

# Práctica 4: Máquinas de vectores soporte

Convocatoria de enero (curso académico 2022/2023)

Asignatura: Introducción a los modelos computacionales  
4º Grado Ingeniería Informática (Universidad de Córdoba)

30 de noviembre de 2022

## Resumen

Esta práctica sirve para familiarizar al alumno con las máquinas de vectores soporte (*Support Vector Machines*, SVMs). De esta forma, utilizaremos las SVMs para varios problemas de clasificación, con objeto de que el alumno entienda mejor su funcionamiento y, sobre todo, el efecto de sus parámetros. Para ello, emplearemos su implementación en `scikit-learn` a través de la librería `libsvm`, una de las implementaciones más eficientes disponibles hoy en día, además de ser gratuita. La entrega se hará utilizando la tarea en Moodle habilitada al efecto. Se deberá subir un único fichero comprimido incluyendo los entregables que se indican en este guión. El día tope para la entrega es el **20 de diciembre de 2022**. En caso de que dos alumnos entreguen prácticas copiadas, no se puntuarán ninguna de las dos.

## 1. Introducción

En esta práctica, vamos a realizar distintos ejercicios para entender mejor el funcionamiento de las máquinas de vectores soporte (*Support Vector Machines*, SVMs) en problemas de clasificación. Las SVMs han sido durante muchos años el estado del arte en problemas de reconocimiento de patrones y es importante conocer su funcionamiento y como responden a sus distintos parámetros. La práctica se realizará utilizando la librería `libsvm` a través de `scikit-learn`, que es una de las implementaciones más populares y eficientes que se disponen actualmente. En todos los casos, consideraremos la versión "C-SVC", incluida en `libsvm`.

La sección 2 describe una serie de experimentos a realizar utilizando bases de datos sintéticas, que nos ayudarán a entender mejor el funcionamiento de las SVMs y de sus parámetros. La sección 3 describe otros experimentos a realizar sobre bases de datos reales. Finalmente, la sección 4 especifica los ficheros a entregar para esta práctica.

Se deben contestar todas las preguntas que vengan en recuadros de este estilo.

## 2. Bases de datos sintéticas

En esta primera parte de la práctica, utilizaremos una serie de experimentos sobre bases de datos sintéticas, que nos ayudarán a entender cómo funcionan las SVMs y qué efectos producen sobre los resultados los dos parámetros que hay que especificar para entrenarlas.

### 2.1. Instalación de `libsvm`

Para esta práctica vamos a utilizar Python. `scikit-learn` ya trae un clasificador que aplica el algoritmo SVC, utilizando la implementación de `libsvm`. En concreto, el clasificador es `sklearn.svm.SVC`.

## 2.2. Representación 2D de las SVMs

Para realizar la representación gráfica en dos dimensiones desde Python, vamos a utilizar el *script* `libsvm.py`.

**Pregunta [1]:** Abre este *script* y explica su contenido. Podrás ver que se utiliza el primer *dataset* de ejemplo y se realiza la representación gráfica del SVM. Comenta que tipo de *kernel* se está utilizando y cuáles son los parámetros de entrenamiento.

## 2.3. Primer dataset de ejemplo

Ejecuta el *script* mencionado anteriormente pero comenta las líneas correspondientes a la representación gráfica de la SVM, para ver solo los puntos del *dataset* (representación similar a la de Figura 1, `dataset1.csv`).

**Pregunta [2]:** Intuitivamente, ¿qué hiperplano crees que incurrirá en un menor error de *test* en la tarea de separar las dos clases de puntos?

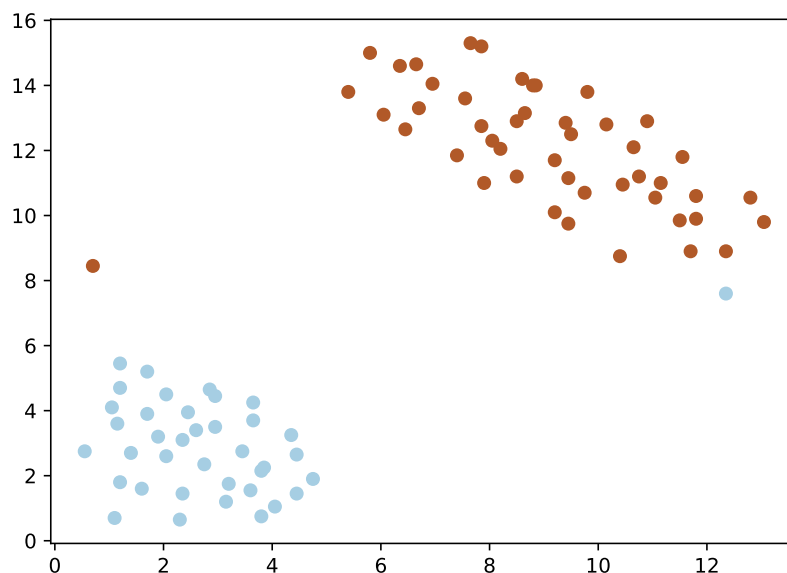


Figura 1: Dataset sintético o tipo *toy* #1.

En esta parte del ejercicio, vamos a utilizar una SVM lineal, es decir, una SVM donde la separación de las dos clases será una línea recta. La forma de especificar los parámetros del algoritmo SVC en Python puede consultarse en <http://scikit-learn.org/stable/modules/generated/sklearn.svm.SVC.html>.

Vamos a probar a utilizar distintos valores para el parámetro  $C$ . De manera informal, podemos decir que el parámetro  $C$  es un parámetro que define cuánto vamos a penalizar el hecho de que se cometan errores en la clasificación. Un valor de  $C$  muy grande resultará en una SVM que cometerá el mínimo número posible de errores, mientras que un valor pequeño de  $C$  favorecerá el que se obtenga un clasificador con el máximo margen, aunque se cometan errores. Por supuesto, esto está relacionado con el posible sobre-entrenamiento (aprendizaje de patrones ruidosos de entrenamiento que puedan despistar al modelo).

**Pregunta [3]:** Modifica el *script* probando varios valores de  $C$ , en concreto,  $C \in \{10^{-2}, 10^{-1}, 10^0, 10^1, 10^2, 10^3, 10^4\}$ . Observa qué sucede, explica porqué y escoge el valor más adecuado.

## 2.4. Segundo dataset de ejemplo

En este caso, trabajaremos con el *dataset* de la Figura 2 (`dataset2.csv`).

**Pregunta [4]:** Prueba a lanzar una SVM lineal con los valores para  $C$  que se utilizaron en la pregunta anterior. ¿Consigues algún resultado satisfactorio en el sentido de que no haya errores en el conjunto de entrenamiento?. ¿Por qué?

Para conseguir hiperplanos de separación no lineales, las SVMs utilizan el llamado *truco del kernel* que, explicado de manera informal, construye un modelo lineal en un espacio de muchas dimensiones y lo proyecta de nuevo al espacio original, resultando en un modelo no lineal. Esto se consigue mediante el uso de funciones de *kernel* aplicadas sobre cada uno de los patrones de entrenamiento. El *kernel* más común es el *kernel* RBF, también llamado Gaussiano. Para poder aplicarlo, el usuario tiene que especificar un parámetro adicional ( $\gamma$ , gamma, en `libsvm`  $\gamma = 1/\text{radio}$ ). Un radio alto tiende a soluciones más suaves, con menor sobre-entrenamiento, mientras que un radio de *kernel* pequeño tiende a producir mayor sobre-entrenamiento.

**Pregunta [5]:** Propón una configuración de SVM no lineal (utilizando el *kernel* tipo RBF o Gaussiano) que resuelva el problema. El resultado debería ser similar al de la Figura 3. Muévete en potencias de 2 utilizando un rango del tipo  $\{2^{-2}, 2^{-1}, 2^0, 2^1, 2^2\}$ . ¿Qué valores has considerado para  $C$  y para  $\gamma$ ? Además, incluye un ejemplo de una configuración de parámetros que produzca sobre-entrenamiento y otra que produzca infra-entrenamiento.

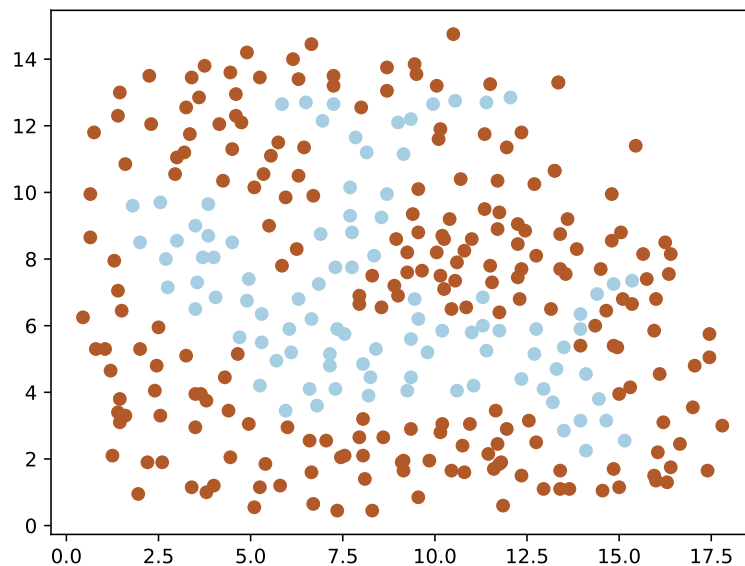


Figura 2: Dataset sintético o tipo *toy* #2.

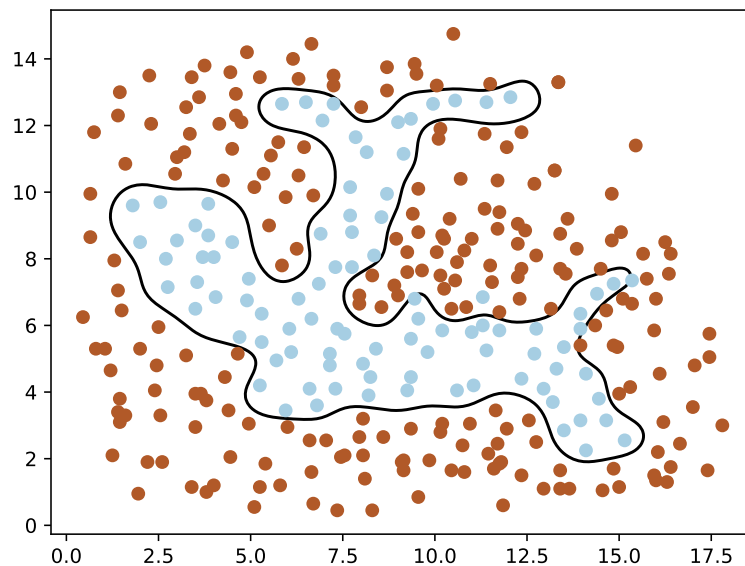


Figura 3: Dataset sintético o tipo *toy* #2 bien clasificado.

## 2.5. Tercer dataset de ejemplo

Por último, vamos a trabajar con el *dataset* de la Figura 4 (`dataset3.csv`).

**Pregunta [6]:** En este caso, ¿es el *dataset* linealmente separable?. A primera vista, ¿detectas puntos que presumiblemente sean *outliers*?, ¿por qué?.

**Pregunta [7]:** Lanza una SVM para clasificar los datos, con objeto de obtener un resultado lo más parecido al de la Figura 5. Muévete en potencias de 2 utilizando un rango del tipo  $\{2^{-2}, 2^{-1}, 2^0, 2^1, 2^2\}$ . Establece el valor de los parámetros óptimos. Además, incluye un ejemplo de una configuración de parámetros que produzca sobre-entrenamiento y otra que produzca infra-entrenamiento.

## 2.6. Interfaz de consola

Aunque estamos utilizando la interfaz de Python, vamos a especificar algunas características de la interfaz de consola de `libsvm`<sup>1</sup>, para conocer mejor el proceso de entrenamiento de una SVM. La interfaz de consola de `libsvm` está compuesta de tres ficheros ejecutables o programas, que pueden ser invocados desde la consola<sup>2</sup>:

1. `svm-scale`: estandariza el fichero de una base de datos, para que sea posible procesarlo con el clasificador SVM. Ten en cuenta que, utilizando un *kernel* tipo RBF, todas las variables de entrada deben estar en la misma escala, sino algunas variables recibirán más importancia que otras a la hora del cálculo de distancias.
2. `svm-train`: entrena un modelo SVM para un fichero de base de datos. El resultado del entrenamiento se guarda en un fichero con extensión `.model`.

<sup>1</sup>Disponible en <http://www.csie.ntu.edu.tw/~cjlin/libsvm/>

<sup>2</sup>Consulte el fichero `README` de la raíz de `libsvm` para comprobar cómo funcionan y qué argumentos reciben cada uno de estos programas.

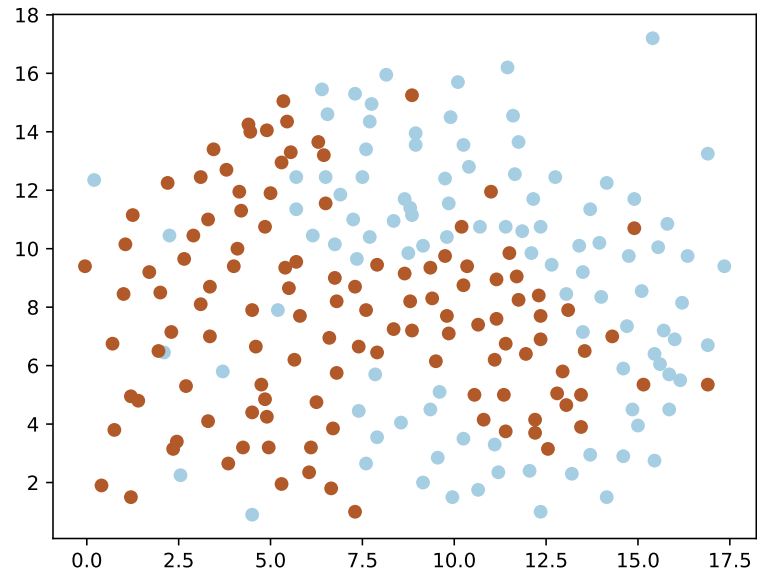


Figura 4: Dataset sintético o tipo *toy* #3.

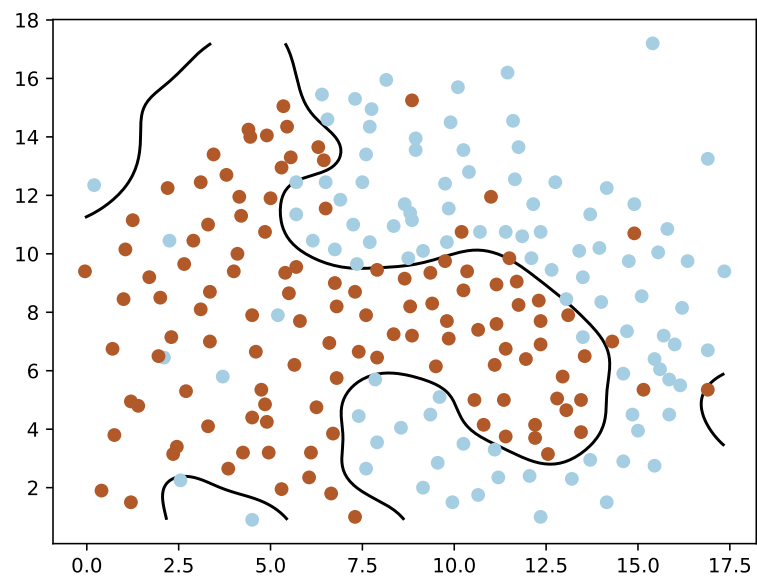


Figura 5: Dataset sintético o tipo *toy* #3 bien clasificado.

3. `svm-predict`: utiliza un modelo ya entrenado para ver como generaliza en una nueva base de datos (que debe tener la misma estructura). Es decir, realiza la fase de *test* del clasificador, comprobando cómo funciona para datos que no se han utilizado durante el entrenamiento.

Supongamos que tenemos dos ficheros de datos `train.libsvm` y `test.libsvm`. El proceso de entrenamiento estándar de una SVM implicaría, en primer lugar, entrenar la SVM utilizando `svm-train` y el fichero `train.libsvm`, generando un modelo SVM. En segundo lugar, utilizaríamos el programa `svm-predict` para ver que tal generaliza el modelo SVM sobre los datos `test.libsvm`.

Si utilizamos la clase `sklearn.svm.SVC` de `scikit-learn`, la estandarización tendremos que aplicarla nosotros, haciendo uso del objeto `StandardScaler`<sup>3</sup>. Aquí tienes un pequeño ejemplo de como se pueden estandarizar los datos:

```
1 from sklearn import preprocessing
2 scaler = preprocessing.StandardScaler()
3 X_train = scaler.fit_transform(X_train)
4 X_test = scaler.transform(X_test)
```

Por otro lado, la clase `sklearn.model_selection.StratifiedShuffleSplit` realiza una o varias particiones de una base de datos de forma “estratificada”, es decir, manteniendo la proporción de patrones de cada clase en la base de datos original<sup>4</sup>.

**Pregunta [8]:** Vamos a reproducir este proceso en Python. Divide el *dataset* sintético `dataset3.csv` en dos subconjuntos aleatorios de forma **estratificada**, con un 75 % de patrones en *train* y un 25 % de patrones en *test*. Realiza el proceso de entrenamiento completo (estandarización, entrenamiento y predicción), volviendo a optimizar los valores de  $C$  y  $\gamma$  que obtuviste en la última pregunta. Comprueba el porcentaje de buena clasificación que se obtiene para el conjunto de *test*. Repite el proceso más de una vez para comprobar que los resultados dependen mucho de la semilla utilizada para hacer la partición.

Este proceso tiene un pequeño inconveniente y es que es necesario especificar un valor para el parámetro  $C$  y un valor para el parámetro  $\gamma$  (ancho del *kernel*, para el caso en que queramos un clasificador no lineal) y, como ya viste en la primera parte de la práctica, estos parámetros son muy sensibles. Para evitar esto, existe una forma de ajustar automáticamente los parámetros, siguiendo el siguiente proceso de *validación cruzada anidada tipo K-fold*:

1. Realizamos una partición tipo  $K$ -fold (donde  $K$  es un parámetro) de los datos de entrenamiento. Es decir, dividimos los datos de entrenamiento en  $K$  subconjuntos disjuntos.
2. Para cada combinación de parámetros, realizamos  $K$  entrenamientos. Para cada entrenamiento  $k$  de los  $K$  totales:
  - a) Utilizamos el subconjunto  $k$  como conjunto de *test* y el resto de subconjuntos como conjunto de entrenamiento.
  - b) Acumulamos el error de *test* en una variable.

Al terminar, calculamos el error medio cometido durante los  $K$  entrenamientos. El proceso de validación cruzada se repite completo para todas las combinaciones posibles de parámetros  $C$  y  $\gamma$ . Por ejemplo, si vamos a darle 9 valores a  $C$  y 9 valores a  $\gamma$  y vamos a realizar una validación cruzada de tipo 5-fold, esto implica que realizaremos  $9 \times 9 \times 5 = 405$  entrenamientos. Nótese que hasta el momento no se han utilizado los datos de test.

3. Escogemos la combinación de parámetros que resultó en el menor error medio en el paso anterior y tomamos esos valores de parámetros como los valores óptimos.

<sup>3</sup><http://scikit-learn.org/stable/modules/generated/sklearn.preprocessing.StandardScaler.html>

<sup>4</sup>[http://scikit-learn.org/stable/modules/generated/sklearn.model\\_selection.StratifiedShuffleSplit.html](http://scikit-learn.org/stable/modules/generated/sklearn.model_selection.StratifiedShuffleSplit.html)

4. Utilizando esos parámetros, repetimos el proceso de entrenamiento pero ahora considerando el conjunto completo de entrenamiento (sin hacer subconjuntos).
5. Con el modelo obtenido en el paso anterior, comprobamos cuál es el error en test (fase real de *test*).

En `scikit-learn`, la búsqueda de parámetros la podemos realizar haciendo uso del objeto `sklearn.grid_search.GridSearchCV`<sup>5</sup>. Aquí tienes un pequeño ejemplo:

```
1 from sklearn.model_selection import GridSearchCV
2 Cs = np.logspace(-5, 15, num=11, base=2)
3 Gs = np.logspace(-15, 3, num=9, base=2)
4 optimo = GridSearchCV(estimator=svm_model, param_grid=dict(C=Cs, gamma=Gs),
5                       n_jobs=-1, cv=5)
6 optimo.fit(X_train, y_train)
7 print optimo.score(X_test, y_test)
```

**Pregunta [9]:** Amplía el código anterior para realizar el entrenamiento de la pregunta 8 sin necesidad de especificar los valores de  $C$  y  $\gamma$ . Compara los valores óptimos obtenidos para ambos parámetros con los que obtuviste a mano. Extiende el rango de valores a explorar, si es que lo consideras necesario.

**Pregunta [10]:** ¿Qué inconvenientes observas en ajustar el valor de los parámetros “a mano”, viendo el porcentaje de buena clasificación en el conjunto de test (lo que se hizo en la pregunta 8)?.

**Pregunta [11]:** Para estar seguros de qué has entendido como se realiza la búsqueda de parámetros, implementa de forma manual (sin usar `GridSearchCV`) la *validación cruzada anidada tipo K-fold* expuesta en esta sección. Te puede ser útil el uso de listas por comprensión y la clase `StratifiedKFold`. Compara los resultados con los que obtienes usando `GridSearchCV`.

### 3. Bases de datos reales

Una vez te has familiarizado con el clasificador SVM y con sus versiones lineal y no lineal, en esta parte de la práctica vamos a abordar distintos problemas reales, con objeto de ver su aplicabilidad en bases de datos más complejas.

#### 3.1. Base de datos *noMNIST*

*ILPD* contiene 405 patrones de entrenamiento y 174 patrones de test. El conjunto de datos se recogió en el noreste de Andhra Pradesh, India<sup>6</sup>. La etiqueta de clase se utiliza para dividir los pacientes en dos grupos (pacientes hepáticos o no). 441 registros corresponden a hombres, mientras que 142 corresponden a mujeres. Cualquier paciente cuya edad supere los 89 años aparece como de edad “90”. Hay un total de 10 variables de entrada que incluyen:

1. Age: Edad del paciente.
2. TB: Bilirrubina total.
3. DB: Bilirrubina directa.
4. AAP: Fosfatasa Alcalina.

<sup>5</sup>[http://scikit-learn.org/stable/modules/generated/sklearn.grid\\_search.GridSearchCV.html](http://scikit-learn.org/stable/modules/generated/sklearn.grid_search.GridSearchCV.html)

<sup>6</sup>Para más información, consultar [https://archive.ics.uci.edu/ml/datasets/ILPD+\(Indian+Liver+Patient+Dataset\)](https://archive.ics.uci.edu/ml/datasets/ILPD+(Indian+Liver+Patient+Dataset))

5. Sgpt: Alamina Aminotransferasa.
6. Sgot: Aspartato Aminotransferasa.
7. TP: Prótidos totales.
8. ALB: Albúmina.
9. A/G Ratio: Ratio de albúmina y globulina.
10. Gender: Género.

Esta base de datos presenta un problema de desequilibrio de clases, ya que hay 167 pacientes con enfermedad hepática y 416 pacientes sanos (aunque se han reducido después de eliminar datos erróneos). Además, un estudio reciente identificó que los modelos entrenados sobre esta base de datos presentan un sesgo de género, ya que los modelos tienden a predecir que los hombres tienen la enfermedad hepática y las mujeres no<sup>7</sup>. En esta práctica, las variables de entrada de esta base de datos han sido previamente estandarizadas en su versión `csv`.

**Pregunta [12]:** Utiliza el *script* que desarrollaste en la pregunta 9 para entrenar esta base de datos, pero sin dividir el dataset ni estandarizarlo. Utiliza como rangos  $C, \gamma \in \{10^{-3}, 10^{-2}, \dots, 10^3\}$ . Observa el valor de  $CCR$  obtenido para el conjunto de generalización y compáralo con el obtenido en prácticas anteriores. Al finalizar, toma nota de los valores óptimos obtenidos para los parámetros.

**Pregunta [13]:** Localiza dónde se especifica el valor de  $K$  para la validación cruzada interna y el rango de valores que se han utilizado para los parámetros  $C$  y  $\gamma$ . ¿Cómo podrías reducir el tiempo computacional necesario para realizar el experimento?. Prueba a establecer  $K = 3$ ,  $K = 5$  y  $K = 10$  y compara, utilizando una tabla, los tiempos computacionales obtenidos y los resultados de  $CCR$  en test.

Ahora vamos a utilizar la biblioteca `fairlearn` que empleaste en la LA3. La idea es analizar el rendimiento del mejor modelo desde diferentes perspectivas.

En primer lugar, ten en cuenta que `fairlearn` espera que la clase positiva sea 1 y la clase negativa sea 0 (y esto es al revés en el archivo `.csv` que os hemos proporcionado). Comprueba este problema en tu conjunto de datos, por ejemplo, identificando la etiqueta minoritaria (que debería ser un 1), y realiza el preprocesamiento adecuado.

Una de las ventajas de la clase `MetricFrame` de `fairlearn` es que puede proporcionar acceso muy fácilmente a las métricas calculadas para todo el conjunto (propiedad `overall`) y para cada grupo de la variable sensible (propiedad `by_group`).

**Pregunta [14]:** Utiliza estos métodos para calcular el rendimiento del modelo para cada clase y por grupos. Es decir, para el mejor modelo, calcula las matrices de confusión para entrenamiento y test, y muestra algunas métricas de rendimiento para test (por ejemplo, el *accuracy* y el FNR, pero puedes utilizar otras que consideres válidas). Analiza el comportamiento esperado según esta información y discute brevemente las limitaciones de las métricas al representar la matriz de confusión en este problema.

Para la siguiente pregunta, se pide que uses la validación cruzada manual que desarrollaste en la pregunta 11. La razón es que el objeto `GridSearchCV` no está preparado para tratar con las métricas de `fairlearn` porque necesitan la característica sensible.

Vamos a optimizar el modelo para obtener justicia en la clasificación cambiando los criterios de validación cruzada. Sabemos que los modelos entrenados en ILPD tienden a discriminar a las mujeres produciendo modelos menos precisos que suelen estar sesgados con mayores tasas de falsos negativos para este grupo. Una estrategia de mitigación consiste en elegir una métrica

<sup>7</sup>Para más información, consultar el artículo Straw, I., and Wu, H. (2022). Investigating for bias in healthcare algorithms: A sex-stratified analysis of supervised machine learning models in liver disease prediction. *BMJ Health & Care Informatics*, 29(1), e100457. <https://doi.org/10.1136/bmjhci-2021-100457>



de error que penalice este comportamiento durante la optimización de los hiperparámetros. En la clasificación desequilibrada, una de las opciones más comunes es evaluar el rendimiento del modelo con la media geométrica de la sensibilidad para cada clase (nos referimos aquí a la clase a predecir, no a la característica sensible):

$$GM = \sqrt[Q]{\prod_{q=1}^Q S_q},$$

donde  $S_q$  es la sensibilidad (ratio de precisión) del clasificador para la clase  $q$ . Si  $GM = 0$ , el clasificador clasifica totalmente mal al menos una clase. Para el caso binario, esta expresión genérica se convierte en:

$$GM = \sqrt{S_+ \cdot S_-}.$$

**Pregunta [15]:**

Podemos adaptar esta idea a nuestro problema calculando los  $GM$  de los grupos (mujeres y hombres), de forma que podamos guiar la optimización por el  $GM$  del grupo de peor rendimiento o por la media de los  $GM$  de todos los grupos. Prueba estas dos estrategias y compara los resultados con el enfoque estándar (*accuracy* global).

### 3.2. Referencia de rendimiento para el conjunto de datos *ILDP*

Aquí mostramos algunas matrices de confusión y métricas de rendimiento para este problema. Hemos utilizado tres criterios de optimización, pero otras propuestas son bienvenidas:

1. *accuracy* global.
2. Media geométrica para las mujeres.
3. Promedio de las medias geométricas para mujeres y para hombres.

```

1      Optimization criteria: global accuracy
2
3      FINAL => C 10.000000, Gamma 10.000000, Validation Score 0.728395, TEST
4      Score 0.712644
5      Overall test CM:
6      [[121  2]
7      [ 48  3]]
8      Grouped test CM:
9      [[93  2]
10     [33  2]]
11     [[28  0]
12     [15  1]]
13     TEST accuracy: 0.712644
14     TEST FNR (men): 0.942857
15     TEST FNR (women): 0.937500
16
17     Optimization criteria: Geometric Mean of Women
18
19     FINAL => C 1000.000000, Gamma 0.100000, Validation Score 0.636903,
20     TEST Score 0.422577
21     Overall test CM:
22     [[89 34]
23     [33 18]]
24     Grouped test CM:
25     [[69 26]
26     [21 14]]
27     [[20  8]
28     [12  4]]
29     TEST accuracy: 0.614943

```

```

29      TEST FNR (men): 0.600000
30      TEST FNR (women): 0.750000
31
32      Optimization criteria: (Geometric Mean of Men + Geometric Mean of
33                               Women) / 2
34
35      FINAL => C 1000.000000, Gamma 0.100000, Validation Score 0.615596,
36               TEST Score 0.480791
37      Overall test CM:
38      [[89 34]
39       [33 18]]
40      Grouped test CM:
41      [[69 26]
42       [21 14]]
43      [[20  8]
44       [12  4]]
45      TEST accuracy: 0.614943
46      TEST FNR (men): 0.600000
47      TEST FNR (women): 0.750000

```

### 3.3. Base de datos clasificación de *spam*

Uno de los campos donde se utiliza el aprendizaje automático con un mayor éxito es la detección automática de *spam* en servidores de correo. Como dato significativo, ten en cuenta que el correo que recibes cada día en tu cuenta de la Universidad de Córdoba (UCO) es solo un 5 % de los correos reales que te llegan, de manera que el 95 % restante es detectado y descartado automáticamente por filtros anti-*spam*. La mayoría de estos filtros utilizan internamente técnicas de clasificación binaria (correo tipo *spam*,  $y = 1$ , o no *spam*,  $y = 0$ ) y cada vez está ganando mayor aceptación el uso de SVM. En esta sección, vamos a explicar los pasos que se llevan a cabo para preprocesar el texto de un correo y convertirlo en un vector de características que podamos utilizar como entrada en cualquier clasificador. Posteriormente, se te pedirá que utilices SVM para clasificar una base de datos real.

```

1  > Anyone knows how much it costs to host a web portal ?
2  >
3  Well, it depends on how many visitors youre expecting. This can be
4  anywhere from less than 10 bucks a month to a couple of $100. You
5  should checkout http://www.rackspace.com/ or perhaps Amazon EC2 if
6  youre running something big..
7  To unsubscribe yourself from this mailing list, send an email to:
8  groupname-unsubscribe@egroups.com

```

Figura 6: Ejemplo de correo

Imaginemos un correo cuyo cuerpo sea como el de la Figura 6. Podemos observar cómo este correo contiene una URL, una dirección de correo electrónico (al final), números y una cantidad en dólares. Muchos e-mails tendrán este tipo de cosas (direcciones de correo, URLs) pero, seguramente, el valor concreto (la URL concreta o el correo concreto) serán diferentes en cada e-mail. Normalmente, para la tarea concreta de clasificación (es *spam* o no), te interesa el hecho de que el correo contenga una dirección URL o un e-mail o un número, pero no su valor concreto. Por tanto, el proceso que vamos a seguir “normaliza” todas las URLs para que se traten del mismo modo, al igual que todos los números, todas las direcciones de correo, etc. Para ello, cada vez que aparezca en nuestro correo una dirección `http`, la reemplazaremos por la palabra reservada `httpaddr` que indica que había una dirección. Esto está demostrado que mejora el porcentaje de buena clasificación en detección de *spam*, porque casi siempre los *spammers* cambian de forma aleatoria las direcciones URL (y así es más difícil detectarlos).

Se te va a proporcionar una base de datos de correos que es parte del repositorio *SpamAssassin*

*Public Corpus*<sup>8</sup>. Cada correo de esta base de datos ha recibido el siguiente preprocesamiento:

1. Todo el correo ha sido escrito en letras minúsculas (es decir, “INDICATE” se va a tratar igual que “Indicate” o que “indicate”). A menudo, los *spammers* usan esta técnica para evitar ser detectados.
2. Eliminar todo código HTML: todos los correos serán procesados como texto plano, de manera que se elimina cualquier etiqueta de formato de tipo HTML.
3. Normalizar las URL: todas las URL se reemplazan por `httpaddr`.
4. Normalizar las direcciones de e-mail: todas las direcciones de e-mail se reemplazan por `emailaddr`.
5. Normalizar los números: todos los números se reemplazan por `number`.
6. Normalizar los dólares: todos los signos de dolar (\$, asociados con frecuencia a correos de *spam*) se reemplazan por `dollar`.
7. *Stemming* de palabras: el *stem* de una palabra es su raíz. Para la detección de *spam* es buena idea utilizar las raíces de las palabras, en lugar de las palabras completas. Por ejemplo, las palabras “discount”, “discounts”, “discounted” y “discounting” se reemplazan todas por “discount”. Hay programas (*stemmers*) que realizan este proceso de manera automática, utilizando para ello diccionarios específicos.
8. Quitar todo lo que no son palabras, por ejemplo, signos de puntuación o símbolos poco frecuentes. Además, todos los espacios, tabuladores o saltos de líneas que vayan consecutivos son sustituidos por un solo carácter de espacio.

El resultado de todos estos pasos sobre el correo anterior puede verse en la Figura 7. Con esta representación del correo es mucho más fácil extraer características que transformen el correo en un vector de números.

```
1 anyone know how much it cost to host a web portal well it depend on how mani visitor your
   expect thi can be anywher from less than number buck a month to a coupl of
   dollarnumb you should checkout httpaddr or perhap amazon ecnumb if your run someth
   big to unsubscrib yourself from thi mail list send an email to emailaddr
```

Figura 7: Ejemplo de correo ya preprocesado

Tras este paso, debemos elegir cuáles de las palabras que tenemos como resultado vamos a utilizar para clasificar *spam*. A esto se le llama lista de vocabulario. En el caso que nos ocupa, se han seleccionado 1899 palabras, que son las que aparecen con más frecuencia. Estas palabras están en el fichero `vocab.txt` y se muestran también en la Figura 8. Estas palabras son aquellas que aparecían al menos 100 veces en el total de e-mails de la base de datos.

Utilizando esta lista de vocabulario, el e-mail de la Figura 7 queda reducido a un conjunto de índices que representan cada una de las palabras (suponiendo que la palabra esté en el vocabulario, sino, simplemente se ignora por ser una palabra muy rara). El resultado de esa transformación puede verse en la Figura 9.

Desde el punto de vista de la detección de *spam*, nos interesa si la palabra del vocabulario aparece, pero no en qué orden. Por tanto, el último paso transforma el vector de índices anterior en un vector de valores binarios. Para nuestro caso tendremos un vector  $\mathbf{x} \in \mathbb{R}^{1899}$  y la posición  $i$  valdrá 1 ( $x_i = 1$ ) si la palabra del  $i$ -ésima del vocabulario aparece en el e-mail ( $x_i = 0$  si no aparece). Este vector será el vector que se use para clasificar si el correo es o no *spam*. Usando esta representación, da igual que las palabras aparezcan más de una vez.

<sup>8</sup><http://spamassassin.apache.org/publiccorpus/>

```

1 1 aa
2 2 ab
3 3 abil
4 ...
5 916 know
6 ...
7 1898 zero
8 1899 zip

```

Figura 8: Lista de vocabulario

```

1 86 916 794 1077 883 370 1699 790 1822 1831 883 431 1171 794 1002 1893 1364 592 1676 238
   162 89 688 945 1663 1120 1062 1699 375 1162 479 1893 1510 799 1182 1237 810 1895
   1440 1547 181 1699 1758 1896 688 1676 992 961 1477 71 530 1699 531

```

Figura 9: Ejemplo de correo preprocesado y del que se han extraído los índices de las palabras.

Dispones de dos ficheros de datos donde todo este proceso ya se ha realizado. El fichero de entrenamiento contiene 4000 correos, mientras que el de *test* tiene 1000 correos. Ambos utilizan la lista de vocabulario de 1899 palabras contenida en `vocab.txt`. Por tanto, cada patrón tiene 1899 valores binarios.

**Pregunta [16]:** Debes entrenar un modelo **lineal** de SVM con valores  $C = 10^{-2}$ ,  $C = 10^{-1}$ ,  $C = 10^0$  y  $C = 10^1$ . Para ello utiliza un *script* similar al que usaste para la pregunta 9. Compara los resultados y establece la mejor configuración.

**Pregunta [17]:** Para la mejor configuración, construye la matriz de confusión y establece cuáles son los correos en los que la SVM se equivoca. Consulta las variables de entrada para los correos que no se clasifican correctamente y razona el motivo. Ten en cuenta que para cada patrón, cuando  $x_i$  es igual a 1 quiere decir que la palabra  $i$ -ésima del vocabulario aparece, al menos una vez, en el correo.

**Pregunta [18]:** Compara los resultados obtenidos con los resultados utilizando una red RBF. Para ello, haz uso del programa desarrollado en la práctica anterior. Utiliza solo una semilla (la que mejor resultado obtenga).

**Pregunta [19]:** Entrena una SVM no lineal y compara los resultados.

## 4. Entregables

Los ficheros a entregar serán los siguientes:

- Memoria de la práctica en un fichero pdf que conteste a todas las cuestiones que se han ido planteando en este guión. Incluye las tablas o imágenes que consideres necesarias para demostrar que se han realizado los ejercicios que se piden. La memoria de la práctica deberá incluir, al menos, el siguiente contenido:
  - Portada con el número de práctica, título de la práctica, asignatura, titulación, escuela, universidad, curso académico, nombre, DNI y correo electrónico del alumno.
  - Índice del contenido de la memoria con numeración de las páginas.
  - Respuesta a todas las cuestiones y ejercicios que se plantean en este guión.
  - Referencias bibliográficas u otro tipo de material distinto del proporcionado en la asignatura que se haya consultado para realizar la práctica (en caso de haberlo hecho).

- Un fichero comprimido con todos los *scripts*, ficheros de datos y ficheros de salida que hayas necesitado para contestar a las cuestiones.

Opcionalmente, podéis presentar tanto la memoria como el código de forma conjunta en un cuaderno de Jupyter.