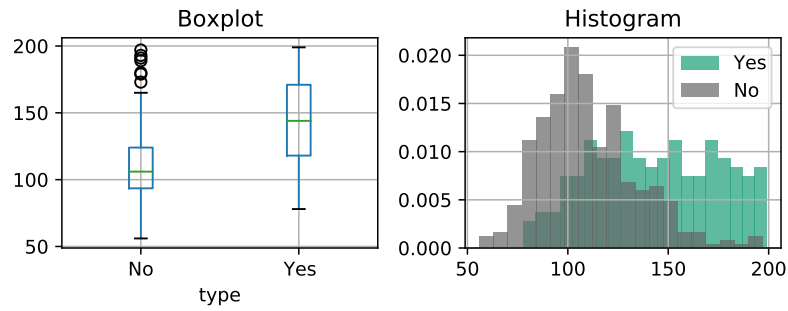


Figura 5: Distribución de la variable “glu”

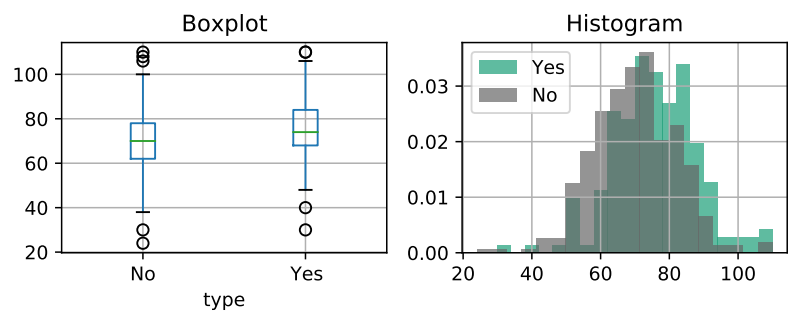


Figura 6: Distribución de la variable “bp”

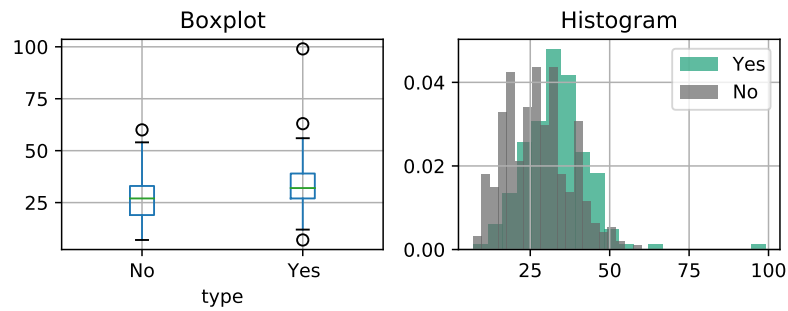


Figura 7: Distribución de la variable “skin”

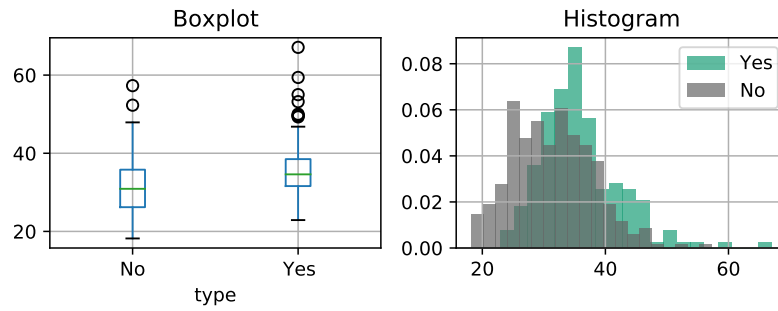


Figura 8: Distribución de la variable “bmi”

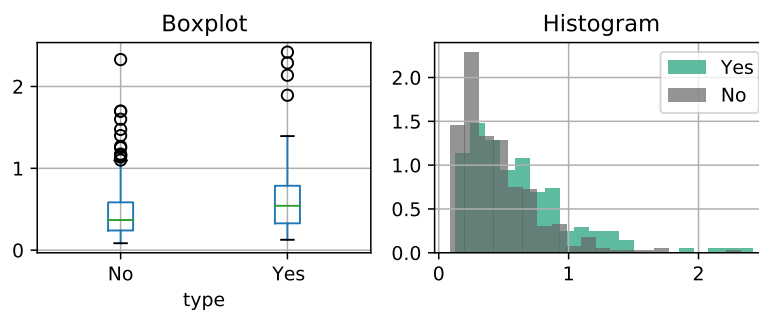


Figura 9: Distribución de la variable “ped”

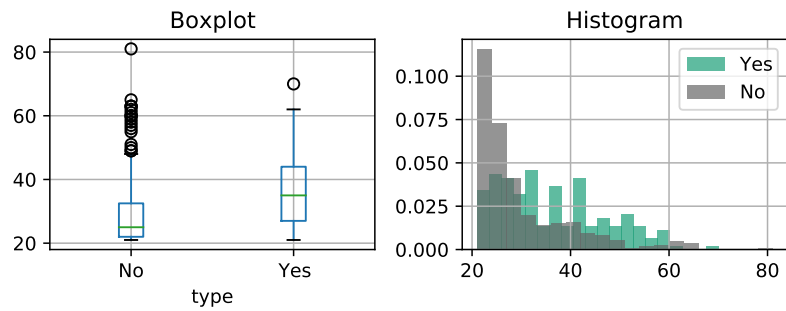


Figura 10: Distribución de la variable “age”

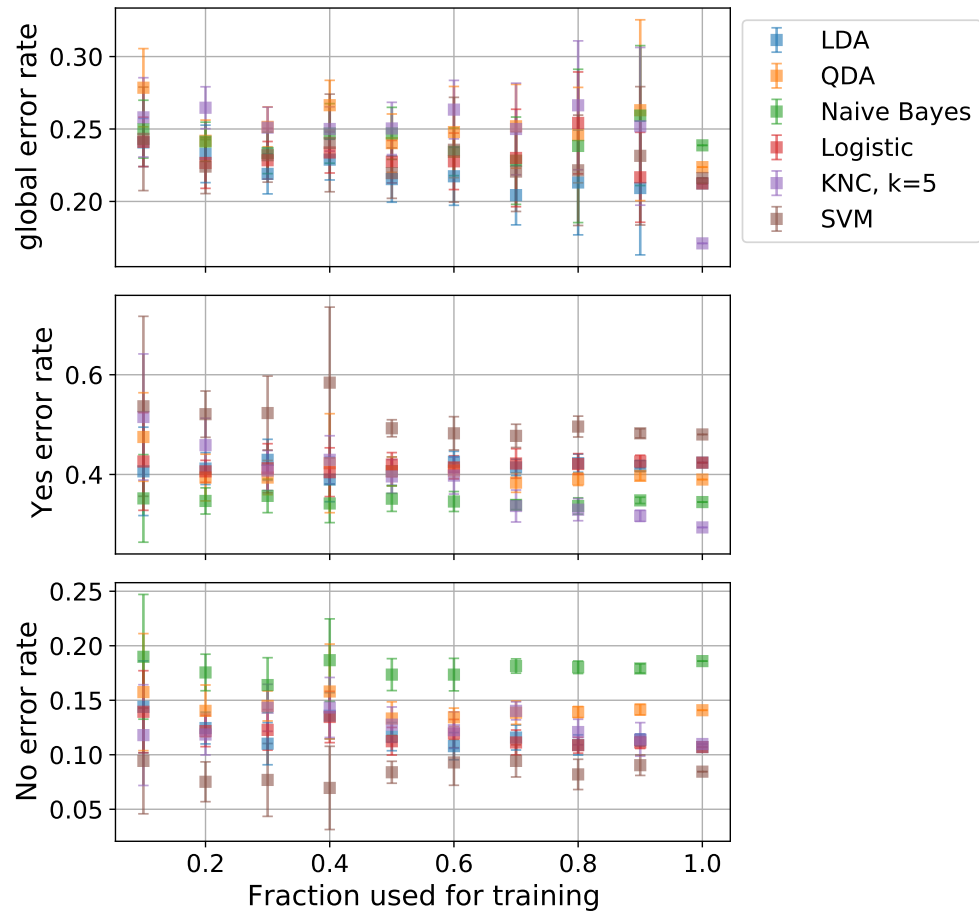


Figura 11: Tasas globales y locales para varios modelos, con conjuntos de entrenamiento divididos sin mantener proporcionalidad

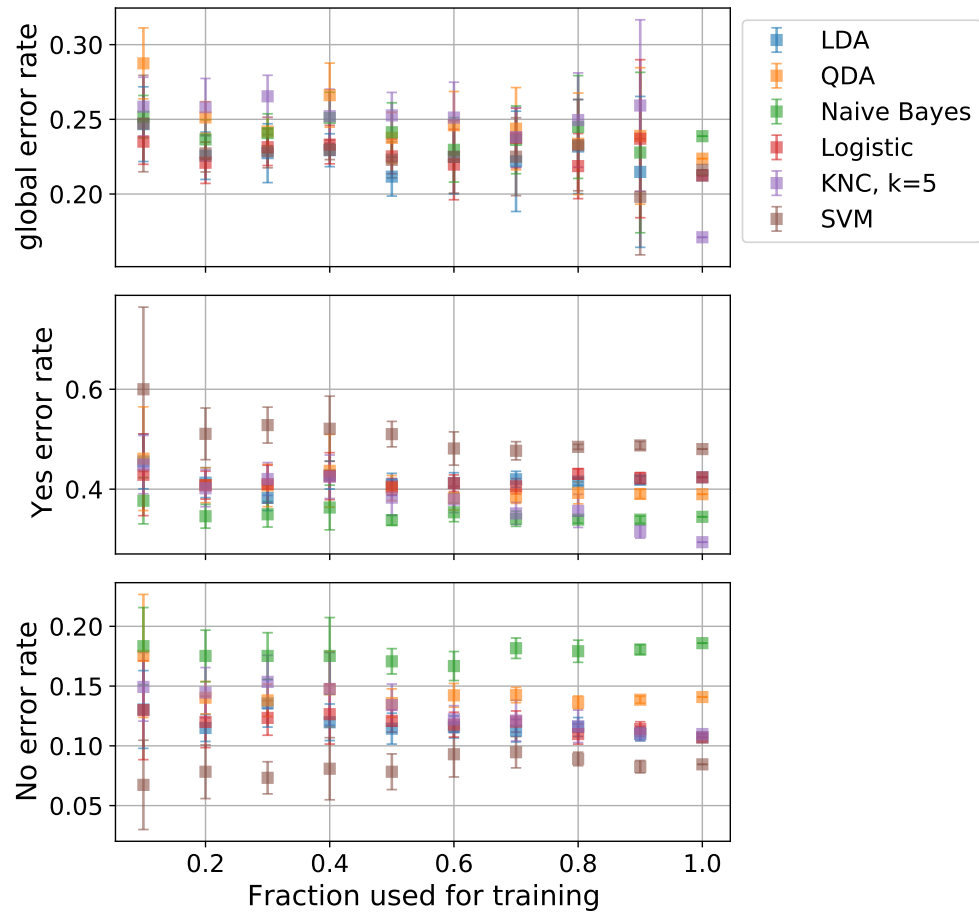


Figura 12: Tasas globales y locales para varios modelos, con conjuntos de entrenamiento divididos manteniendo proporcionalidad

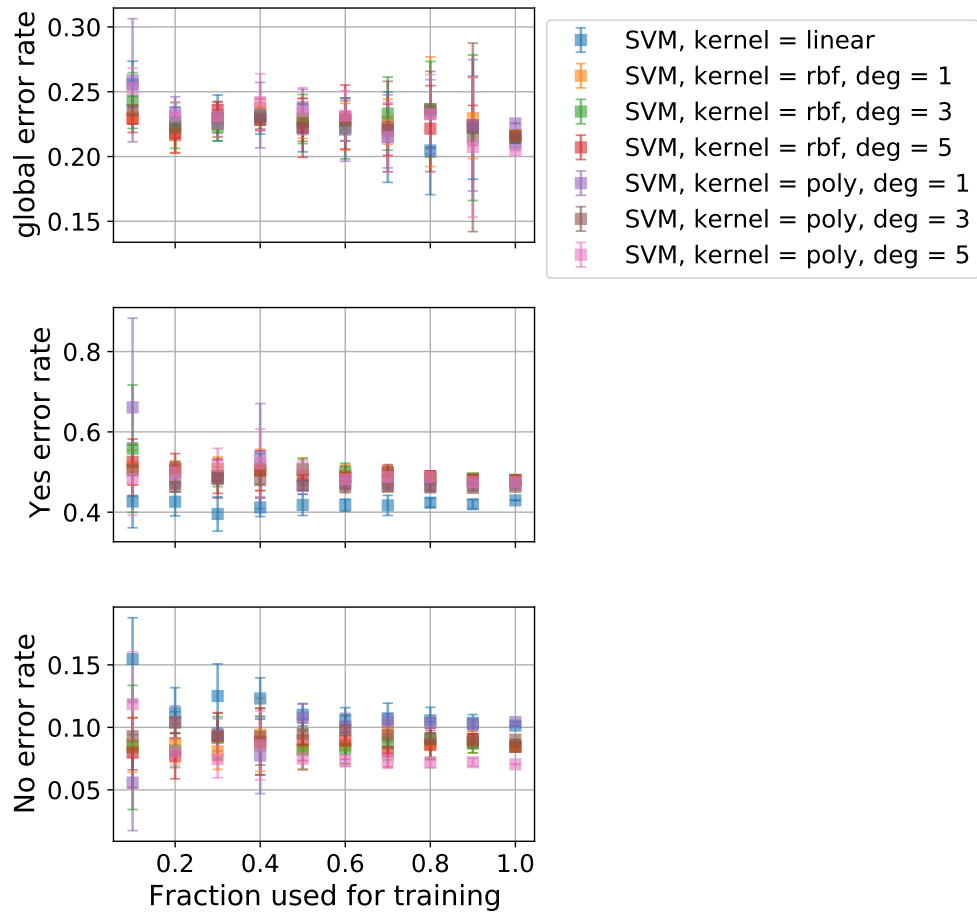


Figura 13: Tasas globales y locales para modelos SVM, con conjuntos de entrenamiento divididos manteniendo proporcionalidad

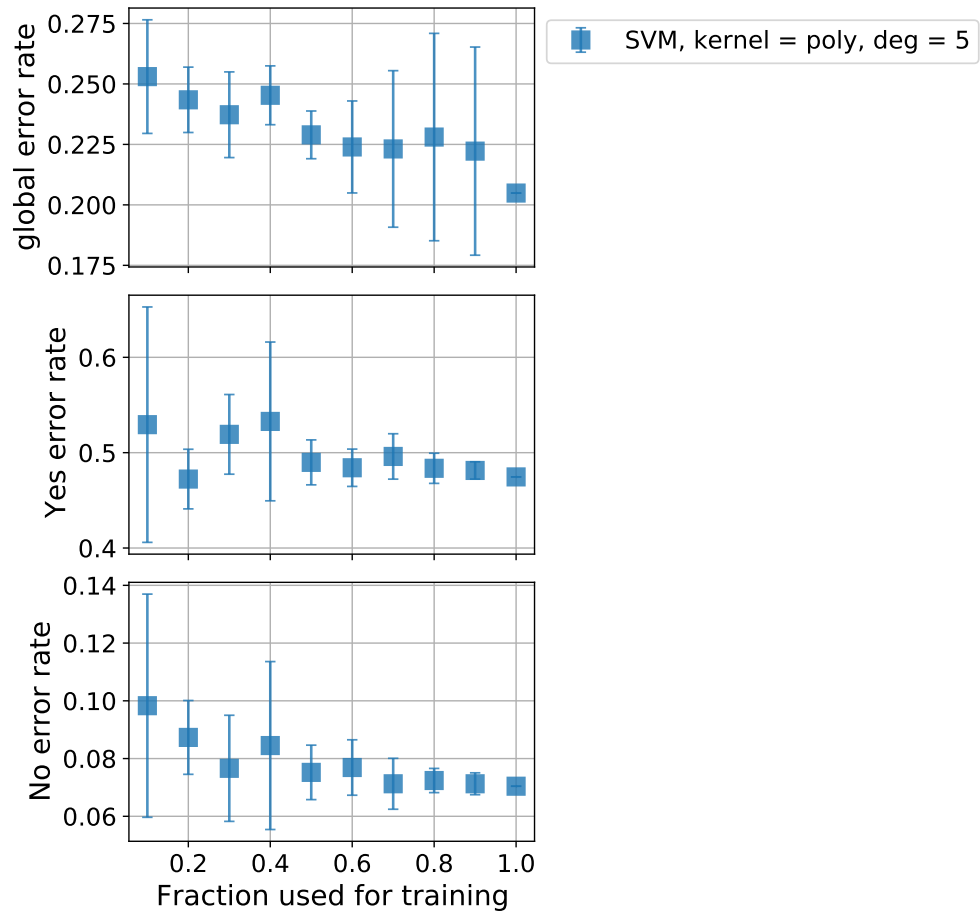


Figura 14: Tasas globales y locales para modelos SVM, con conjuntos de entrenamiento divididos manteniendo proporcionalidad

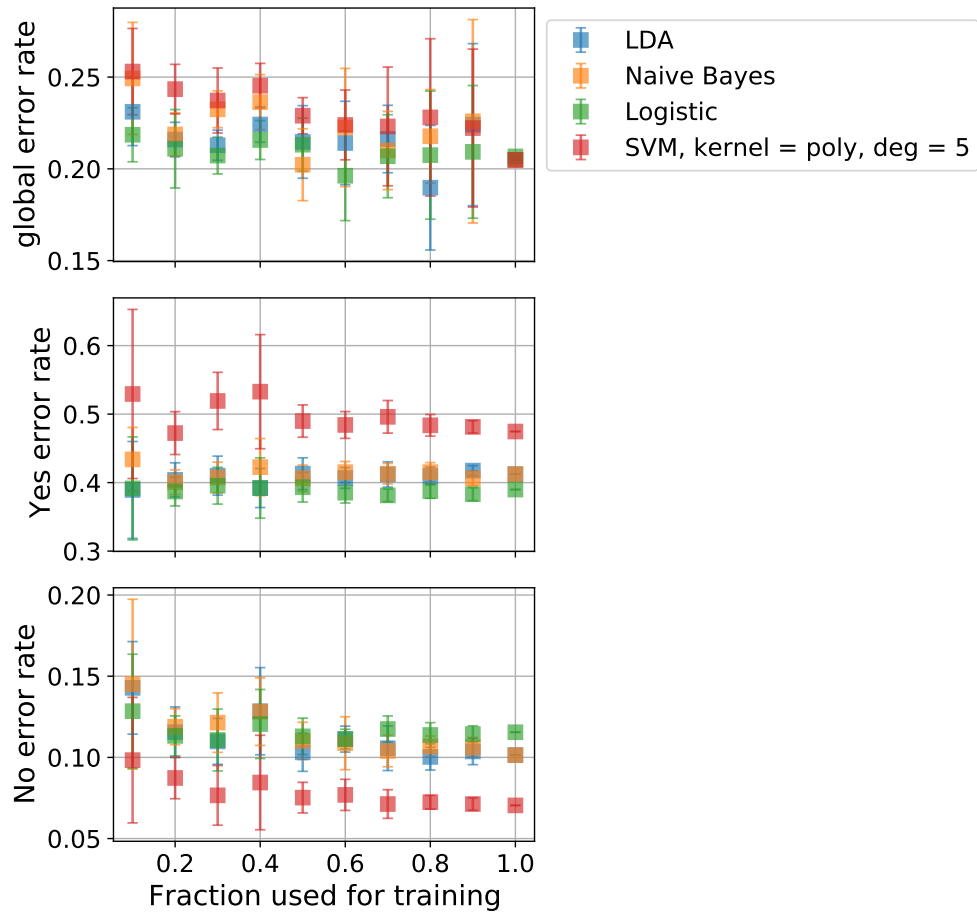


Figura 15: Tasas globales y locales para modelos SVM, con conjuntos de entrenamiento divididos manteniendo proporcionalidad

Problema 2

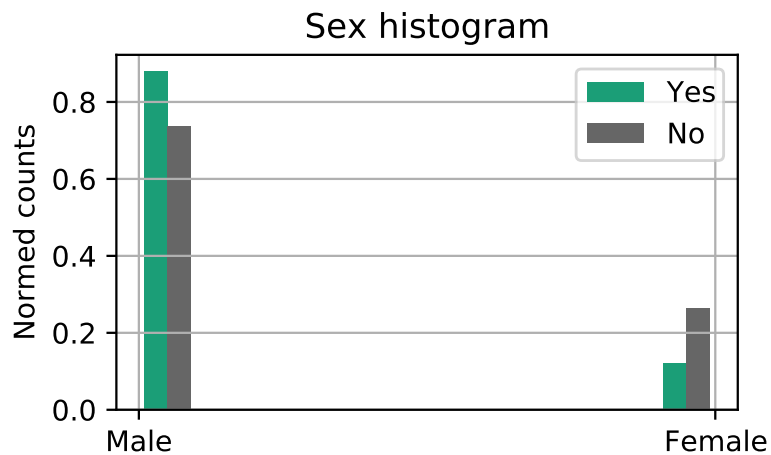


Figura 16: Distribución de la variable "Sex"

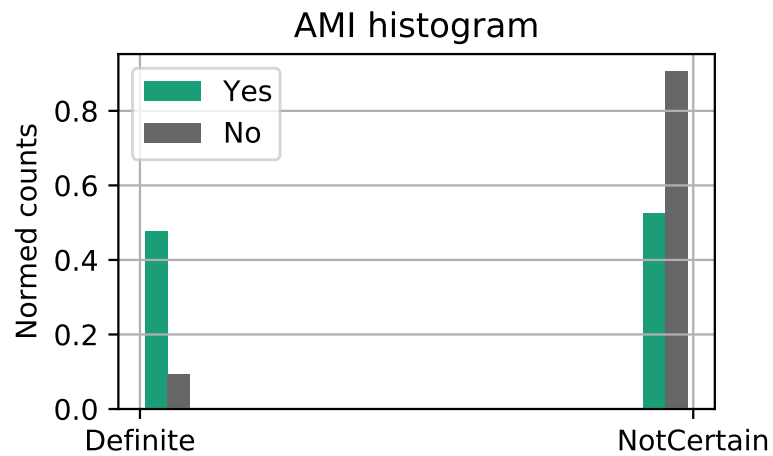


Figura 17: Distribución de la variable “AngPec”

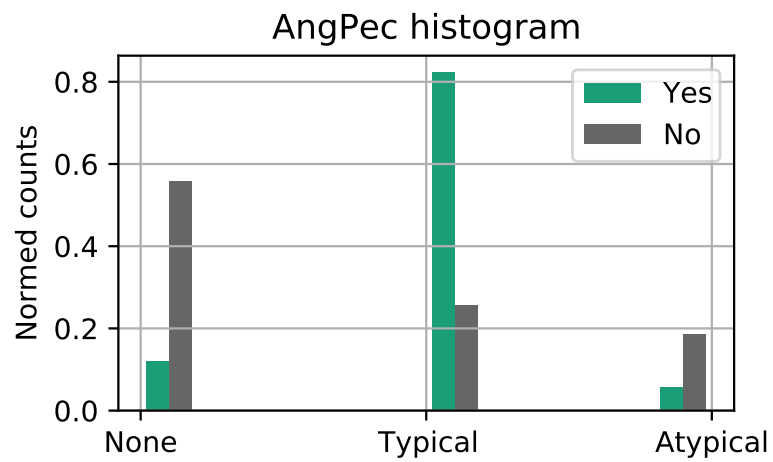


Figura 18: Distribución de la variable “AngPec”

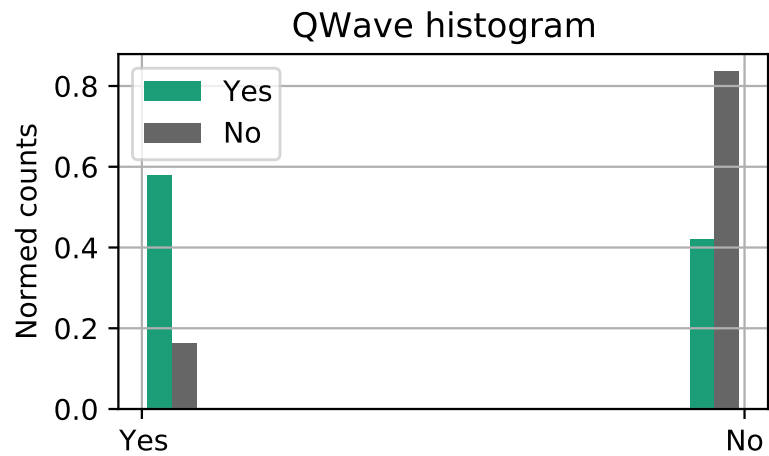


Figura 19: Distribución de la variable “QWave”

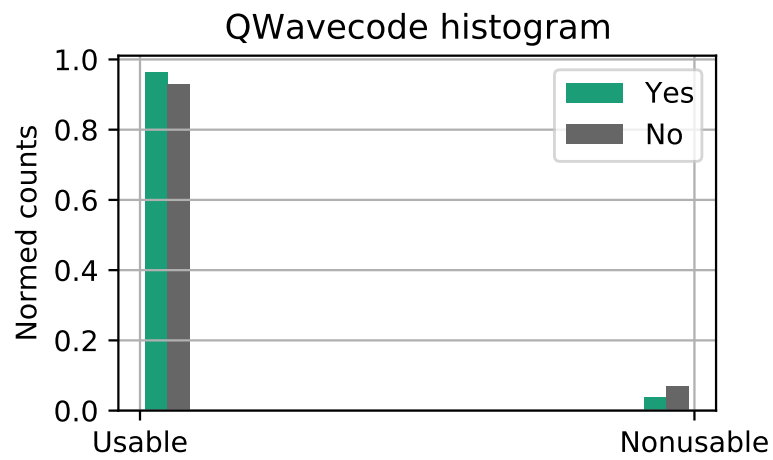


Figura 20: Distribución de la variable “QWavecode”

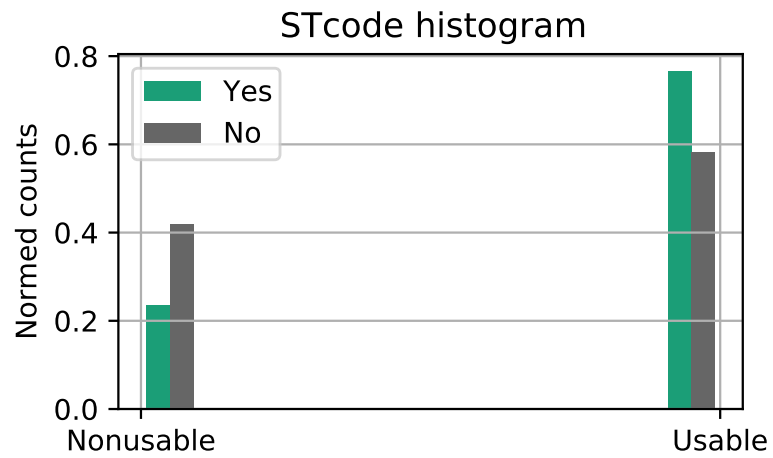


Figura 21: Distribución de la variable “STcode”

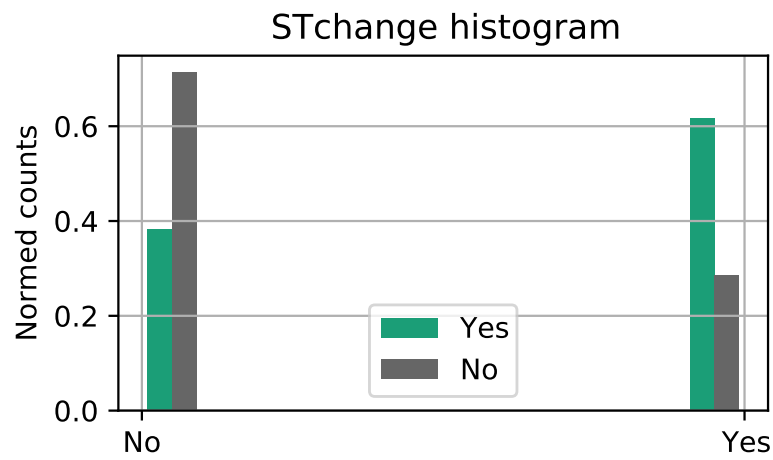


Figura 22: Distribución de la variable “STchange”

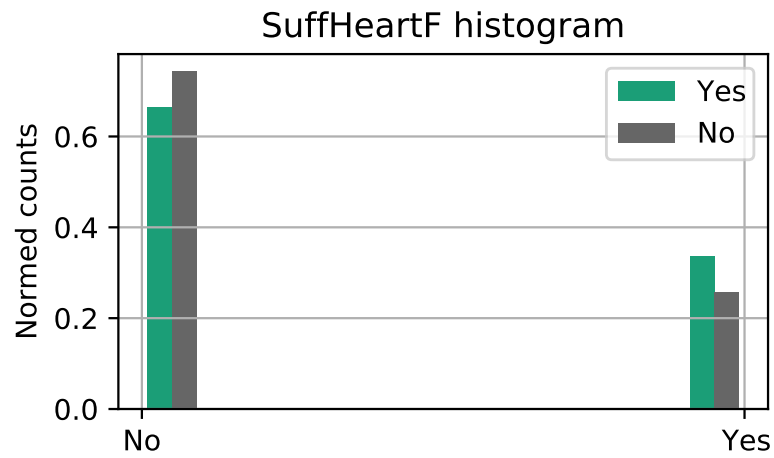


Figura 23: Distribución de la variable “SuffHeartF”

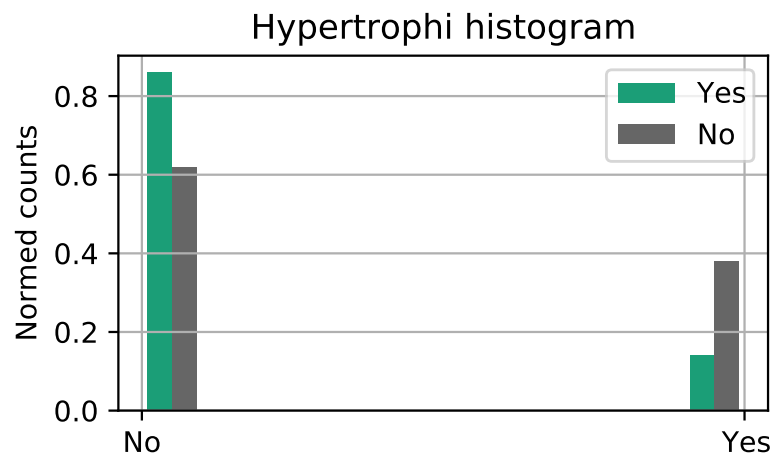


Figura 24: Distribución de la variable “Hypertrophi”

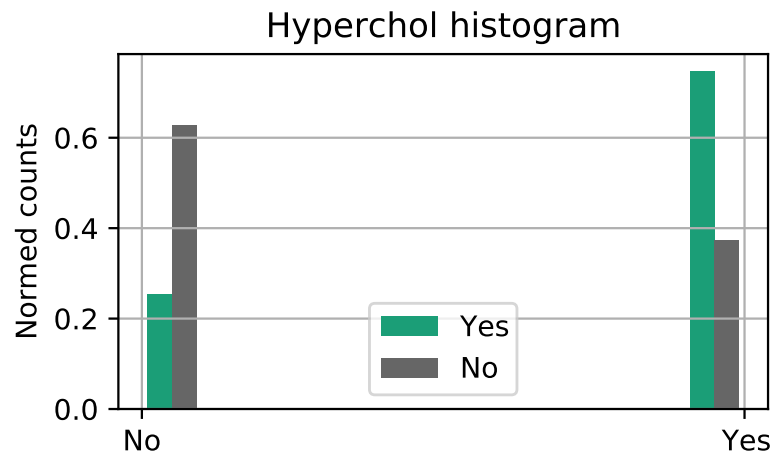


Figura 25: Distribución de la variable “Hyperchol”

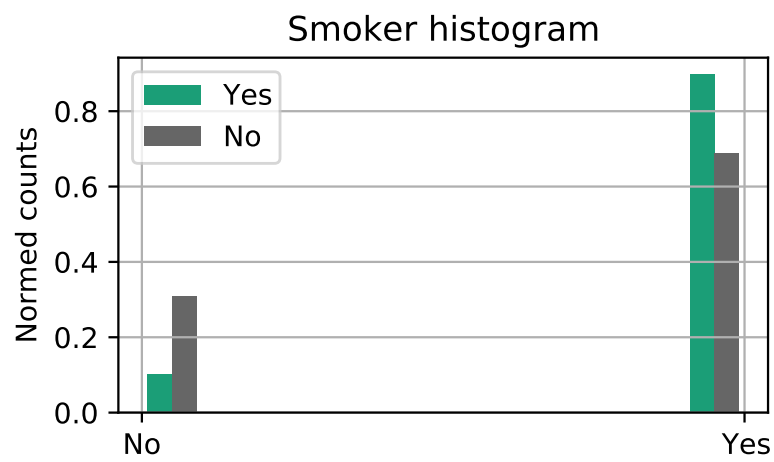


Figura 26: Distribución de la variable “Smoker”

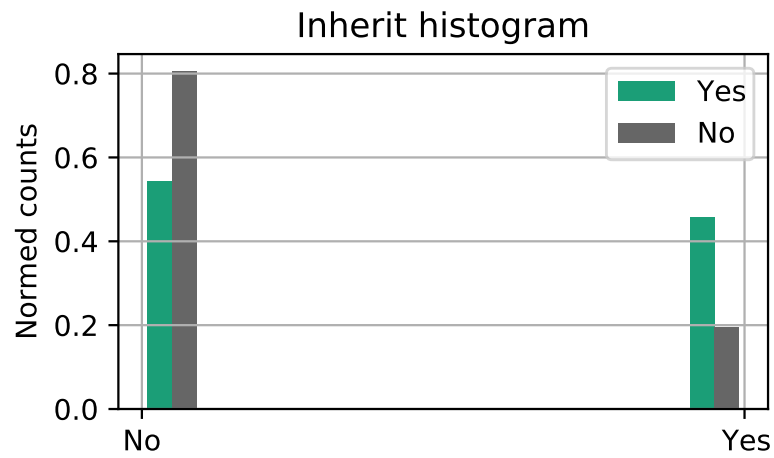


Figura 27: Distribución de la variable “Inherit”

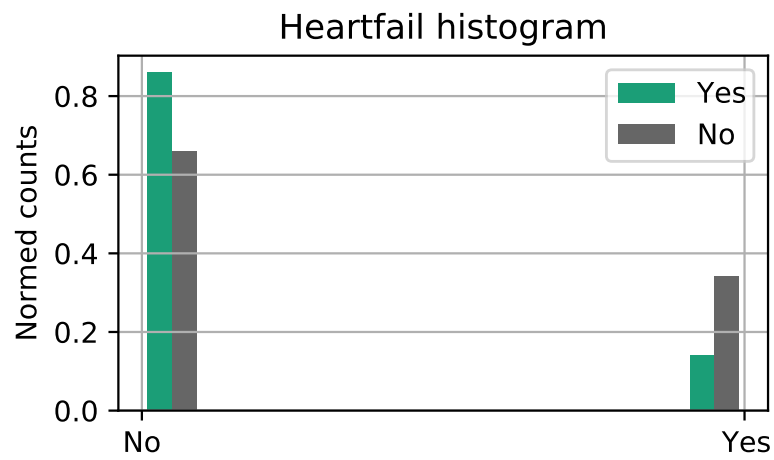


Figura 28: Distribución de la variable “Heartfail”

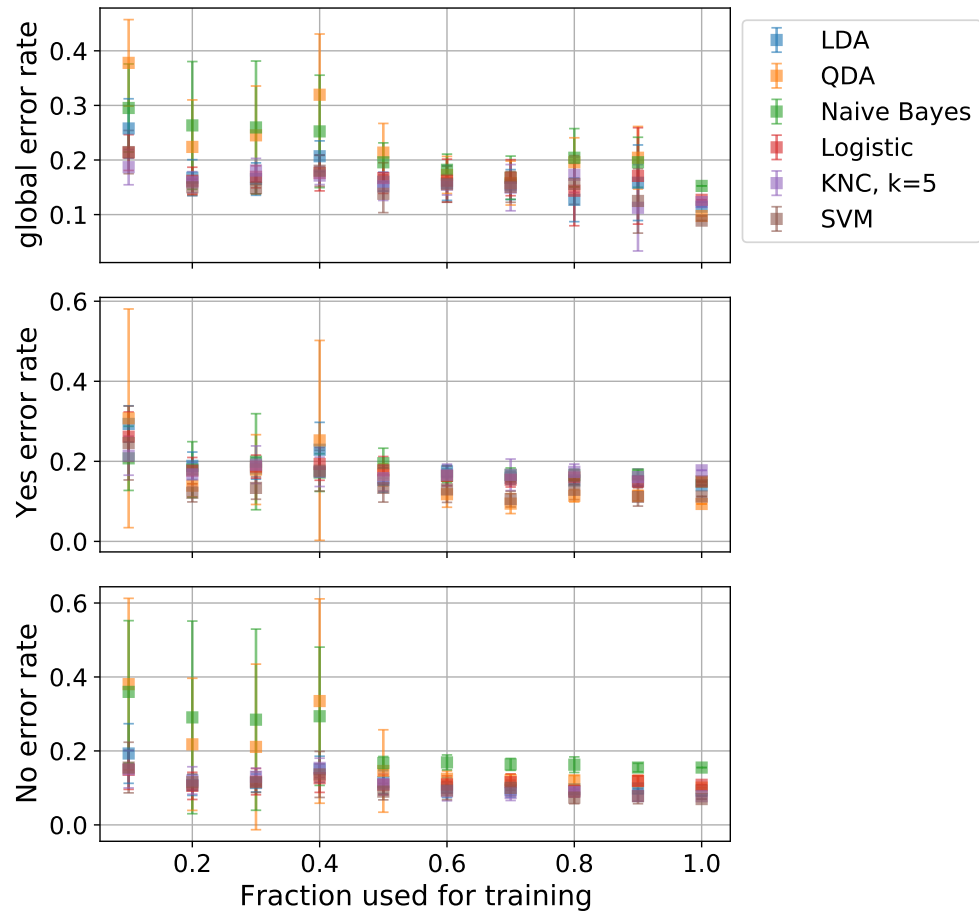


Figura 29: Tasas globales y locales para varios modelos, con conjuntos de entrenamiento divididos manteniendo proporcionalidad

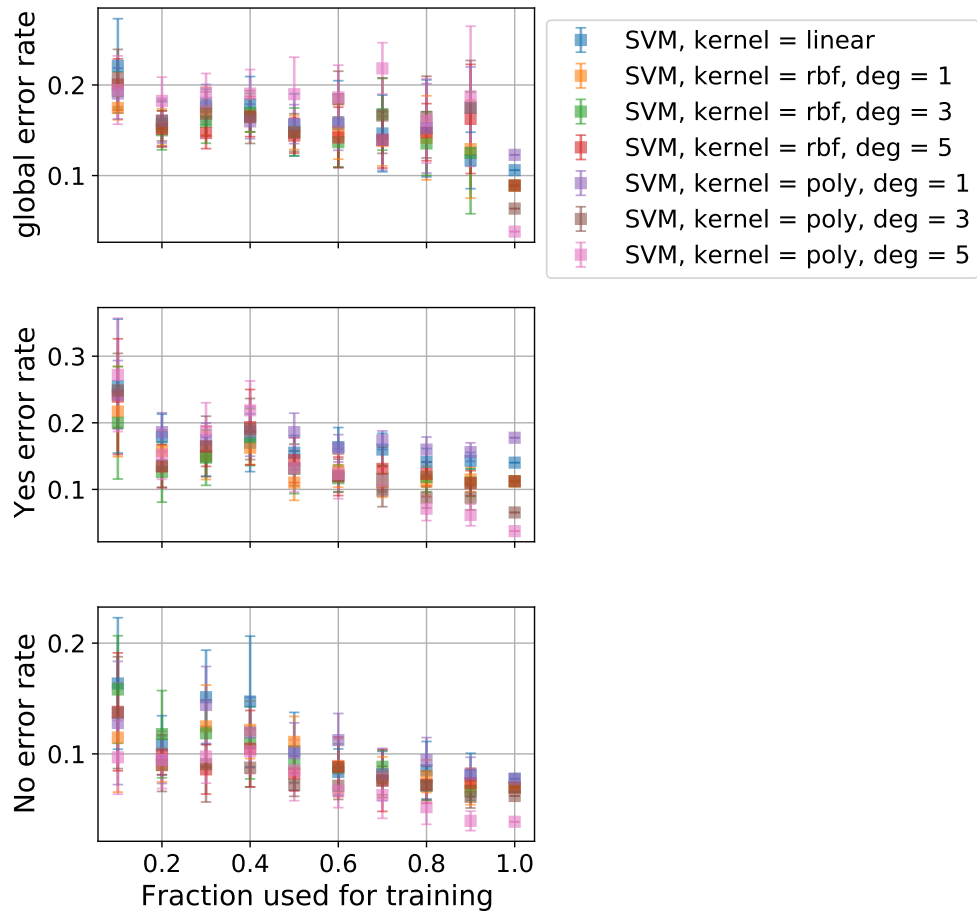


Figura 30: Tasas globales y locales para modelos SVM, con conjuntos de entrenamiento divididos manteniendo proporcionalidad

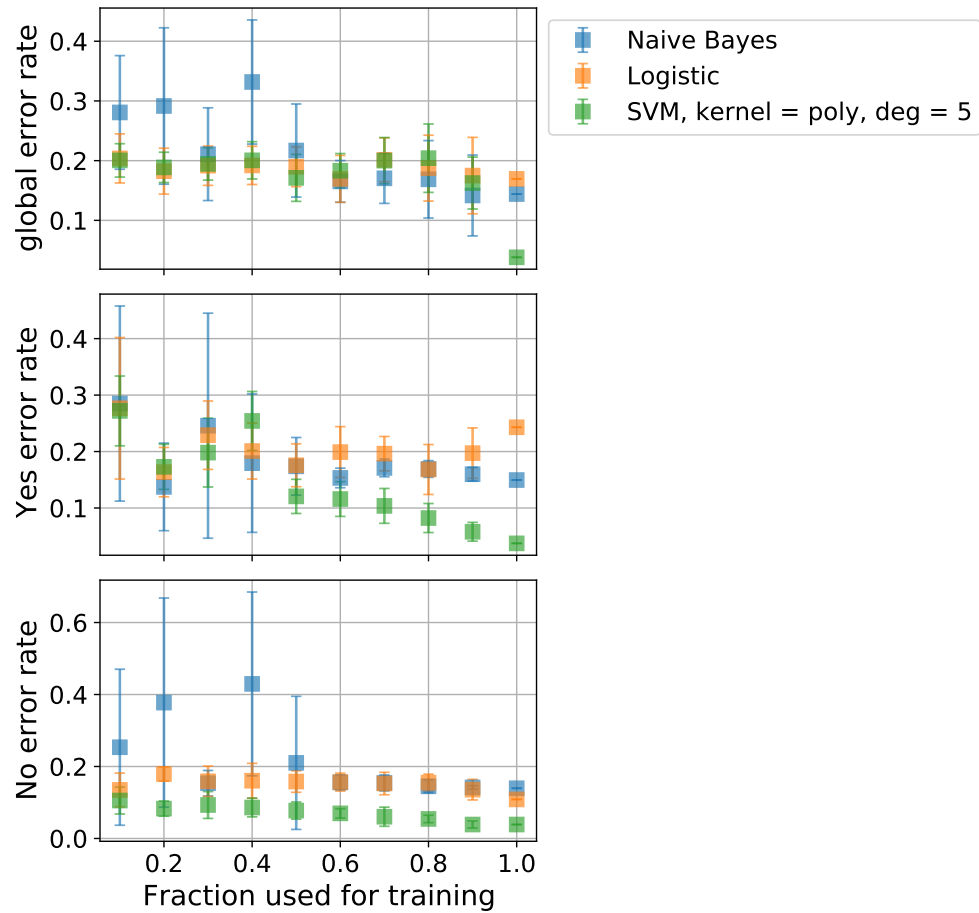


Figura 31: Tasas globales y locales para modelos SVM, con conjuntos de entrenamiento divididos manteniendo proporcionalidad

Problema 3

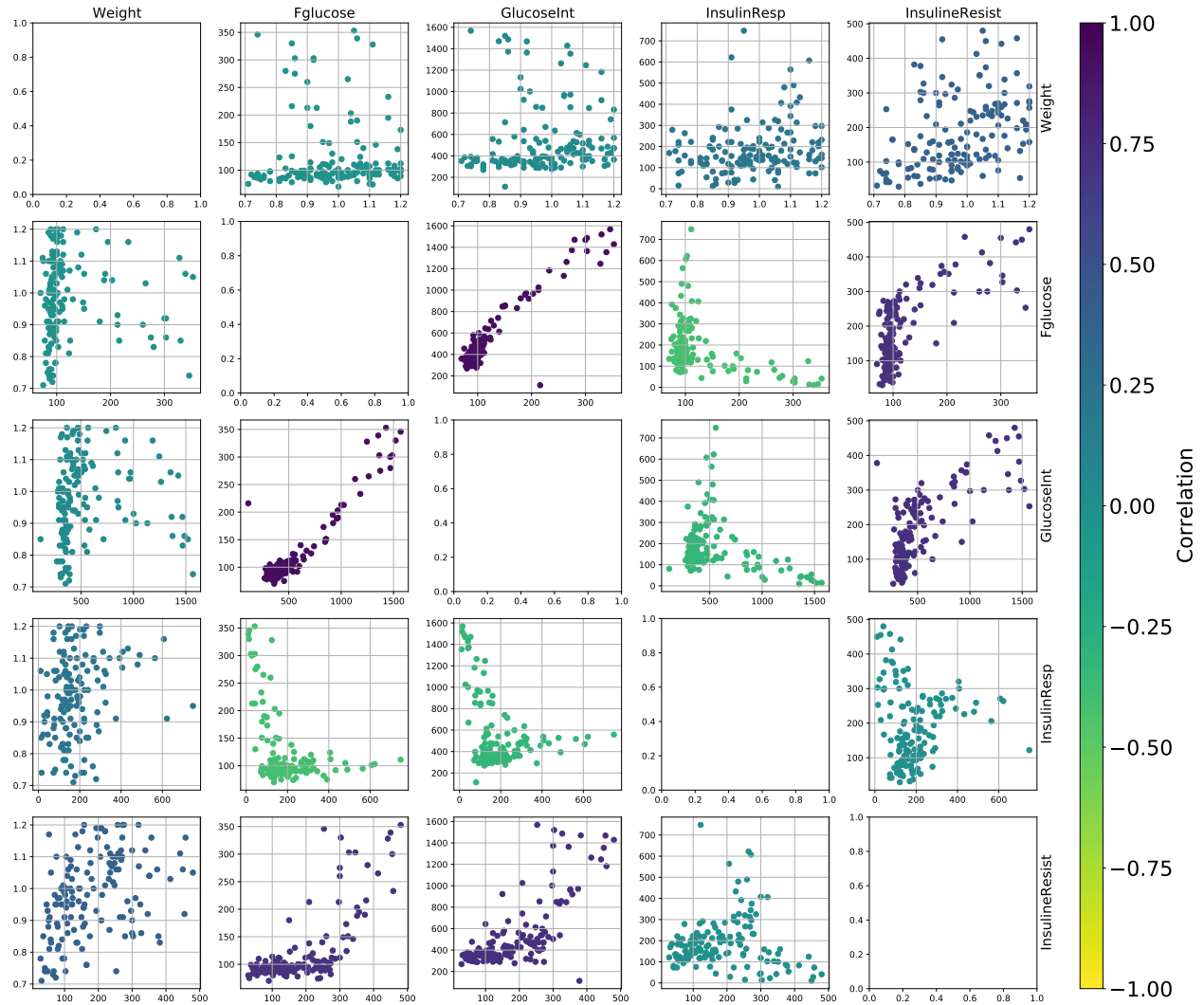


Figura 32: Gráfica de correlación para las variables numericas del problema 3

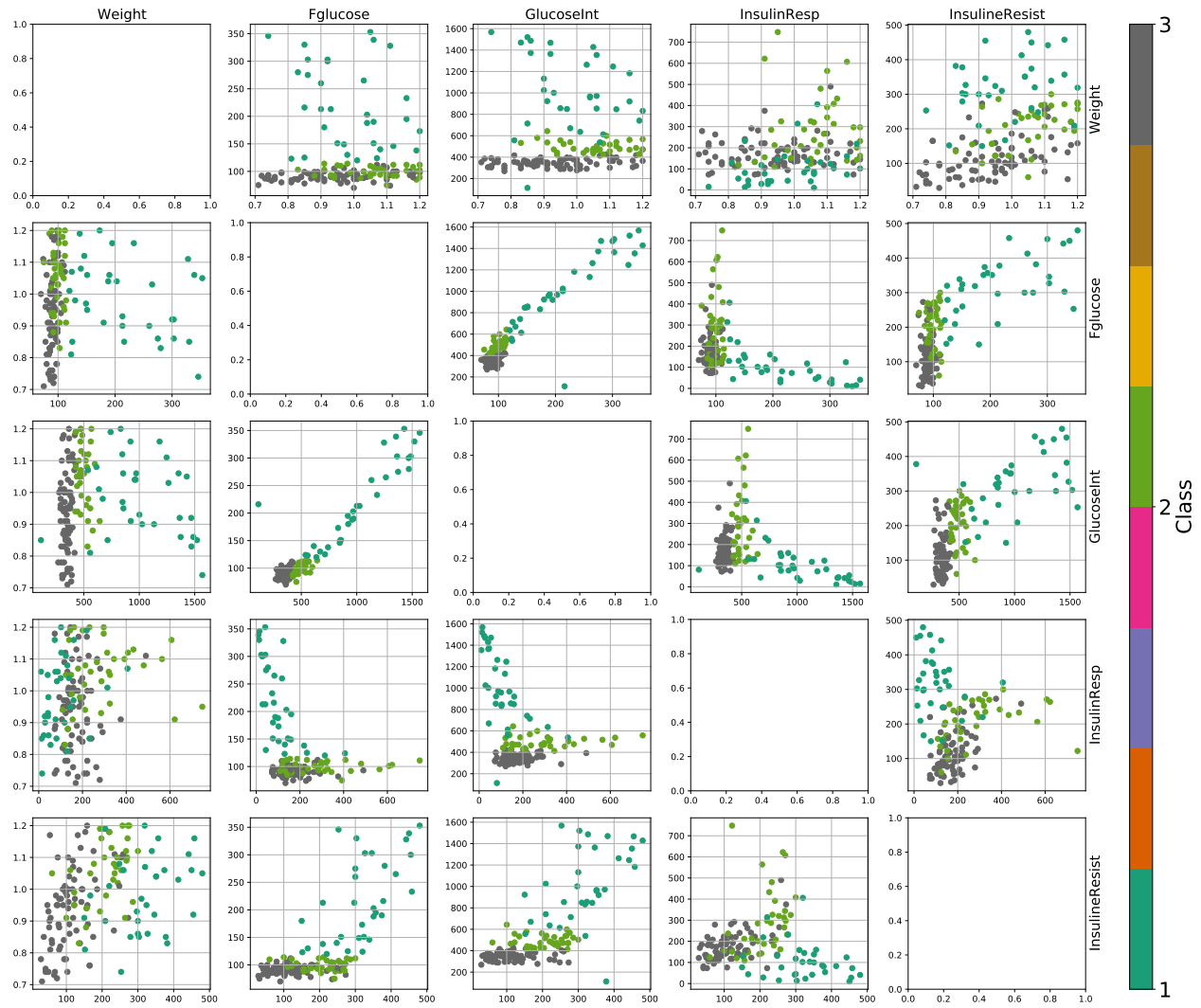


Figura 33: Gráficas de una variable contra otra divididas en clases para las variables numéricas del problema 3

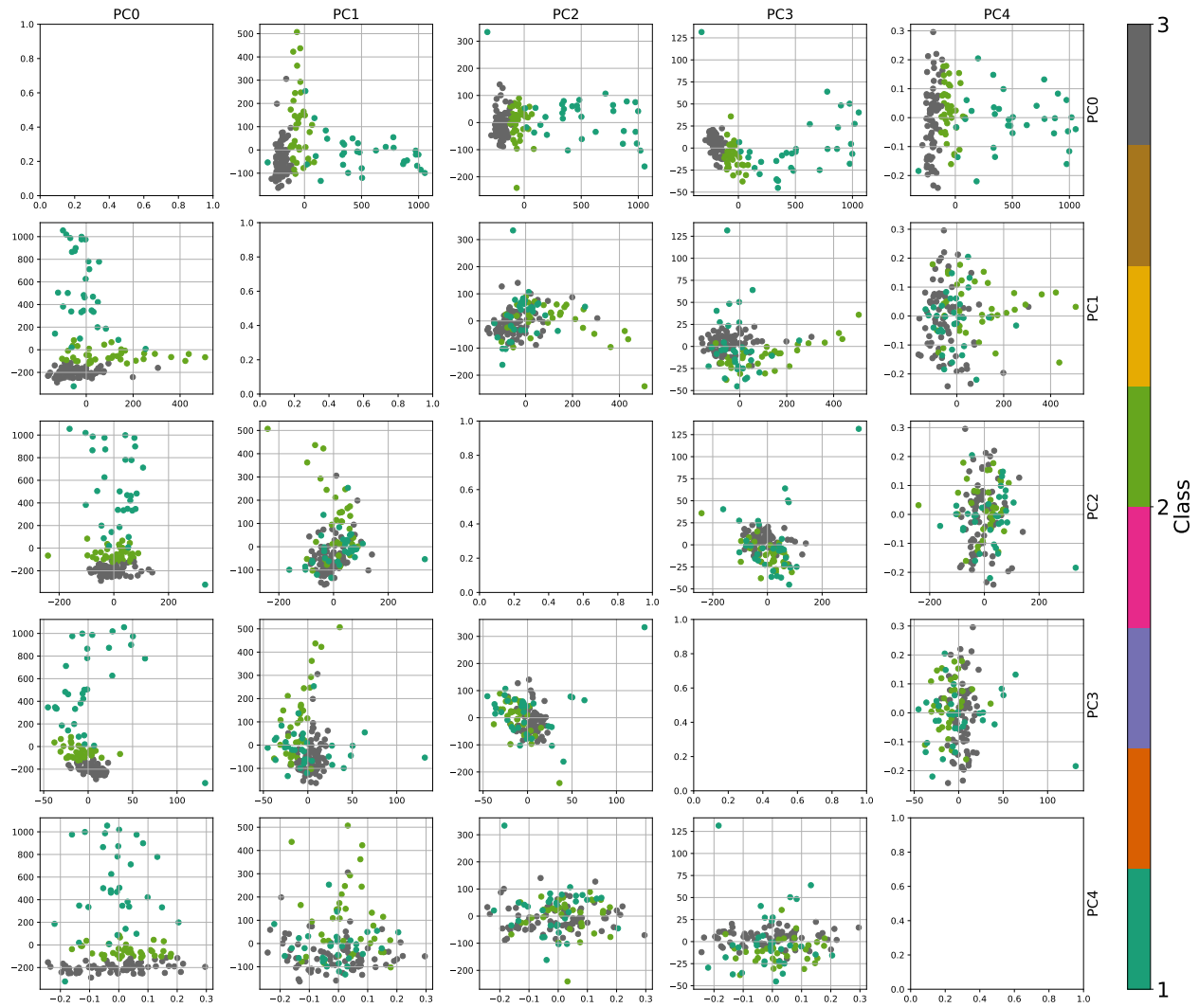


Figura 34: Gráficas de una variable contra otra divididas en clases para las variables numéricas del problema 3

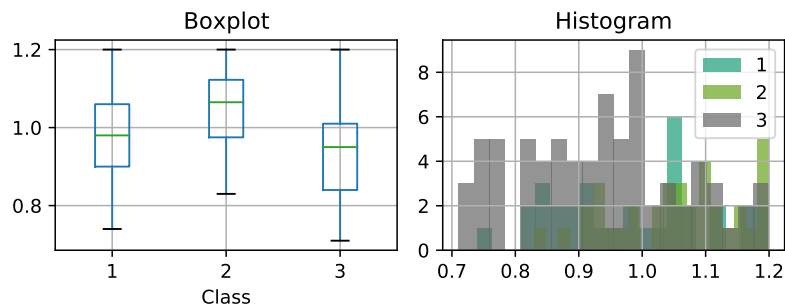


Figura 35: Distribución de la variable “Weight”

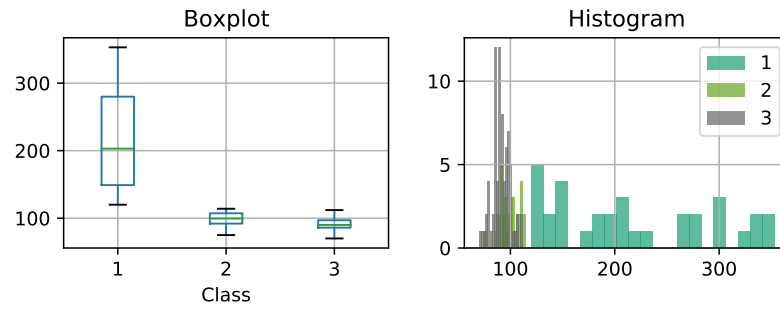


Figura 36: Distribución de la variable “Fglucose”

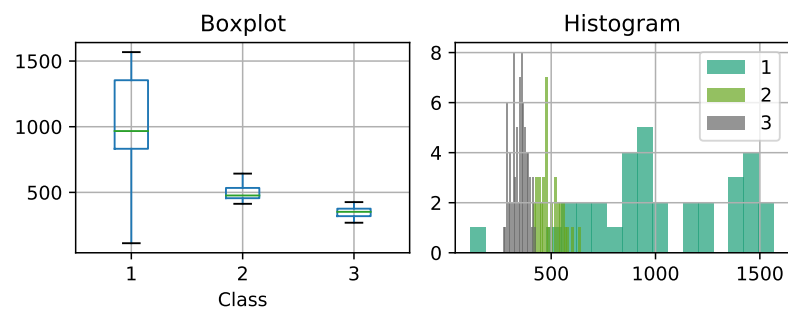


Figura 37: Distribución de la variable “Glucoseint”

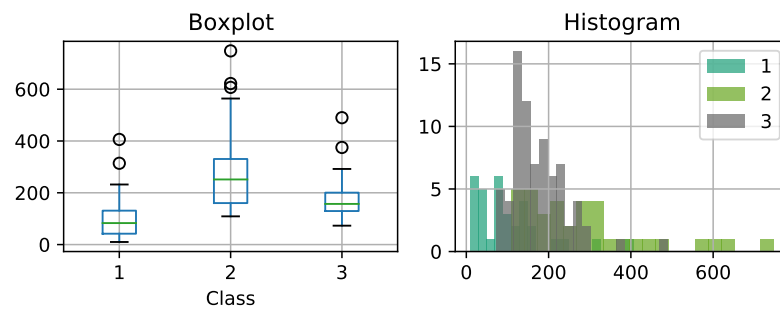


Figura 38: Distribución de la variable “InsulinResp”

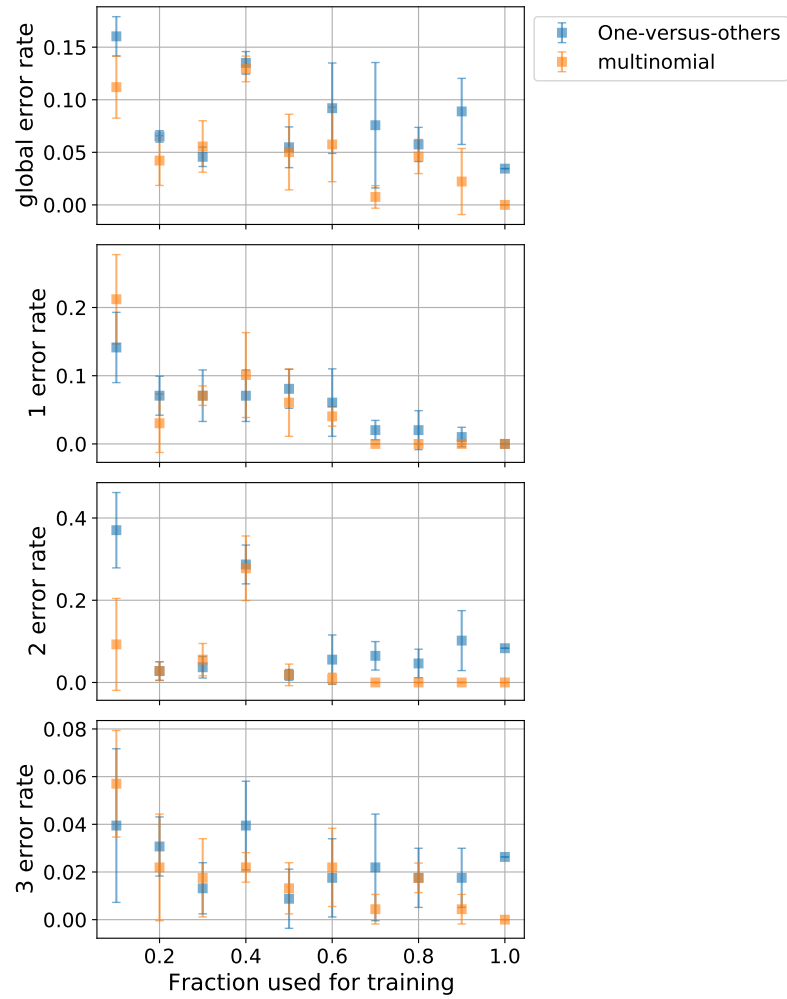


Figura 39: Distribución de la variable “InsulineResist”

Anexo 3: Código en Python de los problemas

```
1 #!/usr/bin/env python
2 # coding: utf-8
3 # Funciones auxiliares
4 import os
5 import numpy as np
6 import pandas as pd
7 import matplotlib.pyplot as plt
8 import matplotlib.colors as mcolors
9 import matplotlib.cm as cm
10 from scipy.stats import boxcox
11 import sklearn.linear_model as lm
12 import sklearn.discriminant_analysis as da
13 import sklearn.svm as svm
14 import sklearn.neighbors as ng
15 import sklearn.naive_bayes as nb
16 import sklearn.model_selection as ms
17 import sklearn.preprocessing as prep
18 import sklearn.decomposition as de
19 import statsmodels.api as sm
20 if not(os.path.isdir("tarea")):
21     mkdir("tarea")
22 def dataScaled(df):
23     # Creating dictionary to store the different data frames
24     data = {"original":df}
25     # Standarizing data to have mean 0 and variance 1
26     scaler = pr.StandardScaler()
27     scaler.fit(df)
28     data["standarized"] = pd.DataFrame(scaler.transform(df),index=df.index
29 ,columns=df.columns)
30     # Centering data to have variance 1
31     scaler = pr.StandardScaler(with_mean=False)
32     scaler.fit(df)
33     data["withmean"] = pd.DataFrame(scaler.transform(df),index=df.index,
34 columns=df.columns)
35     return data
36 def correlationPlots(df,saveName="",cmap=cm.viridis_r,f=16,width=50,**
37 kwargs):
38     variables = [var for var in df.columns if np.issubdtype(df[var].dtype,
39 np.number)]
40     nvar = len(variables)
41     values = [df[var].to_numpy() for var in variables]
42     correlations = np.corrcoef(values)
43     norm = mcolors.Normalize(vmin=-1,vmax=1)
44     fig , ax = plt.subplots(ncols=nvar,nrows=nvar,figsize=(19,16),
45 constrained_layout=True)
46     for i in range(nvar):
47         for j in range(nvar):
48             if i!=j:
49                 ax[i][j].scatter(values[i],values[j],color=cmap(norm(
50 correlations[i][j])),**kwargs)
51                 ax[i][j].grid()
52         ax[0][i].set_title(str(variables[i]),fontsize=f)
```

```

47         ax[i][nvar-1].yaxis.set_label_position("right")
48         ax[i][nvar-1].set_ylabel(str(variables[i]), fontsize=f)
49         cbar = fig.colorbar(cm.ScalarMappable(norm=norm, cmap=cmap), ax=ax.ravel(
50             ).tolist(), aspect=width)
51         cbar.ax.tick_params(labelsize=1.5*f)
52         cbar.ax.yaxis.set_label_position("right")
53         cbar.ax.set_ylabel("Correlation", fontsize=1.5*f)
54         if bool(saveName):
55             plt.savefig(os.path.join("tarea", "{0}-correlation.pdf".format(
56                 saveName)))
57             plt.show()
58 def classPlots(X, y, saveName="", cmap=cm.Dark2, f=16, width=50, **kwargs):
59     variables = [var for var in X.columns if np.issubdtype(X[var].dtype, np
60         .number)]
61     nvar = len(variables)
62     values = [X[var].to_numpy() for var in variables]
63     classes = list(set(y))
64     nclass = len(classes)
65     norm = mcolors.Normalize(vmin=0, vmax=nclass-1)
66     colors = [cmap(norm(classes.index(clas))) for clas in y]
67     fig, ax = plt.subplots(ncols=nvar, nrows=nvar, figsize=(19, 16),
68         constrained_layout=True)
69     for i in range(nvar):
70         for j in range(nvar):
71             if i != j:
72                 ax[i][j].scatter(values[i], values[j], color=colors, **kwargs
73             )
74             ax[i][j].grid()
75             ax[0][i].set_title(str(variables[i]), fontsize=f)
76             ax[i][nvar-1].yaxis.set_label_position("right")
77             ax[i][nvar-1].set_ylabel(str(variables[i]), fontsize=f)
78             cbar = fig.colorbar(cm.ScalarMappable(norm=norm, cmap=cmap), ax=ax.ravel(
79                 ).tolist(), aspect=width, ticks=range(nclass))
80             cbar.ax.tick_params(labelsize=1.5*f)
81             cbar.ax.yaxis.set_label_position("right")
82             cbar.ax.set_ylabel("Class", fontsize=1.5*f)
83             cbar.ax.set_yticklabels(classes)
84             if bool(saveName):
85                 plt.savefig(os.path.join("tarea", "{0}-classes.pdf".format(saveName
86                     )))
87             plt.show()
88 def pcaPlots(X, y, saveName="", cmap=cm.Dark2, f=16, width=50, ncomps=0, **kwargs
89     ):
90     variables = [var for var in X.columns if np.issubdtype(X[var].dtype, np
91         .number)]
92     nvar = len(variables)
93     values = [X[var].to_numpy() for var in variables]
94     aux = np.transpose(values)
95     pca = de.PCA()
96     pca.fit(aux)
97     Z = pca.transform(aux)
98     variables = ["PC{0}".format(i) for i in range(len(variables))]

```

```

92     values = np.transpose(Z)
93     classes = list(set(y))
94     nclass = len(classes)
95     norm = mcolors.Normalize(vmin=0,vmax=nclass-1)
96     colors = [cmap(norm(classes.index(clas))) for clas in y]
97     if ncomps==0:
98         ncomps=nvar
99     fig , ax = plt.subplots(ncols=ncomps,nrows=ncomps,figsize=(19,16),
100 constrained_layout=True)
101     for i in range(ncomps):
102         for j in range(ncomps):
103             if i!=j:
104                 ax[i][j].scatter(values[i],values[j],color=colors,**kwargs
105 )
106                 ax[i][j].grid()
107                 ax[0][i].set_title(str(variables[i]),fontsize=f)
108                 ax[i][nvar-1].axis.set_label_position("right")
109                 ax[i][nvar-1].set_ylabel(str(variables[i]),fontsize=f)
110     cbar = fig.colorbar(cm.ScalarMappable(norm=norm,cmap=cmap),ax=ax.ravel
111 ().tolist(),aspect=width,ticks=range(nclass))
112     cbar.ax.tick_params(labelsize=1.5*f)
113     cbar.ax.yaxis.set_label_position("right")
114     cbar.ax.set_ylabel("Class",fontsize=1.5*f)
115     cbar.ax.set_yticklabels(classes)
116     if bool(saveName):
117         plt.savefig(os.path.join("tarea","{0}-pca.pdf".format(saveName)))
118     plt.show()
119
120 def indVarDistPlots(df,resCol,saveName="",cmap=cm.Dark2,**kwargs):
121     if not(resCol in df.columns):
122         raise ValueError("Name {0} is not in columns".format(resCol))
123     classes = list(set(df[resCol]))
124     nclass = len(classes)
125     norm = mcolors.Normalize(vmin=0,vmax=nclass-1)
126     colors = [cmap(norm(i)) for i in range(nclass)]
127     for col in df.columns:
128         if col != resCol and np.issubdtype(df[col].dtype,np.number):
129             fig, ax = plt.subplots(ncols=2,nrows=1,figsize=(6,2.5))
130             df.boxplot(column=col,by=resCol,ax=ax[0])
131             for i,val in enumerate(classes):
132                 df[df[resCol]==val][col].hist(ax=ax[1],label=str(val),
133 color=colors[i],**kwargs)
134             ax[0].set_title("Boxplot")
135             ax[1].legend()
136             ax[1].set_title("Histogram")
137             fig.suptitle("{0} distribution".format(col),y=1.1)
138         elif col != resCol and (df[col].dtype == np.dtype('O')):
139             fig,ax = plt.subplots(nrows=1,ncols=1,figsize=(4,2.5))
140             labels = [str(val) for val in set(df[resCol])]
141             weights = [[1/df[df[resCol]==val][col].shape[0] for i in range
142 (df[df[resCol]==val][col].shape[0])] for val in set(df[resCol])]
143             weights = np.transpose(weights)
144             ax.hist(np.transpose([df[df[resCol]==val][col].to_numpy() for

```

```

141         val in set(df[resCol]))),
142             weights=weights,
143             label=labels,
144             color=colors,
145             **kwargs)
146     ax.legend()
147     ax.grid()
148     ax.set_ylabel("Normed counts")
149     ax.set_title("{0} histogram".format(col))
150     elif col != resCol:
151         print("Column {0} cannot be interpreted".format(col))
152     if bool(saveName):
153         plt.tight_layout()
154         plt.savefig(os.path.join("tarea", "{0}-{1}-dist.pdf".format(
155             saveName, col)))
156         plt.show()
157
158 def boxcoxLambdaTable(df, resCol, alpha=0.05):
159     names = []
160     lambdas = []
161     intervalsBot = []
162     intervalsTop = []
163     for col in df.columns:
164         if col != resCol and np.issubdtype(df[col].dtype, np.number):
165             if (df[col]>0).prod():
166                 names.append(col)
167                 bx = boxcox(df[col], alpha=alpha)
168                 lambdas.append(bx[1])
169                 intervalsBot.append(bx[2][0])
170                 intervalsTop.append(bx[2][1])
171             else:
172                 print("Can't convert column {0}: not entirely positive".
173                     format(col) )
174     fin = pd.DataFrame.from_dict({"lambda":lambdas,"Lower confidence
175     interval, alpha = {0}".format(alpha):intervalsBot,"Upper confidence
176     interval, alpha = {0}".format(alpha):intervalsTop})
177     fin.index = names
178     return fin.transpose()
179
180 def errorRates(mod,X,y,n=100,size=0.5,equalRatios=False):
181     fin = np.zeros(n)
182     classes = list(set(y))
183     perClass = [np.zeros(n) for val in classes]
184     for i in range(n):
185         if size<1.0:
186             if equalRatios:
187                 X_train, X_test, y_train, y_test = ms.train_test_split(X,y
188                 ,train_size=size,stratify=y)
189             else:
190                 X_train, X_test, y_train, y_test = ms.train_test_split(X,y
191                 ,train_size=size)
192             fit = mod.fit(X_train,y_train)
193             res = fit.predict(X_test)
194             fin[i] = np.mean(y_test != res)

```



```

188         for j, val in enumerate(classes):
189             curr = X[y==val]
190             res = fit.predict(curr)
191             perClass[j][i] = np.mean(y[y==val] != res)
192     elif size==1:
193         fit = mod.fit(X,y)
194         res = fit.predict(X)
195         fin[i] = np.mean(y != res)
196         for j, val in enumerate(classes):
197             curr = X[y==val]
198             res = fit.predict(curr)
199             perClass[j][i] = np.mean(y[y==val] != res)
200     final = [np.mean(fin), np.std(fin)]
201     for cla in perClass:
202         final.append(np.mean(cla))
203         final.append(np.std(cla))
204     return final
205
206 def errorRatesTable(models, X, y, n=100, size=0.5, equalRatios=False, names=[]):
207     if not bool(names):
208         names = [str(mod).split("(")[0] for mod in models]
209     elif len(names) != len(models):
210         raise ValueError("length of names must match length of models")
211     errors = [errorRates(mod, X, y, n=n, size=size, equalRatios=equalRatios)
212               for mod in models]
213     fin = pd.DataFrame(errors)
214     fin.index = names
215     cols = ["Mean global error", "std"]
216     classes = list(set(y))
217     for i in range(len(classes)):
218         cols.append("Mean {0} error".format(classes[i]))
219         cols.append("std")
220     fin.columns = cols
221     return fin.transpose()
222
223 def errorRatesSizeTable(models, X, y, sizes, n=100, equalRatios=False, names=[]):
224     :
225     tables = []
226     indexes = []
227     for i, size in enumerate(sizes):
228         print("Simulating size: {0}".format(size))
229         tab = errorRatesTable(models, X, y, n=n, size=size, equalRatios=
230 equalRatios, names=names)
231         tables.append(tab)
232         newIndex = [{"0", "entrenamiento: {1}"].format(name, np.round(sizes[i
233 ], decimals=2)) for name in tab.index]
234         indexes.extend(newIndex)
235     df = pd.concat(tables, ignore_index=True)
236     df.index = indexes
237     return df
238
239 def errorRatesSizesPlot(errorTable, saveName="", f1=16, f2=14, **kwargs):
240     sizes = [float(name.split(":")[-1]) for name in errorTable.index]
241     sizes = list(set(sizes))

```

```

238     nclass = int(errorTable.shape[0]/(2*len(sizes)))
239     names = [name.split(" ")[1] for name in errorTable.index]
240     names = [names[i] for i in range(0,2*nclass,2)]
241     fig,ax = plt.subplots(nrows=nclass,ncols=1,sharex=True,figsize=(8,2.5*
nclass))
242     for col in errorTable.columns:
243         for j in range(nclass):
244             vals = errorTable.iloc[range(2*j,errorTable.shape[0],2*nclass)
][col].to_numpy()
245             stds = errorTable.iloc[range(2*j+1,errorTable.shape[0],2*
nclass)][col].to_numpy()
246             ax[j].errorbar(sizes,vals,yerr=stds,label=str(col),fmt="o",**
kwargs)
247             ax[j].grid(True)
248             ax[j].set_ylabel("{0} error rate".format(names[j]),fontsize=f1
)
249             ax[j].tick_params(axis="both",labelsize=f2)
250     ax[nclass-1].set_xlabel("Fraction used for training",fontsize=f1)
251     ax[0].legend(bbox_to_anchor=(1,1),fontsize=f2)
252     plt.subplots_adjust(hspace=0.05)
253     if bool(saveName):
254         plt.tight_layout()
255         plt.savefig(os.path.join("tarea","{0}-sizeDependence.pdf".format(
saveName)))
256     plt.show()
257
258 def pairInteractionMin(df,resCol,model):
259     transformed = df.copy()
260     vals = []
261     inters = []
262     for col1 in df.columns:
263         for col2 in df.columns:
264             if col1 != resCol and col2!= resCol and col1!=col2:
265                 transformed["inter"] = transformed[col1]*transformed[col2]
266                 models = [model]
267                 names = ["LDA"]
268                 sizes=[1.0]
269                 res = errorRatesSizeTable(models,transformed.drop(resCol,
axis=1),transformed[resCol],sizes=sizes,n=10,equalRatios=True)
270                 vals.append(res.iloc[0,0])
271                 inters.append(col1 + "*" + col2)
272     m = np.argmin(vals)
273     print("error minimizing interacion: {0}".format(inters[m]))
274     print("Value: {0}".format(vals[m]))
275     return inters[m].split("*")
276
277 def multinomialLogisticRegressionClassifier(X,y,reference,**kwargs):
278     classes = list(set(y))
279     nclasses = len(classes)
280     classifiers = []
281     newY = y.astype("category").cat.codes
282     refCode = newY[y==reference].iloc[0]
283     newY = (newY - refCode) % nclasses
284     newY.name = y.name

```

```
285     conversion = {}
286     for c in classes:
287         code = newY[y==c].iloc[0]
288         conversion.update({c:code})
289     mod = sm.MNLogit(exog=X,endog=newY)
290     fit = mod.fit(maxiters=10**5)
291     return conversion,fit
292
293 # Problema 1
294 df = pd.read_csv("pimate.csv")
295 df = df.append(pd.read_csv("pimatr.csv"),ignore_index=True)
296 print(df.head())
297
298 correlationPlots(df,saveName="1",alpha=0.7,s=5)
299
300 classPlots(df.drop("type",axis=1),df["type"],saveName="1",s=5)
301
302 pcaPlots(df.drop("type",axis=1),df["type"],saveName="1",s=5)
303
304 indVarDistPlots(df,"type",saveName="1",alpha=0.7,density=True,bins=20)
305
306 models = [
307     da.LinearDiscriminantAnalysis(),
308     da.QuadraticDiscriminantAnalysis(),
309     nb.GaussianNB(),
310     lm.LogisticRegression(dual=False,max_iter=10**6),
311     ng.KNeighborsClassifier(),
312     svm.SVC()
313 ]
314 names = [
315     "LDA",
316     "QDA",
317     "Naive Bayes",
318     "Logistic",
319     "KNC, k=5",
320     "SVM"
321 ]
322 sizes=np.linspace(0.1,1.0,10)
323
324 res = errorRatesSizeTable(models,df.drop("type",axis=1),df["type"],sizes=
    sizes,n=10,names=names)
325 tab = "p{4cm}"
326 for col in res.columns:
327     tab += "|p{1.5cm}"
328 res.to_latex(buf=os.path.join("tarea","1-gen-tab.tex"),float_format="{:0.3
    f}".format,longtable=True,column_format=tab)
329 errorRatesSizesPlot(res,saveName="1-gen",alpha=0.6,capsize=3,marker="s")
330
331 sizes=np.linspace(0.1,1.0,10)
332 res = errorRatesSizeTable(models,df.drop("type",axis=1),df["type"],sizes=
    sizes,n=10,equalRatios=True,names=names)
333 tab = "p{4cm}"
334 for col in res.columns:
335     tab += "|p{1.5cm}"
```

```

336 res.to_latex(buf=os.path.join("tarea", "1-gen-eq-tab.tex"), float_format="
    {:.3f}".format, longtable=True, column_format=tab)
337 errorRatesSizesPlot(res, saveName="1-gen-eq", alpha=0.6, capsize=3, marker="s"
    )
338
339 transTable = boxcoxLambdaTable(df, "type", alpha=0.05)
340 tab = "p{4cm}"
341 for col in transTable.transpose().columns:
342     tab += "|p{3cm}"
343 transTable.transpose().to_latex(buf=os.path.join("tarea", "1-bxparam-tab.
    tex"), float_format="{:0.4f}".format, column_format=tab)
344
345 transformed = df.copy()
346 for col in transTable.columns:
347     transformed[col] = boxcox(transformed[col], lmbda = transTable[col]["
    lambda"])
348
349 sizes=np.linspace(0.1, 1.0, 10)
350 res = errorRatesSizeTable(models, transformed.drop("type", axis=1),
    transformed["type"], sizes=sizes, n=10, equalRatios=True)
351 errorRatesSizesPlot(res, saveName="1-boxcox-eq", alpha=0.6, capsize=3, marker="s")
352
353 # LDA analysis
354 mod = da.LinearDiscriminantAnalysis()
355 col1, col2 = pairInteractionMin(df, "type", mod)
356
357 transformed = df.copy()
358 transformed["pairIter"] = transformed[col1]*transformed[col2]
359 models = [mod]
360 names = ["LDA"]
361 sizes=np.linspace(0.1, 1.0, 10)
362 res1 = errorRatesSizeTable(models, transformed.drop("type", axis=1),
    transformed["type"], sizes=sizes, n=10, equalRatios=True, names = names)
363 tab = "p{4cm}"
364 for col in res1.columns:
365     tab += "|p{1.5cm}"
366 res1.to_latex(buf=os.path.join("tarea", "1-lda-eq-tab.tex"), float_format="
    {:.3f}".format, longtable=True, column_format=tab)
367 errorRatesSizesPlot(res1, saveName="1-lda-eq", alpha=0.6, capsize=3, marker="s"
    )
368
369 # Naive Bayes
370 mod = nb.GaussianNB()
371 col1, col2 = pairInteractionMin(df, "type", mod)
372
373 transformed = df.copy()
374 transformed["bmi*ped"] = transformed[col1]*transformed[col2]
375 models = [mod]
376 names = ["Naive Bayes"]
377 sizes=np.linspace(0.1, 1.0, 10)
378 res2 = errorRatesSizeTable(models, transformed.drop("type", axis=1),
    transformed["type"], sizes=sizes, n=10, equalRatios=True, names=names)
379 tab = "p{4cm}"

```

```
380 for col in res2.columns:
381     tab += "|p{1.5cm}"
382 res2.to_latex(buf=os.path.join("tarea", "1-nb-eq-tab.tex"), float_format="{:0.3f}".format, longtable=True, column_format=tab)
383 errorRatesSizesPlot(res2, saveName="1-nb-eq", alpha=0.6, capsize=3, marker="s"
384 )
385 # Logistic regression
386 transformed = df.copy()
387 transformed = transformed[["glu", "bmi", "ped", 'age', "type"]]
388 transformed["age2"] = transformed["age"]*transformed["age"]
389 models = [lm.LogisticRegression(max_iter=10**6)]
390 names = ["Logistic"]
391 sizes=np.linspace(0.1,1.0,10)
392 res3 = errorRatesSizeTable(models, transformed.drop("type", axis=1),
393     transformed["type"], sizes=sizes, n=10, equalRatios=True, names=names)
394 tab = "p{4cm}"
395 for col in res3.columns:
396     tab += "|p{1.5cm}"
397 res3.to_latex(buf=os.path.join("tarea", "1-log-eq-tab.tex"), float_format="{:0.3f}".format, longtable=True, column_format=tab)
398 errorRatesSizesPlot(res3, saveName="1-log-eq", alpha=0.6, capsize=3, marker="s"
399 )
400 # SVM analysis
401 models = [
402     svm.SVC(kernel="linear"),
403     svm.SVC(kernel="rbf", degree=1),
404     svm.SVC(kernel="rbf", degree=3),
405     svm.SVC(kernel="rbf", degree=5),
406     svm.SVC(kernel="poly", degree=1),
407     svm.SVC(kernel="poly", degree=3),
408     svm.SVC(kernel="poly", degree=5),
409 ]
410 names = [
411     "SVM, kernel = linear",
412     "SVM, kernel = rbf, deg = 1",
413     "SVM, kernel = rbf, deg = 3",
414     "SVM, kernel = rbf, deg = 5",
415     "SVM, kernel = poly, deg = 1",
416     "SVM, kernel = poly, deg = 3",
417     "SVM, kernel = poly, deg = 5"
418 ]
419 res = errorRatesSizeTable(models, df.drop("type", axis=1), df["type"], sizes=
420     sizes, n=10, names=names)
421 errorRatesSizesPlot(res, saveName="1-svm", alpha=0.6, capsize=3, marker="s")
422 res = errorRatesSizeTable(models, df.drop("type", axis=1), df["type"], sizes=
423     sizes, n=10, names=names, equalRatios=True)
424 errorRatesSizesPlot(res, saveName="1-svm-eq", alpha=0.6, capsize=3, marker="s"
425 )
426 models = [
```

```
426     svm.SVC(kernel="poly",degree=5)
427 ]
428 names = [ "SVM, kernel = poly, deg = 5"
429 ]
430 sizes=np.linspace(0.1,1,10)
431 res4 = errorRatesSizeTable(models,df.drop("type",axis=1),df["type"],sizes=
    sizes,n=10,names=names,equalRatios=True)
432 tab = "p{4cm}"
433 for col in res4.columns:
434     tab += "|p{1.5cm}"
435 res4.to_latex(buf=os.path.join("tarea","1-svm-fin-eq-tab.tex"),
    float_format="{:0.3f}".format,longtable=True,column_format=tab)
436
437 errorRatesSizesPlot(res4,saveName="1-svm-fin-eq",alpha=0.8,capsize=3,
    marker="s",markersize=10)
438
439 final = pd.concat([res1,res2,res3,res4],axis=1)
440 tab = "p{4cm}"
441 for col in final.columns:
442     tab += "|p{1.5cm}"
443 final.to_latex(buf=os.path.join("tarea","1-final-tab.tex"),float_format="
    {:0.3f}".format,longtable=True,column_format=tab)
444 errorRatesSizesPlot(final,saveName="1-final",alpha=0.65,capsize=3,marker="
    s",markersize=8)
445
446 neds = [res.iloc[range(0,res.shape[0],2)[-3:],:] for res in [res1,res2,
    res3,res4]]
447 final = pd.concat(neds,axis=1)
448 tab = "p{4cm}"
449 for col in final.columns:
450     tab += "|p{1.5cm}"
451 final.to_latex(buf=os.path.join("tarea","1-final-correg-aparent-tab.tex"),
    float_format="{:0.3f}".format,longtable=True,column_format=tab)
452
453 neds = [res.iloc[range(0,res.shape[0],2)[-12:-9],:] for res in [res1,res2,
    res3,res4]]
454 final = pd.concat(neds,axis=1)
455 tab = "p{4cm}"
456 for col in final.columns:
457     tab += "|p{1.5cm}"
458 final.to_latex(buf=os.path.join("tarea","1-final-correg-nonaparent-tab.tex
    "),float_format="{:0.3f}".format,longtable=True,column_format=tab)
459
460 # Problema 2
461 df = pd.read_csv("cad1.csv",index_col=0)
462 print(df.head())
463
464 indVarDistPlots(df,"CAD",saveName="2")
465
466 resCol="CAD"
467 dfCoded = df.copy()
468 for col in df.columns:
469     if col!=resCol and df[col].dtype==np.dtype("O"):
470         dfCoded[col] = df[col].astype("category").cat.codes
```

```
471
472 dfCoded
473
474 models = [
475     da.LinearDiscriminantAnalysis(),
476     da.QuadraticDiscriminantAnalysis(),
477     nb.GaussianNB(),
478     lm.LogisticRegression(dual=False,max_iter=10**6),
479     ng.KNeighborsClassifier(),
480     svm.SVC()
481 ]
482 names = [
483     "LDA",
484     "QDA",
485     "Naive Bayes",
486     "Logistic",
487     "KNC, k=5",
488     "SVM"
489 ]
490
491 sizes=np.linspace(0.1,1.0,10)
492 res = errorRatesSizeTable(models,dfCoded.drop("CAD",axis=1),dfCoded["CAD"
493 ],sizes=sizes,n=10,equalRatios=True,names=names)
494 tab = "p{4cm}"
495 for col in res.columns:
496     tab += "|p{1.5cm}"
497 res.to_latex(buf=os.path.join("tarea","2-gen-eq-tab.tex"),float_format="
498     {:0.3f}".format,longtable=True,column_format=tab)
499 errorRatesSizesPlot(res,saveName="2-gen-eq",alpha=0.6,capsize=3,marker="s"
500 )
501
502 # Naive Bayes
503 mod = nb.GaussianNB()
504 col1,col2 = pairInteractionMin(dfCoded,"CAD",mod)
505
506 transformed = dfCoded.copy()
507 transformed["inter"] = transformed[col2]*transformed[col1]
508 models = [mod]
509 names = ["Naive Bayes"]
510 sizes=np.linspace(0.1,1.0,10)
511 res2 = errorRatesSizeTable(models,transformed.drop("CAD",axis=1),
512     transformed["CAD"],sizes=sizes,n=10,equalRatios=True,names=names)
513 tab = "p{4cm}"
514 for col in res2.columns:
515     tab += "|p{1.5cm}"
516 res2.to_latex(buf=os.path.join("tarea","2-nb-eq-tab.tex"),float_format="
517     {:0.3f}".format,longtable=True,column_format=tab)
518 errorRatesSizesPlot(res2,saveName="2-nb-eq",alpha=0.6,capsize=3,marker="s"
519 )
520
521 # logistic Regresion
522 transformed = dfCoded.copy()
523 transformed = transformed[["AngPec", "AMI", "STcode", "STchange", "Hyperchol
524     ", "CAD"]]
```

```
518 models = [lm.LogisticRegression(max_iter=10**6)]
519 names = ["Logistic"]
520 sizes=np.linspace(0.1,1.0,10)
521 res3 = errorRatesSizeTable(models,transformed.drop("CAD",axis=1),
    transformed["CAD"],sizes=sizes,n=10,equalRatios=True,names=names)
522 tab = "p{4cm}"
523 for col in res3.columns:
524     tab += "|p{1.5cm}"
525 res3.to_latex(buf=os.path.join("tarea","2-log-eq-tab.tex"),float_format="
    {:0.3f}".format,longtable=True,column_format=tab)
526 errorRatesSizesPlot(res3,saveName="2-log-eq",alpha=0.6,capsize=3,marker="s
    ")
527
528 # SVM analysis
529 models = [
530     svm.SVC(kernel="linear"),
531     svm.SVC(kernel="rbf",degree=1),
532     svm.SVC(kernel="rbf",degree=3),
533     svm.SVC(kernel="rbf",degree=5),
534     svm.SVC(kernel="poly",degree=1),
535     svm.SVC(kernel="poly",degree=3),
536     svm.SVC(kernel="poly",degree=5),
537 ]
538 names = [
539     "SVM, kernel = linear",
540     "SVM, kernel = rbf, deg = 1",
541     "SVM, kernel = rbf, deg = 3",
542     "SVM, kernel = rbf, deg = 5",
543     "SVM, kernel = poly, deg = 1",
544     "SVM, kernel = poly, deg = 3",
545     "SVM, kernel = poly, deg = 5"
546 ]
547
548 res = errorRatesSizeTable(models,dfCoded.drop("CAD",axis=1),dfCoded["CAD"
    ],sizes=sizes,n=10,names=names,equalRatios=True)
549 errorRatesSizesPlot(res,saveName="2-svm-eq",alpha=0.6,capsize=3,marker="s"
    )
550
551 models = [
552     svm.SVC(kernel="poly",degree=5)
553 ]
554 names = [ "SVM, kernel = poly, deg = 5"
555 ]
556 sizes=np.linspace(0.1,1,10)
557 res4 = errorRatesSizeTable(models,dfCoded.drop("CAD",axis=1),dfCoded["CAD"
    ],sizes=sizes,n=10,equalRatios=True,names=names)
558 tab = "p{4cm}"
559 for col in res4.columns:
560     tab += "|p{1.5cm}"
561 res4.to_latex(buf=os.path.join("tarea","2-svm-fin-eq-tab.tex"),
    float_format="{:0.3f}".format,longtable=True,column_format=tab)
562 errorRatesSizesPlot(res4,saveName="2-svm-fin-eq",alpha=0.8,capsize=3,
    marker="s",markersize=10)
563
```



```
564 final = pd.concat([res2,res3,res4],axis=1)
565 tab = "p{4cm}"
566 for col in final.columns:
567     tab += "|p{1.5cm}"
568 final.to_latex(buf=os.path.join("tarea","2-final-tab.tex"),float_format="
    {:0.3f}".format,longtable=True,column_format=tab)
569 errorRatesSizesPlot(final,saveName="2-final",alpha=0.65,capsize=3,marker="
    s",markersize=8)
570
571 neds = [res.iloc[range(0,res.shape[0],2)[-3:],:] for res in [res2,res3,
    res4]]
572 final = pd.concat(neds,axis=1)
573 tab = "p{4cm}"
574 for col in final.columns:
575     tab += "|p{1.5cm}"
576 final.to_latex(buf=os.path.join("tarea","2-final-correg-aparent-tab.tex"),
    float_format="{:0.3f}".format,longtable=True,column_format=tab)
577
578 neds = [res.iloc[range(0,res.shape[0],2)[-12:-9],:] for res in [res2,res3,
    res4]]
579 final = pd.concat(neds,axis=1)
580 tab = "p{4cm}"
581 for col in final.columns:
582     tab += "|p{1.5cm}"
583 final.to_latex(buf=os.path.join("tarea","2-final-correg-nonaparent-tab.tex
    "),float_format="{:0.3f}".format,longtable=True,column_format=tab)
584
585 # Problema 3
586 df = pd.read_csv("Glucose1.txt",index_col="Patient")
587 df["Class"] = df["Class"].astype("0")
588 print(df.head())
589
590 indVarDistPlots(df,"Class",saveName="3",alpha=0.7,bins = 20)
591
592 correlationPlots(df,saveName="3")
593
594 pcaPlots(df.drop("Class",axis=1),df["Class"],saveName="3")
595
596 classPlots(df.drop("Class",axis=1),df["Class"],saveName="3")
597
598 models = [lm.LogisticRegression(max_iter=10**6,multi_class="ovr"),lm.
    LogisticRegression(max_iter=10**6,multi_class="multinomial")]
599 names = ["One-versus-others","multinomial"]
600 sizes=np.linspace(0.1,1.0,10)
601
602 res = errorRatesSizeTable(models,df.drop("Class",axis=1),df["Class"].
    astype("int64"),sizes=sizes,n=3,equalRatios=True,names=names)
603 tab = "p{4cm}"
604 for col in res.columns:
605     tab += "|p{1.5cm}"
606 res.to_latex(buf=os.path.join("tarea","3-eq-tab.tex"),float_format="{:0.3f
    }".format,longtable=True,column_format=tab)
607 errorRatesSizesPlot(res,saveName="3-eq",alpha=0.6,capsize=3,marker="s")
608
```

```

609 moddescript = df[["Fglucose", "GlucoseInt", "Class"]].copy()
610 moddescript["Fglucose^2"] = df["Fglucose"]**2
611
612 models = [lm.LogisticRegression(max_iter=10**6, multi_class="ovr")]
613 names = ["One-versus-others"]
614 sizes=np.linspace(0.1,1.0,10)
615
616 modpredict1 = df[["Fglucose", "Class"]].copy()
617 interactions = ["GlucoseInt*GlucoseInt"]
618 for inter in interactions:
619     columns = inter.split("*")
620     modpredict1[inter] = df[columns].product(axis=1)
621
622 res1 = errorRatesSizeTable(models, modpredict1.drop("Class", axis=1),
623     modpredict1["Class"].astype("int64"), sizes=sizes, n=20, equalRatios=True,
624     names=names)
625
626 modpredict2 = df[["InsulinResp", "Class"]].copy()
627 interactions = ["Fglucose*InsulinResp", "GlucoseInt*InsulinResp"]
628 for inter in interactions:
629     columns = inter.split("*")
630     modpredict2[inter] = df[columns].product(axis=1)
631
632 res2 = errorRatesSizeTable(models, modpredict2.drop("Class", axis=1),
633     modpredict2["Class"].astype("int64"), sizes=sizes, n=20, equalRatios=True,
634     names=names)
635
636 modpredict3 = df[["Fglucose", "Class"]].copy()
637 interactions = ["Weight*InsulineResist", "Fglucose*GlucoseInt", "
638     GlucoseInt*GlucoseInt"]
639 for inter in interactions:
640     columns = inter.split("*")
641     modpredict[inter] = df[columns].product(axis=1)
642
643 res3 = errorRatesSizeTable(models, modpredict3.drop("Class", axis=1),
644     modpredict3["Class"].astype("int64"), sizes=sizes, n=20, equalRatios=True,
645     names=names)
646
647 neds = [res.iloc[range(0, res.shape[0], 2)[-4:], 0] for res in [res1, res2,
648     res3]]
649 final = pd.concat(neds, axis=1)
650 final.columns = ["Modelo {0}".format(i) for i in range(1, 4)]
651 tab = "p{4cm}"
652 for col in final.columns:
653     tab += "|p{2cm}"
654 final.to_latex(buf=os.path.join("tarea", "3-descrip-aparent-tab.tex"),
655     float_format="{:0.3f}".format, longtable=True, column_format=tab)
656 final
657
658 neds = [res.iloc[range(0, res.shape[0], 2)[-16:-12], 0] for res in [res1, res2,
659     res3]]
660 final = pd.concat(neds, axis=1)
661 final.columns = ["Modelo {0}".format(i) for i in range(1, 4)]
662 tab = "p{4cm}"

```

```
653 for col in final.columns:
654     tab += "|p{2cm}"
655 final.to_latex(buf=os.path.join("tarea", "3-descrip-nonaparent.tex"),
656               float_format="{:0.3f}".format, longtable=True, column_format=tab)
657 final
658 modpredict = modpredict2
659
660 res1 = errorRatesSizeTable(models, moddescript.drop("Class", axis=1),
661                             moddescript["Class"].astype("int64"), sizes=sizes, n=3, equalRatios=True,
662                             names=names)
661 tab = "p{4cm}"
662 for col in res.columns:
663     tab += "|p{1.5cm}"
664 res1.to_latex(buf=os.path.join("tarea", "3-descrip-tab.tex"), float_format="
665               {:0.3f}".format, longtable=True, column_format=tab)
666 errorRatesSizesPlot(res1, saveName="3-descrip", alpha=0.6, capsize=3, marker="
667               s")
668
669 res2 = errorRatesSizeTable(models, modpredict.drop("Class", axis=1),
670                             modpredict["Class"].astype("int64"), sizes=sizes, n=20, equalRatios=True,
671                             names=names)
672 tab = "p{4cm}"
673 for col in res.columns:
674     tab += "|p{1.5cm}"
675 res2.to_latex(buf=os.path.join("tarea", "3-predict-tab.tex"), float_format="
676               {:0.3f}".format, longtable=True, column_format=tab)
677 errorRatesSizesPlot(res2, saveName="3-predict", alpha=0.6, capsize=3, marker="
678               s")
679
680 neds = [res.iloc[range(0, res.shape[0], 2)[-4:], 0] for res in [res1, res2]]
681 final = pd.concat(neds, axis=1)
682 final.columns = ["Descriptivo", "Predictivo"]
683 tab = "p{4cm}"
684 for col in final.columns:
685     tab += "|p{1.5cm}"
686 final.to_latex(buf=os.path.join("tarea", "3-final-correg-aparent-tab.tex"),
687               float_format="{:0.3f}".format, longtable=True, column_format=tab)
688 final
689
690 neds = [res.iloc[range(0, res.shape[0], 2)[-16:-12], 0] for res in [res1, res2
691                               ]]
692 final = pd.concat(neds, axis=1)
693 final.columns = ["Descriptivo", "Predictivo"]
694 tab = "p{4cm}"
695 for col in final.columns:
696     tab += "|p{2cm}"
697 final.to_latex(buf=os.path.join("tarea", "3-final-correg-nonaparent-tab.tex
698               "), float_format="{:0.3f}".format, longtable=True, column_format=tab)
699 final
```

Referencias

- [1] Jack W Smith y col. “Using the ADAP learning algorithm to forecast the onset of diabetes mellitus”. En: *Proceedings of the Annual Symposium on Computer Application in Medical Care*. American Medical Informatics Association. 1988, pág. 261.
- [2] *C-Support Vector Classification. sklearn.svm.SVC documentation*. 2020. URL: <https://tex.stackexchange.com/questions/3587/how-can-i-use-bibtex-to-cite-a-web-page> (visitado 05-04-2020).
- [3] JF Hansen. “The clinical diagnosis of ischaemic heart disease due to coronary artery disease.” En: *Danish medical bulletin* 27.6 (1980), págs. 280-286.