

APR (E.T.S. de Ingeniería Informática)
Curso 2020-2021

*Proyecto de prácticas. Reconocimiento de dígitos
manuscritos: MNIST*

Jorge Civera, Francisco Casacuberta, Enrique Vidal
Departamento de Sistemas Informáticos y Computación
Universitat Politècnica de València

Última actualización: 16/12/2020- 14:34:27

Índice

1. Objetivos	3
2. Mixtura de gaussianas	3
2.1. Algoritmo EM	4
2.2. Implementación	6
2.3. Tarea MNIST	10
3. Máquinas de Vectores Soporte	11
3.1. Ejemplo de entrenamiento y clasificación con LibSVM	11
3.1.1. Corpus de ejemplo	11
3.1.2. Entrenamiento	12
3.1.3. Clasificación y estimación de la tasa de acierto	13
3.2. Tarea MNIST	15
4. Redes Neuronales Multicapa	16
4.1. Ejemplo de entrenamiento y clasificación con nnet	16
4.1.1. Definición de entrenamiento y validación	16
4.1.2. Preproceso de los datos	17
4.1.3. Entrenamiento	18
4.1.4. Clasificación y estimación de la tasa de error	19
4.2. Tarea MNIST	20
A. Tutorial sobre Octave	22
A.1. Introducción	22
A.2. Órdenes básicas de <i>Octave</i>	22
A.2.1. Aritmética básica	23

A.2.2. Operadores básicos en vectores y matrices	24
A.2.3. Funciones básicas en vectores y matrices	26
A.2.4. Carga y salvado de datos	28
A.2.5. Funciones Octave	29
A.2.6. Programas Octave	30
A.3. Ejercicios propuestos	32
B. Tarea de clasificación: MNIST	33
B.1. Introducción	33
B.2. Carga de datos	33
B.3. Visualización de dígitos	35
C. Clasificador gaussiano	36
C.1. Estimación de parámetros y clasificación	36
C.2. Implementación y experimentación	37
D. Función plot en Octave	43
E. Recordatorio de teoría de SVM	44

1. Objetivos

En el marco de un aprendizaje basado en proyectos, este proyecto de prácticas es la continuación del realizado en la asignatura de Percepción del pasado curso académico. Al igual que en Percepción, el principal objetivo de este proyecto de prácticas es la implementación y evaluación de diversos clasificadores estudiados en teoría sobre la tarea real de reconocimiento de dígitos manuscritos MNIST. Este objetivo principal se desglosa en subobjetivos básicos que se pueden entender como fases o hitos de este proyecto:

1. Implementar el algoritmo EM para mixtura de gaussianas.
2. Evaluar el clasificador de mixtura de gaussianas y su interacción con la técnica de reducción de dimensionalidad *Principal Component Analysis* PCA en la tarea MNIST.
3. Comprender los conceptos teóricos de *Support Vector Machines* (SVM) mediante su aplicación a pequeñas tareas de clasificación.
4. Evaluar el clasificador basado en SVM en la tarea MNIST.
5. Comprender y aplicar el preproceso necesario a un conjunto de datos para entrenar y evaluar una red neuronal.
6. Evaluar un clasificador basado en redes neuronales en la tarea MNIST, y su interacción con la técnica PCA.

Para la realización de este proyecto de prácticas se supone que previamente has adquirido experiencia en el uso de `octave`, `gnuplot` y *shell scripts*, tanto en la asignatura de Sistemas Inteligentes como en la asignatura de Percepción. Si no es el caso, puedes refrescar tus conocimientos y habilidades sobre dichas herramientas recurriendo al tutorial incluido en el Apéndice A. Este tutorial es el mismo que se utilizó en la asignatura de Percepción. Asimismo, te recomendamos que refresques la descripción de la tarea MNIST que encontrarás en el Apéndice B. La tarea MNIST, así como otros pequeños conjuntos de datos de prueba que se utilizarán en este proyecto, están disponibles en PoliformaT a través del fichero `data.tgz`.

Como regla general, se recomienda leer en su totalidad la sección correspondiente al clasificador con el que se trabajará en las siguientes sesiones. Esto permite centrarse en el trabajo a desarrollar en la sesión de laboratorio y aprovechar mejor el tiempo.

Finalmente, la evaluación del proyecto de la asignatura consta de un total de 3 puntos que se distribuyen de manera uniforme con 1 punto para cada tipo de clasificador.

2. Mixtura de gaussianas

En esta primera parte del proyecto estudiaremos una generalización del clasificador gaussiano visto en la asignatura de Percepción (ver Apéndice C), el clasificador de

mixtura de gaussianas. Las mixturas, ya sean de gaussianas o de cualquier otra distribución, nos permiten introducir una instanciación concreta del algoritmo EM estudiado en teoría. Por ello, primero desarrollaremos la teoría asociada a la estimación de los parámetros de una mixtura de gaussianas en el marco del algoritmo EM, y seguidamente veremos su implementación con las peculiaridades prácticas que ello conlleva.

Los ficheros necesario para desarrollar esta parte del proyecto están disponibles en Poliformat en el fichero `mixgaussian.tgz`.

2.1. Algoritmo EM

Las mixturas de distribuciones de probabilidad, en nuestro caso condicionada a la clase $p(\mathbf{x} \mid c)$, se desarrollan introduciendo una variable oculta k que indica la componente de la mixtura que genera la muestra \mathbf{x}

$$p(\mathbf{x} \mid c) = \sum_{k=1}^K p(\mathbf{x}, k \mid c) = \sum_{k=1}^K p(k \mid c) p(\mathbf{x} \mid k, c) \quad (1)$$

donde $p(k \mid c)$ es la probabilidad a priori de la componente condicionada a la clase, y $p(\mathbf{x} \mid k, c)$ es una distribución de probabilidad condicionada no sólo a la clase, sino también a la componente dentro de esa clase. En nuestro caso $p(\mathbf{x} \mid k, c)$ está modelizada mediante una distribución gaussiana

$$p(\mathbf{x} \mid k, c) \sim \mathcal{N}_D(\boldsymbol{\mu}_{ck}, \Sigma_{ck}), \quad c = 1, \dots, C \quad k = 1, \dots, K$$

pero podría haber sido modelizada mediante una distribución multinomial o cualquier otra distribución de probabilidad.

Observad como la componente de cada clase tendrá su propia media y matriz de covarianzas. El conjunto de parámetros de este modelo de mixturas es

$$\boldsymbol{\theta} = (P_1, \dots, P_C, \boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_C),$$

donde

$$\boldsymbol{\theta}_c = (p_{c1}, \dots, p_{cK}, \mu_{c1} \dots \mu_{cK}, \Sigma_{c1}, \dots, \Sigma_{cK}).$$

La estimación de estos parámetros sería trivial si para cada muestra supiéramos que componente de cada clase la generó, es decir, si conociéramos el valor de z_n . En ese caso, la estimación sería:

$$\hat{P}(c) = \frac{N_c}{N} \quad (2)$$

$$\hat{p}_{ck} = \frac{\sum_{n: c_n=c \wedge z_n=k} 1}{N_c} = \frac{N_{ck}}{N_c} \quad (3)$$

$$\hat{\boldsymbol{\mu}}_{ck} = \frac{1}{N_{ck}} \sum_{n: c_n=c \wedge z_n=k} \mathbf{x}_n \quad (4)$$

$$\hat{\Sigma}_{ck} = \frac{1}{N_{ck}} \sum_{n: c_n=c \wedge z_n=k} (\mathbf{x}_n - \hat{\boldsymbol{\mu}}_{ck})(\mathbf{x}_n - \hat{\boldsymbol{\mu}}_{ck})^t \quad (5)$$

donde N_{ck} es el número de muestras de la componente k de la clase c . Es decir, simplemente tendríamos que calcular los parámetros de cada componente de cada clase con las muestras que sabemos, gracias a z_n , que pertenecen a dicha componente de esa clase.

Sin embargo, desconocemos qué componente generó cada muestra, es decir, el valor de cada z_n está oculto. Por ello, como has estudiado en teoría, el algoritmo EM nos permite estimar de manera iterativa los parámetros de un modelo (paso M) donde hay variables ocultas, calculando una distribución de probabilidad sobre dichas variables ocultas (paso E).

El paso E en la iteración t de un modelo de mixturas dada la estimación actual de los parámetros $\theta(t)$ es:

$$\begin{aligned}
 z_{nk}^{(t)} &= p(z_n = k \mid \mathbf{x}_n, c_n = c; \theta(t)) \\
 &= \frac{p(z_n = k, \mathbf{x}_n, c_n = c; \theta(t))}{p(\mathbf{x}_n, c_n = c; \theta(t))} \\
 &= \frac{p(z_n = k, \mathbf{x}_n, c_n = c; \theta(t))}{\sum_{k'=1}^K p(z_n = k', \mathbf{x}_n, c_n = c; \theta(t))} \\
 &= \frac{p(z_n = k \mid c_n = c; \theta(t)) p(\mathbf{x}_n \mid z_n = k, c_n = c; \theta(t))}{\sum_{k'=1}^K p(z_n = k' \mid c_n = c; \theta(t)) p(\mathbf{x}_n \mid z_n = k', c_n = c; \theta(t))} \\
 &= \frac{p(k \mid c_n) p(\mathbf{x}_n \mid k, c_n)}{\sum_{k'=1}^K p(k' \mid c_n) p(\mathbf{x}_n \mid k', c_n)} \tag{6}
 \end{aligned}$$

Es decir, la probabilidad de la muestra \mathbf{x} de acuerdo a la distribución de probabilidad de la componente k de la clase c . También, se puede interpretar como el grado de pertenencia de la muestra \mathbf{x} a la componente k de la clase c . Desde el punto de vista del paso M, que ahora detallaremos, z_{nk} es el peso o contribución parcial de una muestra a una componente, de forma que una misma muestra puede contribuir parcialmente a la estimación de los parámetros de varias componentes.

En el paso M se estiman los parámetros del modelo teniendo una estimación del grado de pertenencia de cada muestra de entrenamiento a cada componente. La estimación de la probabilidad a priori de cada clase queda fuera de la estimación por el algoritmo EM (ver Ec. 2). El resto de parámetros siguen una estimación similar a las Eqs. 3, 4 y 5, pero teniendo en cuenta la idea de que una muestra contribuye parcialmente acorde a z_{nk} a la estimación de los parámetros de la componente k en la clase c :

$$\hat{p}_{ck}^{(t+1)} = \frac{1}{N_c} \sum_{n: c_n=c} z_{nk}^{(t)} \tag{7}$$

$$\hat{\mu}_{ck}^{(t+1)} = \frac{1}{\sum_{n: c_n=c} z_{nk}^{(t)}} \sum_{n: c_n=c} z_{nk}^{(t)} \mathbf{x}_n \tag{8}$$

$$\hat{\Sigma}_{ck}^{(t+1)} = \frac{1}{\sum_{n: c_n=c} z_{nk}^{(t)}} \sum_{n: c_n=c} z_{nk}^{(t)} (\mathbf{x}_n - \hat{\mu}_{ck})(\mathbf{x}_n - \hat{\mu}_{ck})^t \tag{9}$$

2.2. Implementación

Antes de pasar a la implementación es necesario tratar la inicialización del algoritmo EM. La manera más habitual de inicializar el algoritmo EM es definir una inicialización de los parámetros del modelo $\theta^{(0)}$ que dependen de la variable oculta, pero también sería posible definir una inicialización de la distribución de probabilidad sobre la variable oculta.

En nuestro caso, inicializaremos los parámetros correspondientes a la probabilidad a priori, la media y matriz de covarianza de cada componente. Seguidamente se muestra el código utilizado para la inicialización:

```

1 sigma=cell(C,K);
2 for c=classes'
3   ic=find(c==classes);
4   pkGc{ic}(1:K)=1/K;
5   idc=find(xl==c);
6   Nc=rows(idc);
7   mu{ic}=X(idc(randperm(Nc,K)),:)' ;
8   sigma(ic,1:K)=alpha*cov(X(idc,:),1)/K+(1-alpha)*eye(D);
9 end

```

La línea 1 define una matriz de celdas $C \times K$ para almacenar las matrices de covarianza de cada una de las C clases y K componentes por clase. En el bucle desde la línea 2 a la línea 9 se inicializan los parámetros de cada clase. La línea 4 inicializa la probabilidad de cada componente como una distribución uniforme. La línea 7 define la media de cada componente como una muestra aleatoria de la clase c . Finalmente, en la línea 8 se inicializa la matriz de covarianzas de cada componente como la matriz de covarianzas de la clase dividida por el número de componentes $\text{cov}(X(\text{idc},:),1)/K$, que es suavizada mediante *flat smoothing* con la matriz identidad. Como puedes deducir esta inicialización es arbitraria, pero ha demostrado funcionar bien en la práctica.

Al igual que en el clasificador gaussiano haremos uso de una función auxiliar para calcular la probabilidad de cada muestra¹, pero en este caso para cada componente k . Esto se corresponde con el logaritmo del numerador de la Ec. 6, es decir, $\log p(k | c) + \log p(\mathbf{x} | k, c)$:

```

10 function [zk] = compute_zk(pkGc,mu,sigma,X)
11   D=columns(X);
12   cons=log(pkGc);
13   cons=cons-0.5*D*log(2*pi);
14   cons=cons-0.5*logdet(sigma);
15   cons=cons-0.5*mu'*pinv(sigma)*mu;
16   lin=X*(mu'*pinv(sigma))';

```

¹De cada muestra, pero todas ellas a la vez, independientemente y en paralelo, mediante operaciones matriciales aprovechando las capacidades de Octave.

```

17 qua=-0.5*sum((X*pinv(sigma)).*X,2);
18 zk=qua+lin+cons;
19 end

```

Observarás que es casi idéntica a la función auxiliar del clasificador gaussiano, a excepción de la línea 13 que incluye el término constante porque vamos a hacer una estimación de la probabilidad a posteriori y no únicamente una clasificación como en el clasificador gaussiano.

Una vez han sido inicializados los parámetros, el algoritmo EM ejecuta los pasos E y M de manera iterativa hasta convergencia (o un número máximo de iteraciones). En nuestra implementación la condición de convergencia es que el incremento relativo de la log verosimilitud entre dos iteraciones consecutivas no supere cierto umbral

$$\frac{|L(\boldsymbol{\theta}; \mathbf{X})^{(t+1)} - L(\boldsymbol{\theta}; \mathbf{X})^{(t)}|}{|L(\boldsymbol{\theta}; \mathbf{X})^{(t)}|} < \epsilon$$

donde

$$L(\boldsymbol{\theta}; \mathbf{X}) = \sum_n \log p(c_n) + \log p(\mathbf{x}_n | c_n)$$

donde $p(\mathbf{x}_n | c_n)$ se define en la Ec. 1. De manera similar, $p(\mathbf{x}_n | c_n)$ es el denominador de la Ec. 6, es decir, que calcularemos $p(\mathbf{x}_n | c_n)$ como el sumatorio sobre las K componentes

$$p(\mathbf{x}_n | c_n) = \sum_{k'=1}^K p(k' | c_n) p(\mathbf{x}_n | k', c_n)$$

invocando la función `compute_zk` para cada componente. En cada iteración del algoritmo EM se ejecutan los pasos E y M por cada clase, estimando z_{nk} y $\boldsymbol{\theta}^{(t+1)}$, respectivamente:

```

20 for c=classes'
21
22     % E step: Estimate znk
23     ic=find(c==classes);
24     idc=find(xl==c);
25     Nc=rows(idc);
26     Xc=X(idc,:);
27     z=[];
28     for k=1:K
29         z(:,k)=compute_zk(pkGc{ic}(k),mu{ic}(:,k),sigma{ic,k},Xc);
30     end
31
32     % Robust computation of znk and log-likelihood
33     maxz=max(z,[],2);
34     z=exp(z-maxz);
35     sumz=sum(z,2);
36     z=z./sumz;

```

```

37     L=L+Nc*log(pc(ic))+sum(maxz+log(sumz));
38
39     % M step: parameter update
40     % HERE YOUR CODE FOR PARAMETER ESTIMATION
41
42 end
43
44 % Likelihood divided by the number of training samples
45 L=L/N;

```

Las líneas 22 a 36 definen el paso E, mientras que las líneas de código que implementarás para el paso M debes desarrollarlas a partir de la línea 40. En cuanto al paso E, la línea 29 estima la log probabilidad de cada muestra de la clase c para cada componente k dando lugar a cada vector columna de la matriz \mathbf{z} . Seguidamente en las líneas 33 a 36 se realiza la estimación de z_{nk} de la Ec. 6, pero de manera *robusta* para manejar valores pequeños de log probabilidad².

$$\begin{aligned}
z_{nk}^{(t)} &= \frac{p(k | c) \cdot p(\mathbf{x} | k, c)}{\sum_{k'=1}^K p(k' | c) \cdot p(\mathbf{x} | k', c)} \\
&= \frac{\frac{p(k|c) \cdot p(\mathbf{x}|k, c)}{\max_{k''} p(k''|c) \cdot p(\mathbf{x}|k'', c)}}{\sum_{k'=1}^K \frac{p(k'|c) \cdot p(\mathbf{x}|k', c)}{\max_{k''} p(k''|c) \cdot p(\mathbf{x}|k'', c)}} \\
&= \frac{\exp \left(\log \left(\frac{p(k|c) \cdot p(\mathbf{x}|k, c)}{\max_{k''} p(k''|c) \cdot p(\mathbf{x}|k'', c)} \right) \right)}{\sum_{k'=1}^K \exp \left(\log \left(\frac{p(k'|c) \cdot p(\mathbf{x}|k', c)}{\max_{k''} p(k''|c) \cdot p(\mathbf{x}|k'', c)} \right) \right)} \\
&= \frac{\exp (\log p(k | c) + \log p(\mathbf{x} | k, c)) - \max_{k''} \log p(k'' | c) + \log p(\mathbf{x} | k'', c))}{\sum_{k'=1}^K \exp (\log p(k' | c) + \log p(\mathbf{x} | k', c) - \max_{k''} \log p(k'' | c) + \log p(\mathbf{x} | k'', c))}
\end{aligned}$$

Es decir, para cada muestra se divide la log probabilidad de cada componente por la log probabilidad de aquella componente con máxima log probabilidad. Después se calcula el logaritmo compensado con la exponenciación y se normaliza por la suma del total de las componentes.

Sobre el código, la línea 33 calcula la componente de máxima log probabilidad, restamos el máximo de la log probabilidad en la línea 34, calculamos la suma para todas las componentes en la línea 35 y normalizamos en la línea 36. Nótese que estas operaciones se hacen a la vez para todas las muestras de la clase c aprovechando la capacidad de Octave de trabajar con matrices y vectores de manera más eficiente.

La log verosimilitud para las muestras de la clase c se calcula en la línea 37 aprovechando el cálculo robusto realizado

$$L(\boldsymbol{\theta}; \mathbf{X}) = \sum_c L(\boldsymbol{\theta}_c; \mathbf{X}_c)$$

²<https://en.wikipedia.org/wiki/LogSumExp>

donde

$$\begin{aligned} L(\boldsymbol{\theta}_c; \mathbf{X}_c) &= \sum_{n: c_n=c} \log P(c_n) + \log p(\mathbf{x}_n | c_n) \\ &= N_c \log P(c_n) + \sum_{n: c_n=c} \log \sum_{k=1}^K p(\mathbf{x}_n, k | c_n) \end{aligned}$$

siendo $p(\mathbf{x}_n | c_n)$ el denominador de z_{nk} que se calcula en la línea 35. Sin embargo, debido al cálculo robusto a $p(\mathbf{x}_n | c_n)$ le ha sido restado el máximo de cada componente, y por tanto para compensar debemos sumar el máximo de cada componente como se observa en la línea 37. En la línea 45, una vez se han procesado todas las muestras, se puede observar como la log verosimilitud se normaliza por el número de muestras, pero es algo opcional.

En cuanto al paso M, será una implementación directa de las Ecs. 7, 8 y 9 aprovechando las operaciones matriciales en Octave (ver Ejercicio 2.1).

La clasificación de cada muestra es un cálculo muy similar a la estimación de la log verosimilitud, ya que la función discriminante que se calcula es la probabilidad de cada muestra de acuerdo a los parámetros de cada una de las clases involucradas.

$$\begin{aligned} c^*(\mathbf{x}) &= \operatorname{argmax}_{c=1,\dots,C} \log P(c) + \log p(\mathbf{x} | c) \\ &= \operatorname{argmax}_{c=1,\dots,C} \log P(c) + \log \sum_{k=1}^K p(\mathbf{x}, k | c) \end{aligned}$$

Finalmente, el fichero `mixgaussian.m` de PoliformaT contiene la implementación descrita anteriormente. Por favor, dedica unos minutos a leer el código y comprobar que comprendes la implementación ya realizada.

Ejercicio 2.1 (obligatorio: 0.3 puntos). Implementa el paso M de estimación de los parámetros que gobiernan la mixtura de gaussianas y que se detallan en las Ecs. 7, 8 y 9. Ten en cuenta que no debes utilizar bucles sino operaciones matriciales, a excepción de la estimación de la matriz de covarianzas de cada componente, que necesitarás un bucle que recorra las componentes. Tras la estimación de la matriz de covarianzas de cada componente, recuerda suavizarla mediante *flat smoothing* con la matriz identidad.

Ejercicio 2.2 (opcional: no entregable). Para realizar una primera comprobación del correcto funcionamiento del clasificador de mixtura de gaussianas, compara su funcionamiento para el caso de una única componente por mixtura con el clasificador gaussiano. Para ello, realiza la misma partición que se define en el experimento del clasificador gaussiano descrita en el Apéndice C. Te recomendamos que copies el fichero `gaussian-exp.m` (disponible en PoliformaT) como `mixgaussian-exp.m` y lo modifiques adecuadamente.

2.3. Tarea MNIST

En esta sección vamos a aplicar el clasificador de mixtura de gaussianas a la tarea MNIST para estudiar si sus resultados son competitivos respecto a otros clasificadores.

Ejercicio 2.3 (obligatorio: 0.5 puntos). Realiza un experimento donde se evalúe el error de clasificación en función del número de componentes por mixtura del clasificador ($K = 1, 2, 5, 10, 20, 50, 100$), para un número de dimensiones PCA variable ($D = 1, 2, 5, 10, 20, 50, 100$) y probando algunos de los valores de α de suavizado *flat smoothing* que mejores resultados han proporcionado en el clasificador gaussiano. Para ello, puedes partir del script `mixgaussian-exp.m` para crear el script `pca+mixgaussian-exp.m` que mantenga la misma partición en entrenamiento y validación, y realice la exploración de parámetros propuesta.

Recuerda que el objetivo de este experimento es determinar los valores de los parámetros del clasificador (K , D y α) que minimizan el error de clasificación en el conjunto de validación. En el ejercicio 2.4, estos valores óptimos de los parámetros serán utilizados para entrenar y evaluar un clasificador final en los conjuntos oficiales MNIST de entrenamiento y test, respectivamente.

Representa gráficamente las tasas de error obtenidas en el conjunto de validación. Como punto de partida puedes utilizar el fichero `gaussian-exp.gnp` que utiliza el fichero de resultados `gaussian-exp.out` para generar la gráfica `gaussian-exp.eps` (todos ellos disponibles en PoliformaT). Te recomendamos que generes una gráfica independiente por cada valor de suavizado α que hayas evaluado. A su vez, cada gráfica tendrá tantas curvas como valores de PCA diferentes hayas probado. Finalmente, cada curva mostrará la evolución del error de clasificación (eje y) en función del número de componentes por mixtura (eje x). Para mejorar la legibilidad de la representación gráfica de los resultados puedes descartar aquellas curvas asociadas a valores de PCA cuyas tasas de error sean muy elevadas.

Pista: Las tasas de error en MNIST en el conjunto de validación para una partición 90 % entrenamiento - 10 % validación (semilla aleatoria para el barajado número 23) con proyección a 20 dimensiones de PCA, y 1, 2 y 5 componentes por mixtura en convergencia con suavizado $\alpha = 1e-4$ son 6.17 %, 5.35 % y 4.35 %, respectivamente.

Ejercicio 2.4 (obligatorio: 0.2 puntos). Como hemos comentado anteriormente, en este ejercicio utilizaremos los valores óptimos de los parámetros del clasificador para entrenar y evaluar un clasificador final en los conjuntos oficiales MNIST de entrenamiento y test, respectivamente. Para ello te recomendamos que tomes como punto de partida el script `pca+mixgaussian-exp.m`, modificándolo adecuadamente para generar el script `pca+mixgaussian-eva.m` que también deberá recibir como entrada el conjunto de test de MNIST. Recuerda que toda estimación de (la probabilidad de) error de un clasificador final, debe ir acompañada de sus correspondientes intervalos de confianza al 95 %. Discute los resultados obtenidos comparándolos con los obtenidos con el clasificador gaussiano y con los reportados en la tarea MNIST con clasificadores que conozcas.

3. Máquinas de Vectores Soporte

En esta parte del proyecto se experimenta con SVM mediante la librería **LibSVM**³ desarrollada para Octave. Los ejecutables de esta librería compilados para la versión de Octave instalada en los laboratorios está disponible en PoliformaT en el fichero `svm_apr.tgz`.

Para poder emplear dicha librería, descomprime el fichero `svm_apr.tgz`, que generará el directorio `svm_apr` con los ejecutables `svm-train.mex` y `svm-predict.mex`. Simplemente necesitarás añadir la ruta del directorio `svm_apr` al `PATH` desde Octave mediante el comando:

```
addpath("svm_apr");
```

Si deseas utilizar esta librería en tu propio ordenador, deberás compilar tu propia versión de ambos ejecutable a partir del código fuente.

3.1. Ejemplo de entrenamiento y clasificación con LibSVM

Según la teoría de SVM, dado un conjunto de aprendizaje S , formado por N vectores y sus correspondientes etiquetas de clase, el sistema de aprendizaje obtiene los *multiplicadores de Lagrange*, α_n óptimos asociados a cada vector x_n de S , $1 \leq n \leq N$. Los multiplicadores no nulos definen los vectores de S que constituyen el *soporte* de la función discriminante lineal del clasificador (en dos clases, por el momento). Concretamente, a partir de los multiplicadores no-nulos de los correspondientes vectores soporte y de sus etiquetas de clase, se obtienen fácilmente el vector de pesos θ y el término independiente o *umbral*, θ_0 , que definen la función discriminante lineal del clasificador aprendido.

Por tanto, una vez obtenidos los multiplicadores no nulos, se puede construir de forma trivial un clasificador que prediga la etiqueta de clase asociada a cualquier vector de validación o test.

LibSVM, implementa el proceso de aprendizaje mediante `svmtrain.mex`, y la clasificación mediante `svmpredict.mex`. Un detalle importante de esta librería es que los multiplicadores obtenidos por el proceso de aprendizaje, y los que usan en clasificación, pueden ser positivos o negativos. Se asume que los positivos corresponden a vectores soporte de la clase $+1$ y los negativos, a los de la -1 . A todos los efectos, cualquiera de estos multiplicadores (que el toolkit denota como `sv_coef`), puede considerarse como el producto del verdadero multiplicador de Lagrange, por la etiqueta de clase del vector soporte correspondiente, es decir, $c_n \alpha_n$.

3.1.1. Corpus de ejemplo

Para introducir la funcionalidad básica de la librería **LibSVM** en Octave se propone emplear el corpus de `hart` que se encuentra en el paquete `data.tgz` disponible en PoliformaT. Es un corpus de dos clases donde los datos se representan en dos dimensiones

³<http://www.csie.ntu.edu.tw/~cjlin/libsvm>

(\mathbb{R}^2). Este conjunto de datos *no* es linealmente separable en su espacio de representación original, pero sí lo es en un espacio transformado mediante un *kernel de base radial* (RBF), cuya dimensionalidad es mucho mayor. Los pasos a seguir podrían ser los siguientes.

Primero ejecutamos Octave en el directorio `svm_apr` resultado de descomprimir el paquete `svm_apr.tgz`. Una vez en Octave, cargamos los ficheros con los datos y las etiquetas de clase del conjunto de entrenamiento:

```
octave:1> load data/hart/tr.dat ; load data/hart/trlabels.dat
```

Una vez cargados los datos podemos visualizarlos mediante:

```
octave:2> plot (X(:,1),X(:,2),"x")
```

donde el primer parámetro son los valores de la primer dimensión, el segundo parámetro son los valores de la segunda dimensión, y el tercer parámetro es una cadena de formato que especifica la representación gráfica de los puntos. En el Apéndice D encontrarás más información sobre el comando `plot`.

Si se examinan `trlabels.dat` y `tslabels.dat`, puede observarse que las etiquetas de clase son “1” y “2”. LibSVM admite más de dos clases, que pueden denotarse mediante etiquetas numéricas. En el caso de dos clases, estas etiquetas se convierten internamente en +1 y -1 y esto se refleja en los signos de los multiplicadores (`sv_coef`), como se ha explicado anteriormente. Podemos visualizar separadamente las muestras de cada clase mediante:

```
octave:3> plot(X(xl==1,1),X(xl==1,2),"x",X(xl==2,1),X(xl==2,2),"s")
```

donde los tres primeros parámetros corresponden a la representación gráfica de la clase “1”, y los tres siguientes parámetros, a la de la clase “2”. Se puede observar cómo estas muestras *no* son linealmente separables en su espacio original \mathbb{R}^2 .

3.1.2. Entrenamiento

Como se ha mencionado anteriormente, la función `svmtrain` implementa el proceso de aprendizaje. Se puede obtener una breve ayuda sobre el uso de `svmtrain` invocando la función sin argumentos:

```
octave:4> svmtrain
Usage: model = svmtrain(training_label_vector,
                        training_instance_matrix, 'libsvm_options');
libsvm_options:
[...]
```

Como resultado podemos ver los diferentes parámetros que podemos usar para el aprendizaje: tipo de kernel, parámetro C , grado del kernel si es polinomial, etc. Por ejemplo, para realizar un entrenamiento con kernel tipo RBF (-t 2) con parámetro $C = 1$ (-c 1):

```
octave:5> res = svmtrain(xl, X, '-t 2 -c 1');
.*.*
optimization finished, #iter = 2298
nu = 0.174579
obj = -108.403744, rho = -0.096434
nSV = 398, nBSV = 91
Total nSV = 398
```

El resultado del proceso de entrenamiento se ha almacenado en la variable **res**, que es un tipo estructurado que contiene, entre otros: los parámetros del modelo, los índices de los vectores que han resultado ser vectores soporte, el multiplicador de Lagrange (multiplicado por la correspondiente etiqueta de clase, +1 o -1) asociado a cada vector soporte, etc. El contenido completo de **res** puede visualizarse mediante:

```
octave:6> res
```

Se puede acceder a cada uno de los campos de **res** mediante el operador “.”. Por ejemplo se pueden mostrar los índices de los vectores soporte mediante:

```
octave:7> res.sv_indices
```

Y para mostrar los multiplicadores de Lagrange (multiplicados por las correspondientes etiquetas de clase, +1 o -1, en el caso de 2 clases):

```
octave:8> res.sv_coef
```

Ejercicio 3.1 (opcional: no entregable). Para la tarea **hart** de este ejemplo, representa gráficamente los vectores soporte obtenidos mediante la función **svmtrain**, superponiéndolos a las muestras de entrenamiento.

3.1.3. Clasificación y estimación de la tasa de acierto

A partir del modelo entrenado en la etapa anterior almacenado en la variable **res**, se pueden clasificar los vectores de un conjunto de validación o test. Primeramente, ese necesario cargar este conjunto:

```
octave:9> load data/hart/ts.dat ; load data/hart/tslabels.dat
```

La clasificación se realiza mediante la función **svmpredict**. Como en el caso de **svmtrain**, para ver todas las opciones de **svmpredict** basta invocar **svmpredict** sin argumentos. Tanto en **svmtrain** como en **svmpredict** podremos poner entre comillas simples cualquiera de las opciones que admite la librería LibSVM.

```
octave:10> svmpredict(yl,Y, res,' ');
Accuracy = 98.1% (981/1000) (classification)
```

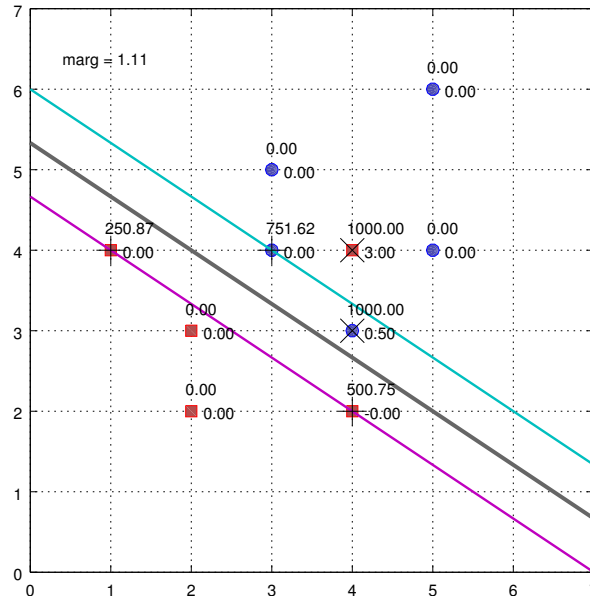
Ejercicio 3.2 (obligatorio: 0.4 puntos). En el subdirectorio `data/mini` se encuentran dos pequeños conjuntos de datos de entrenamiento en dos dimensiones: (`trSep.dat`, `trSeplabels.dat`) y (`tr.dat`, `trlabls.dat`). El primero es linealmente separable (no es necesario *kernel*) y el segundo no. Para cada uno de estos conjuntos⁴:

- Obtén los SVM sin kernel (es decir, kernel tipo lineal). Para simular la optimización estándar del caso separable basta usar un valor grande de \mathcal{C} ($\mathcal{C} = 1000$).
- Determina a) los multiplicadores de Lagrange, α , asociados a cada dato de entrenamiento, b) los vectores soporte, c) el vector de pesos y umbral de la función discriminante lineal, y d) el margen correspondiente.
- Calcula los parámetros de la frontera lineal (recta) de separación;
- Representa gráficamente los vectores de entrenamiento, marcando los que son vectores soporte, y la recta separadora correspondiente.

Además, para el conjunto no-separable utilizando diversos valores relevantes de \mathcal{C} :

- Determina los valores de tolerancia de margen, ζ , asociados a cada dato de entrenamiento.
- Marca los vectores soporte “erróneos” en la representación gráfica.

Pista: Un ejemplo “ideal” de representación gráfica para el caso no-separable, obtenido con $\mathcal{C} = 1000$, se muestra en la siguiente figura:



⁴Te recomendamos que consultes el Apéndice E para un breve recordatorio sobre la teoría de SVM.

3.2. Tarea MNIST

Como se puede observar en la web de MNIST la aplicación de SVM a MNIST proporciona mejores resultados que la combinación PCA con un clasificador cuadrático.

Ejercicio 3.3 (obligatorio: 0.4 puntos). Realiza un experimento donde se evalúe el error de clasificación en función de los parámetros del clasificador basado en SVM. Más concretamente, explora los valores del parámetro C (`-c` 1, 10, 100...) y el tipo de kernel (`-t` 0, 1, 2, 3). Para aquellos tipos de kernel que lo permitan, explora sus parámetros específicos. Por ejemplo, en el caso del kernel polinomial (`-t` 1), explora el grado del polinomio (`-d` 1, 2, 3, 4, 5).

En función del número de resultados obtenidos como consecuencia de la exploración de los valores de los parámetros, utiliza una representación adecuada de los mismos, ya sea gráfica o tabular, que muestre no solo el mejor resultado obtenido, sino también otros resultados relevantes que permitan poner de manifiesto la tendencia a mejorar o empeorar del modelo según varían los valores de los parámetros considerados.

Al igual que en el ejercicio 2.3, te recomendamos que elabores un script `svm-exp.m` a partir del script `pca+mixgaussian-exp.m`, que realice la exploración de parámetros descrita utilizando la misma partición entrenamiento-validación que en los experimentos de mixtura de gaussianas.

Pista: Se recomienda empezar por el kernel polinomial (`-t` 1) con un valor bajo de C (i.e., `-c` 1) y explorar de manera creciente los grados del polinomio en el rango recomendado. En función de los resultados obtenidos, se pueden probar más valores del grado del polinomio o pasar a aumentar el valor de C . Después, en función de la capacidad de cómputo y del tiempo disponible, se pueden estudiar otros tipos de kernel.

Ejercicio 3.4 (obligatorio: 0.2 puntos). Una vez determinado los valores óptimos de los parámetros del clasificador basado en SVM, entrena y evalúa un clasificador final en los conjuntos oficiales MNIST de entrenamiento y test, respectivamente. Para ello te recomendamos que tomes como punto de partida el script `svm-exp.m`, modificándolo adecuadamente para generar el script `svm-eva.m` que también deberá recibir como entrada el conjunto de test de MNIST. Recuerda que toda estimación de (la probabilidad de) error de un clasificador final, debe ir acompañada de sus correspondientes intervalos de confianza al 95 %. Discute los resultados obtenidos comparándolos con los obtenidos con el clasificador de mixtura de gaussianas y con los reportados en la tarea MNIST, especialmente los basados en SVM.

4. Redes Neuronales Multicapa

En esta última parte del proyecto trabajaremos con redes neuronales multicapa dentro de Octave, en concreto con redes *Multilayer Perceptron* (MLP) implementadas mediante la librería **nnet**⁵.

Se ha dejado disponible en PoliformaT el fichero **nnet_apr.tgz** que contiene el directorio **nnet_apr** con todos los ficheros de la librería de redes neuronales. La instalación es muy sencilla. Simplemente se necesita añadir la ruta de la librería **nnet** al **PATH** desde Octave mediante el comando:

```
addpath("nnet_apr");
```

Nótese que la ruta anterior es relativa al directorio de trabajo donde ejecutas los scripts Octave que desarrolles.

4.1. Ejemplo de entrenamiento y clasificación con nnet

Para ilustrar el funcionamiento de la librería **nnet** utilizaremos el corpus **hart** que ya hemos visto anteriormente. Seguidamente, definiremos los conjuntos de entrenamiento, validación y test, y su preproceso para el entrenamiento, clasificación y evaluación.

Al igual que en ejercicios previos, te recomendamos que elabores un script **mlp-exp.m** a partir del script **mixgaussian-exp.m**, que realizará la exploración de parámetros de entrenamiento. Recuerda añadir la librería **nnet** al **PATH** mediante la función **addpath** al principio de tu script Octave.

4.1.1. Definición de entrenamiento y validación

El entrenamiento de **nnet** tiene dos posibles criterios de parada: entrenar hasta un número máximo de *epochs* o entrenar hasta que el error de clasificación de un conjunto de validación no mejore. Normalmente se suele utilizar el segundo criterio, por lo que debemos dedicar una parte del conjunto de entrenamiento a conjunto de validación.

Cómo hemos hecho en experimentos previos, tras la lectura de parámetros de entrada y la carga de datos, barajamos los datos y hacemos una partición en entrenamiento (**trper**=90) y validación (**dvper**=10):

```
N=rows(X);
seed=23; rand("seed",seed); permutation=randperm(N);
X=X(permutation,:); xl=xl(permutation,:);

Ntr=round(trper/100*N);
Ndv=round(dvper/100*N);
Xtr=X(1:Ntr,:); xltr=xl(1:Ntr);
Xdv=X(N-Ndv+1:N,:); xldv=xl(N-Ndv+1:N);
Y=Xdv; yl=xldv;
```

⁵<http://octave.sourceforge.net/nnet>

4.1.2. Preproceso de los datos

Como sugerencia de diseño se recomienda encapsular las siguientes instrucciones dentro de una función `mlp` que reciba como entrada los conjuntos de datos (entrenamiento, validación y test) junto con los parámetros relacionados con el entrenamiento de la red neuronal, y como salida devuelva el error de clasificación calculado en el conjunto de test:

```
function [errY] = mlp(Xtr,xltr,Xdv,xldv,Y,y1,nHidden,epochs,show,seed)
```

Esta función se invocará desde el script `mlp-exp.m`.

Primeramente, las funciones de la librería `nnet` requieren que los datos y las etiquetas de clase estén por columnas en lugar de por filas, por lo que es necesario trasponer los datos:

```
Xtr = Xtr'; xltr=xltr'; Xdv=Xdv'; xldv=xldv'; Y=Y'; y1=y1';
```

Seguidamente, como estamos trabajando con redes neuronales, es necesario codificar cada una de las etiquetas de clase en su formato *one-hot encoding*. Dado un problema de clasificación en C clases, la capa de salida de la red neuronal suele tener C neuronas de forma que la neurona c -ésima se activa para la clase c . Por ejemplo, en un problema de clasificación en 3 clases, éstas clases se codificarían como $[0,0,1]$, $[0,1,0]$ y $[1,0,0]$.

Ejercicio 4.1. Implementa una función `xloh= onehot(xl)` que realice la codificación *one-hot* de un vector de etiquetas de clase, devolviendo la correspondiente matriz donde cada columna es el vector *one-hot* de cada etiqueta de clase.

Seguidamente, es conveniente normalizar los datos de modo que tengan media cero y desviación típica uno. Esto se consigue con la función `prestd`:

```
[Xtrnorm,Xtrmean,Xtrstd] = prestd(Xtr);
```

tal que `Xtrnorm` son los datos normalizados, y `Xtrmean` y `Xtrstd` son la media y desviación típica de los datos originales, respectivamente.

Por último, el conjunto de validación debe representarse en una estructura especial con dos campos `P` y `T`. El campo `P` para los datos de validación normalizados mediante la función `trastd`:

```
XdvNN.P = trastd(Xdv,Xtrmean,Xtrstd);
```

y el campo `T` para la salida, es decir, las etiquetas de clase del conjunto de validación en codificación *one-hot*:

```
XdvNN.T = onehot(xldv);
```

4.1.3. Entrenamiento

Antes entrenar una red neuronal debemos especificarla mediante la siguiente función:

```
initNN = newff (Pr,ss,trf,btf,blf,pf)
```

siendo

- **Pr**: una matriz $D \times 2$ con los valores máximo y mínimo de los datos de entrenamiento en cada dimensión. Se puede calcular mediante la función `minmax` disponible en la librería `nnet`.
- **ss**: un vector fila con el número de neuronas en cada capa oculta y en la capa de salida.
- **trf**: la lista de funciones de activación de cada capa.
- **btf**: el algoritmo de entrenamiento de la red neuronal. En la versión actual, el único algoritmo implementado es `trainlm` que es el algoritmo *backpropagation*.
- **blf**: un parámetro que no se utiliza en la versión actual.
- **pf**: la función objetivo a minimizar. En la versión actual, la única función objetivo es `mse` que es el error cuadrático medio.

Una posible configuración para una topología de una capa oculta de `nHidden` neuronas y con `nOutput` neuronas en la capa de salida, y funciones de activación `tansig` $\in [-1, 1]$ y `logsig` $\in [0, 1]$ para la capa oculta y de salida, respectivamente, la crearíamos mediante:

```
initNN = newff(minmax(Xtrnorm),[nHidden nOutput],
               {"tansig","logsig"},"trainlm","", "mse");
```

Ten en cuenta que `nOutput` debe ser igual al número de clases. Con la red neuronal inicializada en `initNN` podemos establecer algunos parámetros de entrenamiento:

```
initNN.trainParam.show = show;
initNN.trainParam.epochs = epochs;
```

donde el parámetro `show` indica cada cuántas *epochs* de entrenamiento queremos que se imprima información, y el parámetro `epochs` es el número máximo de epochs que queremos que realice, adicionalmente al criterio de parada del conjunto de validación. Puedes definir el valor de `show` a 10 y `epochs` a 300 en tu script `mlp-exp.m`.

Finalmente, podemos entrenar la red mediante:

```
rand("seed",seed);
NN = train(initNN,Xtrnorm,onehot(xltr),[],[],XdvNN);
```

De esta forma obtenemos en `NN` la red entrenada. La razón de incluir la instrucción `rand("seed",seed)` antes del entrenamiento es fijar la semilla del generador de números aleatorios para que la inicialización de la red neuronal sea la misma y los experimentos sean reproducibles. Si ejecutamos el script `mlp-exp.m` que estamos desarrollando para que entrene una red neuronal con `nHidden=5` neuronas en la capa oculta, y 90 % entrenamiento y 10 % validación sobre el conjunto de entrenamiento de la tarea `hart`:

```
./mlp-exp.m data/hart/tr.dat data/hart/trlabels.dat "[5]" 90 10
```

obtendremos la siguiente salida:

```
TRAINLM, Epoch 0/300, MSE 0.296476/0, Gradient 119.882/1e-10
TRAINLM, Epoch 10/300, MSE 0.0545277/0, Gradient 20.566/1e-10
TRAINLM, Epoch 20/300, MSE 0.0178138/0, Gradient 7.4839/1e-10
TRAINLM, Epoch 21/300, MSE 0.0176538/0, Gradient 2.29594/1e-10
TRAINLM, Validation stop.
```

Además se mostrará una gráfica con la evolución del *Mean Square Error* (MSE) en entrenamiento y validación en función del número de *epochs*. Para poder apreciarla, es necesario añadir al final del script la llamada `pause(10)` para pausar la ejecución del script 10 segundos.

4.1.4. Clasificación y estimación de la tasa de error

A continuación podemos calcular la clasificación de las muestras de test normalizadas a partir de la función `sim` que devuelve una matriz donde para cada muestra (columnas) proporciona la distribución de probabilidad a posteriori sobre las clases (filas):

```
Ynorm = trastd(Y,Xtrmean,Xtrstd);
Yout = sim(NN,Ynorm);
```

siendo el valor de la variable `Yout` para las 5 primeras muestras de validación (`Yout(:,1:5)`):

```
ans =
```

```
    9.9899e-01    1.3711e-02    9.8884e-01    9.2127e-01    1.4303e-02
    7.6718e-04    9.8627e-01    8.6588e-03    7.8947e-02    9.8406e-01
```

Por tanto, para realizar la estimación de la clase por cada muestra de test, debemos obtener la etiqueta estimada como la clase con mayor puntuación. A continuación, comparar esta etiqueta estimada con la etiqueta real y calcular el error de clasificación correspondiente.

Ejercicio 4.2. Implementa en la función `mlp` la obtención de la clase estimada a partir de la salida de la función `sim` y la estimación del error de clasificación `errY` que devuelve.

Ejercicio 4.3 Como has hecho con otros parámetros de clasificadores anteriormente, el script `mlp-exp.m` debe poder evaluar un vector de valores de parámetros, en este caso será un número variables de neuronas en la capa oculta. Por ejemplo:

```
./mlp-exp.m data/hart/tr.dat data/hart/trlabels.dat "[1 2 5]" 90 10
```

Acaba de implementar el script `mlp-exp.m` y comprueba su correcto funcionamiento.

Pista: Las tasas de error en el conjunto de validación para `hart` particionando el fichero de entrenamiento en 90 % entrenamiento y 10 % validación con semilla aleatoria el número 23 y utilizando 1, 2 y 5 neuronas en la capa oculta con función `tansig`, y en la capa de salida con función `logsig` son 13 %, 11 % y 5 %, respectivamente.

Ejercicio 4.4 Habrás observado que el script `mlp-exp.m` utiliza el conjunto de validación tanto para el criterio de parada durante el entrenamiento, como test para estimar la tasa de error. Esto es así porque estamos optimizando los valores de los parámetros de la red neuronal.

Sin embargo, una vez hayamos estimado los valores de los parámetros de la red neuronal en el conjunto de validación, la estimación del error del clasificador final se realizará en el conjunto de test. Para ello, como hemos hecho en clasificadores estudiados previamente, implementa un script `mlp-eva.m`. Asumiendo que el número óptimo de neuronas en la capa oculta es 5, la ejecución del script `mlp-eva.m` sería:

```
./mlp-eva.m data/hart/tr.dat data/hart/trlabels.dat data/hart/ts.dat
data/hart/tslabels.dat 5 90 10
```

Pista: La tasa de error en el conjunto de test para `hart` particionando el fichero de entrenamiento en 90 % entrenamiento y 10 % validación con semilla aleatoria el número 23 y utilizando 5 neuronas en la capa oculta con función `tansig`, y en la capa de salida con función `logsig` es 3.90 %.

4.2. Tarea MNIST

Como se puede observar en la web de MNIST la utilización de redes neuronales puede llegar a proporcionar los mejores resultados en esta tarea. Sin embargo, por cuestiones de eficiencia, la librería `nnet` no permite entrenar redes neuronales complejas con un elevado número de muestras cuyos vectores de características de entrada tengan un número elevado de dimensiones, o una configuración de red con un elevado número de neuronas en la capa oculta, o incluso varias capas ocultas. Por ello, será necesario la proyección previa de los datos con PCA a una dimensionalidad reducida que permita realizar el entrenamiento de la red.

Ejercicio 4.5 (obligatorio: 0.75 puntos). Realiza un experimento donde se evalúe el error de clasificación en función del número de neuronas en la capa oculta de la red neuronal (`nHidden=1, 2, 5, 10, 20, 30, 40, 50...`), para un número de dimensiones PCA variable ($D = 1, 2, 5, 10, 20, 30$). También puedes utilizar un subconjunto del conjunto de entrenamiento si tienes limitaciones de uso de memoria o si el tiempo de cómputo consideras que es excesivo (i.e. en lugar del 90 % usa un 40 %).

Para ello, puedes partir de los scripts `pca+mixgaussian-exp.m` y `mlp-exp.m` para crear el script `pca+mlp-exp.m` que mantenga la misma partición en entrenamiento y validación, y realice al menos la exploración de parámetros propuesta.

Recuerda que el objetivo de este experimento es determinar los valores de los parámetros del clasificador (n_{Hidden} y D) que minimizan el error de clasificación en el conjunto de validación.

Representa gráficamente las tasas de error obtenidas en el conjunto de validación de forma similar a la utilizada en mixturas de gaussianas. Es decir, cada gráfica tendrá tantas curvas como valores de PCA diferentes hayas probado y cada curva mostrará la evolución del error de clasificación (eje y) en función del número de neuronas en la capa oculta (eje x). Para mejorar la legibilidad de la representación gráfica de los resultados puedes descartar aquellas curvas asociadas a valores de PCA cuyas tasas de error sean muy elevadas.

Pista: Las tasas de error en MNIST en el conjunto de validación para una partición 40 % entrenamiento y 10 % validación (semilla aleatoria número 23) con proyección a 20 dimensiones de PCA, y 10, 20 y 30 neuronas en la capa oculta son aproximadamente 11.15 %, 6.82 % y 5.32 %, respectivamente.

Ejercicio 4.6 (obligatorio: 0.25 puntos). Tras el ajuste de parámetros en el conjunto de validación, en este ejercicio utilizaremos los valores óptimos de los parámetros del clasificador para entrenar y evaluar un clasificador final en los conjuntos oficiales MNIST de entrenamiento y test, respectivamente. Para ello te recomendamos que tomes como punto de partida el script `pca+mlp-exp.m`, modificándolo adecuadamente para generar el script `pca+mlp-eva.m` que también deberá recibir como entrada el conjunto de test de MNIST. Recuerda que toda estimación de (la probabilidad de) error de un clasificador final, debe ir acompañada de sus correspondientes intervalos de confianza al 95 %. Discute los resultados obtenidos comparándolos con los obtenidos en los clasificadores estudiados y con otros clasificadores basados en redes neuronales reportados en la tarea MNIST.

A. Tutorial sobre Octave

A.1. Introducción

Varias de las técnicas empleadas en Aprendizaje Automático y Reconocimiento de Formas (RF) emplean cálculos vectoriales y matriciales, como por ejemplo en las implementaciones de clasificadores basados en la distribución gaussiana. Por tanto, una herramienta que implemente de forma sencilla estos cálculos matriciales puede simplificar notablemente la implementación de sistemas de RF.

Una de las herramientas comerciales más potentes en cálculo matricial es MATLAB. Asimismo existe una herramienta de código libre que presenta capacidades semejantes: *GNU Octave*.

GNU Octave es un lenguaje de alto nivel interpretado definido inicialmente para computación numérica. Entre otras, posee capacidades de cálculo numérico para solucionar problemas lineales y no lineales. También dispone de herramientas gráficas para visualizar datos y resultados. Se puede usar de forma interactiva y/o programada mediante *scripts* en un lenguaje interpretado. La sintaxis y semántica de Octave es prácticamente idéntica a MATLAB, lo que hace que los programas sean fácilmente portables entre ambas plataformas.

Octave está en continuo crecimiento y puede descargarse y consultarse su documentación y estado en su web <http://www.gnu.org/software/octave/>. Aunque está definido para funcionar en GNU/Linux, también es portable a otras plataformas como OS X y MS-Windows (los detalles pueden consultarse en la web mencionada).

A.2. Órdenes básicas de *Octave*

Octave puede ejecutarse desde la línea de órdenes o desde el menú de aplicaciones del entorno gráfico. Para ejecutar desde la línea de órdenes se abre un terminal y se escribe:

```
octave
```

Generalmente, obtenemos una salida similar a:

```
GNU Octave, version 3.8.2
Copyright (C) 2014 John W. Eaton and others.
This is free software; see the source code for copying conditions.
There is ABSOLUTELY NO WARRANTY; not even for MERCHANTABILITY or
FITNESS FOR A PARTICULAR PURPOSE. For details, type 'warranty'.

Octave was configured for "x86_64-redhat-linux-gnu".

Additional information about Octave is available at http://www.octave.org.

Please contribute if you find this software useful.
```

For more information, visit <http://www.octave.org/get-involved.html>

Read <http://www.octave.org/bugs.html> to learn how to submit bug reports. For information about changes from previous versions, type 'news'.

```
octave:1>
```

La última línea tendrá un cursor indicando que se esperan órdenes de Octave. Estamos pues en el modo **interactivo**.

A.2.1. Aritmética básica

Octave acepta a partir de este momento expresiones aritméticas sencillas (operadores `+`, `-`, `*`, `/` y `^`, este último para exponenciación), funciones trigonométricas (`sin`, `cos`, `tan`, `arcsin`, `arccos`, `arctan`), logaritmos (`log`, `log10`), exponencial neperiana (`e^n` ó `exp(n)`) y valor absoluto (`abs`). Como respuesta a estas expresiones Octave da valor a la variable predefinida `ans`, y la muestra. Pero los resultados también pueden asignarse a otras variables. Por ejemplo:

```
octave:1> sin(1.71)
ans =  0.99033
octave:2> b=sin(2.16)
b =  0.83138
```

Para consultar el valor de una variable, basta con escribir su nombre, aunque también puede emplearse la función `disp`, que muestra el contenido de la variable omitiendo su nombre:

```
octave:3> b
b =  0.83138

octave:4> ans
ans =  0.99033

octave:5> disp(b)
0.83138
```

Las variables pueden usarse en otras expresiones:

```
octave:6> c=b*ans
c =  0.82334
```

Puede evitarse que se muestre el resultado de cada operación añadiendo `(;)` al final de la operación:

```
octave:7> d=ans*b*5;
octave:8> disp(d)
4.1167
```

A.2.2. Operadores básicos en vectores y matrices

Para la notación matricial en Octave se usan los corchetes (`[]`); en su interior, las filas se separan por punto y coma (`;`) y las columnas por espacios en blanco () o por comas (`,`). Por ejemplo, para crear un vector fila de dimensión 3, un vector columna de dimensión 2 y una matriz de 3×2 , se puede hacer:

```
octave:9> v1=[1 3 -5]
```

```
v1 =
```

```
1    3   -5
```

```
octave:10> v2=[4;2]
```

```
v2 =
```

```
4
```

```
2
```

```
octave:11> m=[3,-4;2 1;-5 0]
```

```
m =
```

```
3   -4
```

```
2    1
```

```
-5    0
```

Siempre y cuando las dimensiones de los elementos vectoriales y matriciales implicados sean apropiados, sobre ellos se pueden aplicar operadores de suma (+), diferencia (-) o producto (*). El operador potencia (^) puede aplicarse sobre matrices cuadradas. Los operadores producto, división y potencia tienen la versión “elemento a elemento” (`.*`, `./`, `.^`). Por ejemplo:

```
octave:12> mv=v1*m
```

```
mv =
```

```
34   -1
```

```
octave:13> mx=m.*5
```

```
mx =
```

```
15  -20
```

```
10    5
```

```
-25    0
```

```
octave:14> v3=v1*v2
```

```
error: operator *: nonconformant arguments (op1 is 1x3, op2 is 2x1)
```

```
octave:15> v3=v1*[3;5;6]
```

```
v3 = -12
```

```
octave:16> mvv=[3;5;6]*v1
```



```
mvv =  
    3    9   -15  
    5   15   -25  
    6   18   -30
```

Los operadores de comparación (>, <, >=, <=, ==, !=) se pueden aplicar “elemento a elemento”. Como resultado se obtiene una matriz binaria con 1 en las posiciones en las que se cumple la condición y 0 en caso contrario:

```
octave:17> m>=0  
ans =  
    1    0  
    1    1  
    0    1
```

```
octave:18> m!=0  
ans =  
    1    1  
    1    1  
    1    0
```

Esos operadores se pueden emplear en la comparación de vectores y matrices de dimensiones congruentes. La matriz binaria resultante contiene los resultados de las comparaciones de los pares de elementos en la misma posición en ambas estructuras.

Para tranponer una matriz o vector se usa el operador de transposición (’):

```
octave:19> m2=m’  
m2 =  
    3    2   -5  
   -4    1    0
```

El indexado de los elementos se hace entre paréntesis. Para vectores puede indicarse una posición o lista de posiciones, mientras que para una matriz se espera una fila o secuencia de filas seguida de una columna o secuencia de columnas:

```
octave:20> v1(2)  
ans = 3
```

```
octave:21> v1([2 3])  
ans =  
    3   -5
```

```
octave:22> v2(2)  
ans = 2
```

```
octave:23> m3=[1 2 3 4;5 6 7 8;9 10 11 12]
m3 =
     1     2     3     4
     5     6     7     8
     9    10    11    12
```

```
octave:24> m3([1 3], [1 4])
ans =
     1     4
     9    12
```

Para indicar todas las filas o columnas, se puede emplear (:):

```
octave:25> m3(:, [1 3])
ans =
     1     3
     5     7
     9    11
```

Los rangos se denotan como (*i:f*), donde *i* es el índice inicial y *f* el final. Se puede emplear la notación (*i:inc:f*), donde *inc* indica el incremento, que por omisión es 1.

```
octave:26> m3([1 3], 1:3)
ans =
     1     2     3
     9    10    11
```

```
octave:27> m3([1 3], 1:2:4)
ans =
     1     3
     9    11
```

Para indicar el último índice de una dimensión se puede emplear (**end**):

```
octave:28> m3([1 3], end)
ans =
     4
    12
```

A.2.3. Funciones básicas en vectores y matrices

Octave aporta múltiples funciones para operar con vectores y matrices. Las más importantes son:

- **size(m)**: devuelve número de filas y columnas de la matriz (en el caso de un vector, una de las dimensiones tendrá tamaño 1)

- `eye(f,c)`, `ones(f,c)`, `zeros(f,c)`: dan la matriz identidad, todo unos y nula, respectivamente, de tamaño $f \times c$; si se pone un solo número, da la matriz cuadrada correspondiente
- `sum(v)`, `sum(m)`: da la suma de los elementos del vector o matriz; en el caso de la matriz, devuelve el vector resultante de las sumas por columnas; si se le pasa un segundo argumento (`sum(m,n)`), éste indica la dimensión a sumar (1 para columnas, 2 para filas).
- `max(v)`, `max(m)`: indica el valor máximo del vector o el vector con los máximos por columna de la matriz; si se pide que devuelva dos resultados (`[r1,r2]=max(v)`, `[r1,r2]=max(m)`), el primer resultado almacena los valores y el segundo su posición
- `det(m)`: determinante de `m`
- `eig(m)`: vector de valores propios de `m` o su versión matricial diagonal
- `diag(v)`: crear matriz diagonal con los valores de `v`
- `inv(m)`: inversa de la matriz `m` si esta es no singular
- `trace(m)`: traza de la matriz `m`
- `sort(v)`: vector ordenado con los valores del vector `v`
- `repmat(m,f,c)`: crea una matriz de $f \times c$ bloques de `m`; si `c` se omite, será de $f \times f$
- `find(v)`, `find(m)`: se le pasa un vector o matriz e indica aquellos elementos que no son cero (índices absolutos empezando en 1 y haciendo el recorrido por cada columna y por filas ascendentes); si se le piden dos resultados (`[r1,r2]=find(v)`, `[r1,r2]=find(m)`) el primer resultado almacena fila y el segundo columna; se puede aplicar sobre resultados de operaciones lógicas a fin de verificar elementos de la matriz o vector que cumplen una condición. Por ejemplo, mediante la función (`rem`) que obtiene el resto del primer operador dividido por el segundo, se pueden obtener las posiciones de los elementos pares:

```
octave:29> [r,c]=find(rem(m3,2)==0)
r =
    1
    2
    3
    1
    2
    3

c =
    2
```

```
2
2
4
4
4
```

A.2.4. Carga y salvado de datos

La introducción de datos de forma manual no es apropiada para grandes cantidades de datos. Por tanto, Octave aporta funciones que permiten cargar de y salvar en ficheros. El salvado se hace mediante la orden `save`:

```
octave:30> save "m3.dat" m3
```

El fichero tiene el siguiente contenido:

```
# Created by Octave 3.0.5, DATE <user@machine>
# name: m3
# type: matrix
# rows: 3
# columns: 4
1 2 3 4
5 6 7 8
9 10 11 12
```

Los ficheros de datos a cargar deben seguir este formato, indicando en la línea “# name:” el nombre de la variable en la que se cargarán los datos. Por ejemplo, ante un fichero `maux.dat` cuyo contenido es:

```
# Created by Octave 3.0.5, DATE <user@machine>
# name: A
# type: matrix
# rows: 4
# columns: 3
1 2 -3
5 -6 7
-9 10 11
-5 2 -1
```

Su carga definirá la matriz (**A**) como:

```
octave:31> load "maux.dat"
```

```
octave:32> A
A =
```

```
1    2   -3
5   -6    7
-9   10   11
-5    2   -1
```

La orden `save` puede usarse con opciones como `-text` (grabar en formato texto con cabecera, por omisión), `-ascii` (graba en formato texto sin cabecera), `-z` (graba en formato comprimido), o `-binary` (graba en binario). Por ejemplo:

```
octave:33> save -ascii "m3woh.dat" m3
```

En ocasiones, el salvado de datos puede provocar pérdida de precisión, pues por omisión se salva hasta el cuarto decimal. Para modificar esta precisión de salvado, se puede emplear `save_precision(n)`, donde `n` es el número de posiciones decimales (incluyendo el `.`) que se grabarán.

A.2.5. Funciones Octave

En Octave se pueden definir funciones que hagan tareas específicas y/o complejas. Las funciones Octave pueden recibir varios parámetros y pueden devolver varios valores de retorno (que pueden incluir vectores y matrices). La sintaxis básica es:

```
function [ lista_valores_retorno ] = nombre ( [ lista_parametros ] )
    cuerpo
end
```

Por ejemplo:

```
octave:34> function [m1,m2] = addsub(ma,mb)
> m1=ma+mb
> m2=ma-mb
> end
```

```
octave:35> mat1=[1,2;3,4]
mat1 =
    1    2
    3    4
```

```
octave:36> mat2=[-1,2;3,-4]
mat2 =
   -1    2
    3   -4
```

```
octave:37> addsub(mat1,mat2)
m1 =
    0    4
```

```
6 0

m2 =
2 0
0 8

ans =
0 4
6 0

octave:38> [mr1,mr2]=addsub(mat1,mat2)
mr1 =
0 4
6 0

m2 =
2 0
0 8

mr1 =
0 4
6 0

mr2 =
2 0
0 8
```

Es habitual definir las funciones en ficheros de código Octave cuyo nombre debe ser el mismo que la función con el sufijo “.m” (en nuestro ejemplo sería `addsub.m`). Estos ficheros deben situarse en el mismo directorio en el que se ejecuta Octave. De esa forma, se puede acceder a las funciones sin tener que definir las cada vez.

A.2.6. Programas Octave

Octave se puede usar de forma *no interactiva* mediante *scripts* que son interpretados por Octave. En estos *scripts* o “*programas Octave*” se pueden emplear las mismas instrucciones que en el modo interactivo. Por ejemplo, suponiendo que tenemos en el directorio actual el fichero `addsub.m` con la función previamente definida, desde cualquier terminal podemos crear (con algún editor) el fichero `test.m` con el siguiente contenido:

```
#!/usr/bin/octave -qf
a=[1,2,3;4,5,6;7,8,9]
b=[9,8,7;6,5,4;3,2,1]
c=a+b
```

```
[d,e]=addsub(a,b)
disp(c)
```

Si desde la línea de órdenes le damos permisos de ejecución (`chmod +x test.m`), podremos ejecutarlo como cualquier programa ejecutable o *shell script*:

```
$_ ./test.m
a =
    1    2    3
    4    5    6
    7    8    9

b =
    9    8    7
    6    5    4
    3    2    1

c =
   10   10   10
   10   10   10
   10   10   10

m1 =
   10   10   10
   10   10   10
   10   10   10

m2 =
   -8   -6   -4
   -2    0    2
    4    6    8

d =
   10   10   10
   10   10   10
   10   10   10

e =
   -8   -6   -4
   -2    0    2
    4    6    8

   10   10   10
   10   10   10
```

10 10 10

Estos programas también pueden ejecutarse desde la línea interactiva de Octave escribiendo su nombre (sin el sufijo “.m”). En este caso, se pueden poner argumentos en la línea de órdenes y puede usarse la variable `nargin` (número de argumentos) y la función `argv()` (da una lista de los valores alfanuméricos de los argumentos recibidos). Estos argumentos están en formato de cadena de caracteres; por tanto los valores numéricos deben convertirse a formato numérico, por ejemplo empleando la función `str2num(c)`.

A.3. Ejercicios propuestos

1. Realiza el producto escalar de los vectores $v_1 = (1, 3, 8, 9)$ y $v_2 = (-1, 8, 2, -3)$.
2.
 - a) Obtén la matriz de dimensión 4×4 mediante producto de los vectores v_1 y v_2
 - b) Calcula su determinante
 - c) Calcula sus valores propios.
3. Sobre la matriz del ejercicio 2:
 - a) Calcula su submatriz 2×2 formadas por las filas 1 y 3 y las columnas 2 y 3.
 - b) Súmale la matriz todo unos a la submatriz resultante.
 - c) Calcula el determinante de la matriz resultante.
 - d) Calcula la inversa de la matriz resultante.
4. Sobre la matriz inversa resultado del ejercicio 3:
 - a) Calcula su máximo valor y su posición.
 - b) Calcula las posiciones (fila y columna) de los elementos mayores que 0.
 - c) Calcula la suma de cada columna.
 - d) Calcula la suma de cada fila.
5. Salva la matriz inversa resultante del ejercicio 3 con hasta 7 dígitos decimales; comprueba que se ha salvado correctamente.
6. Define una función Octave que reciba una matriz y devuelva la primera fila y la primera columna de esa matriz.
7. Implementa un *script* Octave que lea una matriz de un fichero `data` y grabe su traspuesta en el fichero `data_trans`.

B. Tarea de clasificación: MNIST

B.1. Introducción

La base de datos MNIST⁶ consiste en una colección de imágenes de dígitos manuscritos (10 clases) con unas dimensiones de 28 x 28 píxeles en escala de 256 niveles de grises. Las dígitos que aparecen en las imágenes han sido normalizados en tamaño (20 x 20 píxeles) y centrados. Esta base de datos es un subconjunto de una base de datos más grande disponible desde el *National Institute of Standards and Technology* (NIST). Ha sido particionada en 60.000 imágenes de entrenamiento y 10.000 de test, que corresponde con los siguientes cuatro ficheros en formato Octave⁷:

- Imágenes de entrenamiento: 60000 x 784 (`train-images-idx3-ubyte.mat.gz`)
- Etiquetas de clase de entrenamiento: 60000 x 1 (`train-labels-idx1-ubyte.gz`)
- Imágenes de test: 10000 x 784 (`t10k-images-idx3-ubyte.mat.gz`)
- Etiquetas de clase de test: 10000 x 1 (`t10k-labels-idx1-ubyte.mat.gz`)

En la Figura 1 se muestra una representación de un dígito cero que se corresponde con la segunda fila de datos del fichero `train-images-idx3-ubyte.mat.gz` que ha sido formateada a 28 filas y 28 columnas. Se puede apreciar como el tamaño del dígito está normalizado a 20 x 20 píxeles centrado sobre un fondo blanco.

La tarea MNIST ha sido cuidadosamente elaborada para que el conjunto de escritores de entrenamiento y test sea disjunto. De esta forma, no hay dígitos del mismo escritor en el entrenamiento y test. Asimismo, existen dos tipos de escritores que conviven en el conjunto de entrenamiento y en el test, que se corresponde con estudiantes de instituto y trabajadores de la oficina del censo, respectivamente. Estos últimos poseían una escritura más clara y fácil de reconocer.

En la web de la base de datos MNIST se proporcionan muchos más detalles sobre su elaboración. Asimismo, en esta misma página web se muestra una tabla de clasificadores (y preproceso aplicado) con la tasa de error conseguida sobre esta tarea y la referencia en forma de enlace a una descripción más detallada sobre el resultado conseguido por el investigador correspondiente.

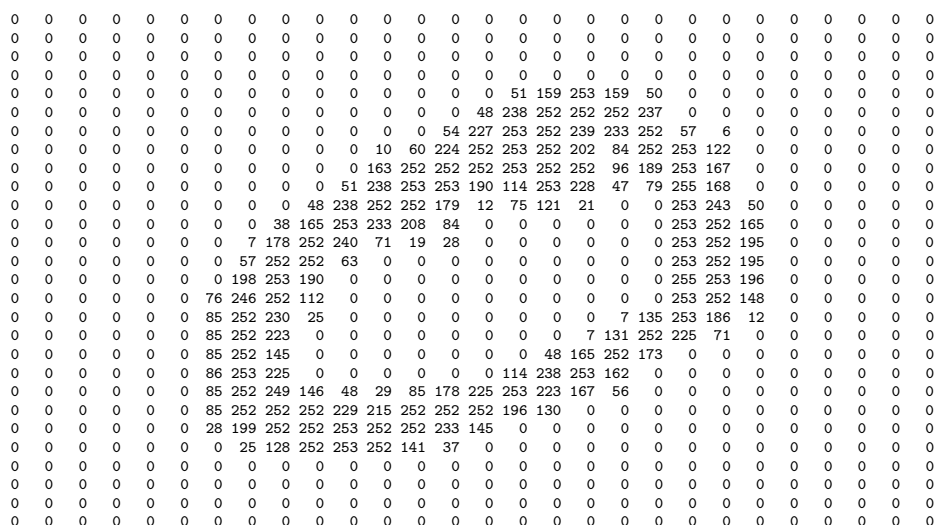
MNIST es una base de datos adecuada para aquellos que desean probar técnicas de aprendizaje automático y métodos de procesamiento de patrones en datos reales dedicando un esfuerzo mínimo al procesamiento de las imágenes y el formato.

B.2. Carga de datos

Los ficheros Octave mencionados anteriormente son ficheros *ascii* comprimidos en formato Octave y están disponibles en PoliformaT.

⁶<http://yann.lecun.com/exdb/mnist>

⁷Estos ficheros se generan al ejecutar el bash script `00-preprocess.sh` disponible en PoliformaT.



Si examinamos `train-images-idx3-ubyte.mat.gz` desde el intérprete de comando (por ejemplo con `zless`) veremos que contiene una cabecera como esta:

Donde `#name` indica el nombre de la variable (matriz Octave) en la que se cargarán los datos, `#type` es el tipo de variable, `#rows` es el número de filas y finalmente `#columns` es el número de columnas. Por lo tanto el fichero `train-images-idx3-ubyte.mat.gz` contiene una matriz representando 60.000 imágenes por filas y 784 dimensiones por columnas.

En Octave cargaremos esta matriz con la siguiente orden:

```
octave:1> load train-images-idx3-ubyte.mat.gz
```

Comprobaremos que la variable X se ha cargado, así como su tamaño:

```
octave:2> size(X)
ans =
```

60000 784

El fichero de datos de test `t10k-images-idx3-ubyte.mat.gz` cuando se cargue se instanciará en la variable `Y`.

Dado que el objetivo es evaluar la tasa de error de clasificación disponemos de las etiquetas de clase de las imágenes, tanto de entrenamiento como de test. Estas etiquetas están en los ficheros `train-labels-idx1-ubyte.gz` y `t10k-labels-idx1-ubyte.mat.gz`, respectivamente. Cargad estos ficheros para comprobar las dimensiones de las variables.

B.3. Visualización de dígitos

Podemos visualizar la imagen de la Fig. 1 que se corresponde con la fila 2 de la variable `X` teniendo en cuenta que es una imagen de 28 x 28 píxeles:

```
octave:3> x=X(2,:);  
octave:4> xr=reshape(x,28,28);  
octave:5> imshow((255-xr)',[])
```

También podemos visualizar las 20 primeras imágenes de la base de datos MNIST para hacernos una idea de la dificultad de esta tarea real:

```
octave:6> for i=1:20  
> xr=reshape(X(i,:),28,28); imshow((255-xr)',[]); pause(1);  
> end
```

C. Clasificador gaussiano

En la asignatura de Percepción se implementó un clasificador gaussiano cuya teoría e implementación refrescamos en esta sección. El código está disponible en PoliformaT.

C.1. Estimación de parámetros y clasificación

Antes de abordar la implementación del clasificador gaussiano es necesario recordar que este clasificador es una instanciación del clasificador de Bayes:

$$\begin{aligned} c^*(x) &= \operatorname{argmax}_{c=1,\dots,C} P(c \mid x) \\ &= \operatorname{argmax}_{c=1,\dots,C} P(c) p(x \mid c) \\ &= \operatorname{argmax}_{c=1,\dots,C} \log P(c) + \log p(x \mid c) \end{aligned}$$

donde la f.d. condicional $p(x \mid c)$ se modeliza mediante una distribución concreta. En esta sección, la f.d. condicional $p(x \mid c)$ será una distribución gaussiana D-dimensional

$$p(\mathbf{x} \mid c) \sim \mathcal{N}_D(\boldsymbol{\mu}_c, \Sigma_c), \quad c = 1, \dots, C$$

Por tanto,

$$\begin{aligned} c^*(\mathbf{x}) &= \operatorname{argmax}_{c=1,\dots,C} \log P(c) + \log p(\mathbf{x} \mid c) \\ &= \operatorname{argmax}_{c=1,\dots,C} \log P(c) - \frac{D}{2} \log 2\pi - \frac{1}{2} \log |\Sigma_c| - \frac{1}{2} (\mathbf{x} - \boldsymbol{\mu}_c)^t \Sigma_c^{-1} (\mathbf{x} - \boldsymbol{\mu}_c) \\ &= \operatorname{argmax}_{c=1,\dots,C} \log P(c) - \frac{1}{2} \log |\Sigma_c| - \frac{1}{2} \mathbf{x}^t \Sigma_c^{-1} \mathbf{x} + \boldsymbol{\mu}_c^t \Sigma_c^{-1} \mathbf{x} - \frac{1}{2} \boldsymbol{\mu}_c^t \Sigma_c^{-1} \boldsymbol{\mu}_c \\ &= \operatorname{argmax}_{c=1,\dots,C} \log P(c) - \frac{1}{2} \mathbf{x}^t \Sigma_c^{-1} \mathbf{x} + \boldsymbol{\mu}_c^t \Sigma_c^{-1} \mathbf{x} + \left(-\frac{1}{2} \log |\Sigma_c| - \frac{1}{2} \boldsymbol{\mu}_c^t \Sigma_c^{-1} \boldsymbol{\mu}_c \right). \end{aligned}$$

Si lo expresamos en términos de función discriminante cuadrática con \mathbf{x} :

$$c^*(\mathbf{x}) = \operatorname{argmax}_{c=1,\dots,C} g_c(\mathbf{x}) = \operatorname{argmax}_{c=1,\dots,C} \log P(c) + \mathbf{x}^t W_c \mathbf{x} + \mathbf{w}_c^t \mathbf{x} + w_{c0}$$

donde

$$W_c = -\frac{1}{2} \Sigma_c^{-1} \tag{10}$$

$$\mathbf{w}_c = \Sigma_c^{-1} \boldsymbol{\mu}_c \tag{11}$$

$$w_{c0} = -\frac{1}{2} \log |\Sigma_c| - \frac{1}{2} \boldsymbol{\mu}_c^t \Sigma_c^{-1} \boldsymbol{\mu}_c \tag{12}$$

Para simplificar esta implementación hemos dejado el término $\log P(c)$ fuera de la componente constante de la probabilidad condicional gaussiana. La estimación máximo verosímil de los parámetros del clasificador gaussiano es ampliamente conocida:

$$\begin{aligned}\hat{P}(c) &= \frac{N_c}{N} \\ \hat{\boldsymbol{\mu}}_c &= \frac{1}{N_c} \sum_{n:c_n=c} \mathbf{x}_n \\ \hat{\Sigma}_c &= \frac{1}{N_c} \sum_{n:c_n=c} (\mathbf{x}_n - \hat{\boldsymbol{\mu}}_c)(\mathbf{x}_n - \hat{\boldsymbol{\mu}}_c)^t\end{aligned}$$

C.2. Implementación y experimentación

A continuación se muestra una implementación en Octave del clasificador gaussiano `gaussian.m` disponible en PoliformaT. Primero, veremos el código correspondiente a la estimación de parámetros:

```
1 function [err] = gaussian(X,xl,Y,yl,alphas)
2
3 classes=unique(xl);
4 N=rows(X);
5 M=rows(Y);
6 D=columns(X);
7
8 % Parameter estimation
9 for c=classes'
10   ic=find(c==classes);
11   idx=find(xl==c);
12   Xc=X(idx,:);
13   Nc=rows(Xc);
14   pc(ic)=Nc/N;
15   muc=sum(Xc)/Nc;
16   mu(:,ic)=muc';
17   sigma{ic}=(Xc-muc)'*(Xc-muc)/Nc;
18 end
```

Las líneas 3-6 obtienen el vector de etiquetas de clases `[0 1 ... 9]`, el número de muestras de entrenamiento y evaluación, y la dimensionalidad de los datos.

A continuación sigue el bucle (líneas 8-18) donde se estimarán los parámetros de cada clase, teniendo en cuenta que `ic` es el índice (o posición) de la clase `c` dentro del vector `classes` y `idx` es el vector de índices a muestras de la clase `c` sobre las filas de la matriz `X`. Recuerda que `ic` es un índice que tomará valores en `[1 2 ... 10]`, mientras `c` es una etiqueta de clase en el rango de valores `[0 1 ... 9]`.

En las líneas 13-15 se estima la probabilidad a priori y media de la clase c . El vector de probabilidades a priori pc será un vector fila, y la matriz de medias μ dispone la media de cada clase como un vector columna.

En la línea 17 se estima la matriz de covarianzas de la clase c y se almacena en la posición ic de un vector de celdas. Estos vectores de celdas permiten almacenar tipos diferentes en el mismo vector, pero en nuestro caso nos será útil para crear un vector de matrices de covarianzas e indexarlo fácilmente.

Las líneas 20-36 de la función `gaussian` contienen el bucle que realizará el cálculo del error en el conjunto de evaluación para los valores de α de suavizado que contiene el vector `alphas`.

```

20 for i=1:length(alphas)
21
22     % Smoothing with identity matrix
23     for c=classes'
24         ic=find(c==classes);
25         ssigma{ic}=alphas(i)*sigma{ic}+(1-alphas(i))*eye(D);
26     end
27
28     % Compute g for each sample in the evaluation set
29     for c=classes'
30         ic=find(c==classes);
31         gY(:,ic)=log(pc(ic))+compute_pxGc(mu(:,ic),ssigma{ic},Y);
32     end
33
34     [~,idy]=max(gY');
35     err(i)=mean(classes(idy)!=y1)*100;
36 end
37
38 end

```

Primeramente, las líneas 22-26 implementan el suavizado *flat smoothing* de la matriz de covarianzas de cada clase con la matriz identidad:

$$\tilde{\Sigma}_c = \alpha \cdot \hat{\Sigma}_c + (1 - \alpha) \cdot I$$

Seguidamente, las líneas 28-32 muestran el código que estima la función discriminante de cada clase calculando $\log P(c) + \log p(\mathbf{x} | c)$ para todas las muestras del conjunto de evaluación. Observa como la función auxiliar `compute_pxGc` implementa $\log p(\mathbf{x} | c)$ estimando para cada muestra del conjunto de evaluación la log probabilidad condicional modelizada mediante una distribución gaussiana. Por último, la línea 34 obtiene el índice de la clase [1 2 ... 10] cuyo valor de su función discriminante es máximo para cada muestra del conjunto de validación. La línea 35 es el cálculo de la probabilidad de error empírica, es decir, el número de errores promedio sobre el conjunto de evaluación. Observa como la indexación del vector `classes` mediante el vector de índices `idy`, es

decir, `classes(idy)`, realiza la conversión de índices de clases [1 2 ... 10] a etiquetas de clase [0 1 ... 9].

A continuación se muestra la implementación de la función auxiliar `compute_pxGc`, invocada en la línea 31:

```

41 function [pxGc] = compute_pxGc(mu,sigma,X)
42     qua=-0.5*sum((X*pinv(sigma)).*X,2);
43     lin=X*(mu'*pinv(sigma))';
44     cons=-0.5*logdet(sigma);
45     cons=cons-0.5*mu'*pinv(sigma)*mu;
46     pxGc=qua+lin+cons;
47 end

```

Esta función calcula la log probabilidad condicional para la gaussiana con media `mu` y matriz de covarianzas `sigma` de cada muestra en `X`. La línea 42 se corresponde con el término cuadrático que pre y posmultiplica `X` en la Ec. 10, la línea 43 implementa el término lineal, que multiplica `X` en la Ec. 11, y las líneas 44 y 45 implementan el término constante en la Ec. 12.

Como puedes observar, para calcular la inversa de la matriz de covarianzas utilizamos la función pseudoinversa `pinv` en lugar de la inversa convencional `inv`. Esto se debe a que no se puede calcular la inversa de una matriz cuyo rango no es completo, pero sí la pseudoinversa. Es decir, nuestra matriz de covarianzas para MNIST es singular, y por tanto tiene filas (o columnas) que son linealmente dependientes unas de otras. Puedes comprobarlo rápidamente ejecutando `rank(cov(X,1))` y viendo que el rango es inferior a la dimensionalidad de los datos $D = 784$. Esto se debe mayormente a que tenemos muchas dimensiones que son siempre cero, al ser parte del fondo sobre el que está centrado el dígito.

Como consecuencia de que la matriz de covarianzas sea singular, su determinante, que se calcula en la línea 44, será cero y su logaritmo `-Inf` imposibilitando el cálculo de esta log probabilidad. Por ello, es necesario una versión robusta del cálculo del logaritmo del determinante de la matrix de covarianzas que se muestra a continuación:

```

55 function v = logdet(X)
56     lambda = eig(X);
57     if any(lambda<=0)
58         v=log(realmin);
59     else
60         v=sum(log(lambda));
61     end
62 end

```

La función `logdet` reemplaza el logaritmo de cero por el logaritmo de la constante `realmin` en Octave. Además, podrás comprobar que la función `logdet` no calcula la `log(det(sigma))` directamente, sino que para evitar valores excesivamente grandes

(Inf) del determinante, este cálculo se realiza de forma robusta como la suma de logaritmos de los valores propios de la matriz de covarianzas⁸.

El clasificador gaussiano se puede invocar directamente desde Octave habiendo cargado previamente los ficheros de datos (ver Apéndice B), pero lo más funcional es hacerlo desde un script Octave. Por ejemplo, el siguiente script, `gaussian-exp.m` disponible en PoliformaT implementa un barrido del parámetro de suavizado α y la correspondiente evaluación de la tasa de error sobre un conjunto de evaluación:

```

1  #!/usr/bin/octave -qf
2
3  if (nargin!=5)
4  printf("Usage: gaussian-exp.m <trdata> <trlabs> <alphas> <%%trper>...
5  exit(1);
6  end;
7
8  arg_list=argv();
9  trdata=arg_list{1};
10 trlabs=arg_list{2};
11 alphas=str2num(arg_list{3});
12 trper=str2num(arg_list{4});
13 dvper=str2num(arg_list{5});
14
15 load(trdata);
16 load(trlabs);
17
18 N=rows(X);
19 seed=23; rand("seed",seed); permutation=randperm(N);
20 X=X(permutation,:); xl=xl(permutation,:);
21
22 Ntr=round(trper/100*N);
23 Ndv=round(dvper/100*N);
24 Xtr=X(1:Ntr,:); xltr=xl(1:Ntr);
25 Xdv=X(N-Ndv+1:N,:); xldv=xl(N-Ndv+1:N);
26
27 [edv] = gaussian(Xtr,xltr,Xdv,xldv,alphas);
28
29 printf("\n  alpha dv-err");
30 printf("\n----- \n");
31
32 for i=1:length(alphas)
33   printf("%.1e %6.3f\n",alphas(i),edv(i));
34 end

```

⁸<https://www.adelaide.edu.au/mathlearning/play/seminars/evaluate-magic-tricks-handout.pdf>

Las líneas 3-6 comprueban que el número de parámetros sea el correcto, para después parsear esos parámetros:

- **trdata**: fichero de imágenes.
- **trlabs**: fichero de etiquetas de clase.
- **alphas**: vector con el rango de valores de suaviado α a evaluar.
- **trper**: Porcentaje del conjunto de imágenes dedicadas a entrenamiento.
- **dvper**: Porcentaje del conjunto de imágenes dedicadas a validación.

Las líneas 15 y 16 realizan la carga de los ficheros de imágenes y etiquetas de clase, respectivamente. Las líneas 18-20 se emplean para barajar aleatoriamente las imágenes con un semilla predefinida (valor 23) que garantiza que el barajado, aunque aleatorio, será siempre el mismo con esa semilla. Las líneas 22-25 realizan la partición en entrenamiento y validación, dedicando el **trper** % de las muestras desde el principio del fichero a entrenamiento, y el **dvper** % desde el final del fichero a validación. Seguidamente, en la línea 27 se invoca a la función **gaussian** con los conjuntos de entrenamiento y validación, y el vector de valores de α . Para cada uno de estos valores de α se devuelve la tasa de error del clasificador gaussiano suavizado con ese valor de α calculada en el conjunto de validación. Finalmente, las líneas 29-34 imprimen la tasa de error en validación para cada valor de α .

La ejecución del script desde línea de comandos del terminal:

```
./gaussian-exp.m train-images-idx3-ubyte.mat.gz
train-labels-idx1-ubyte.mat.gz
"[1e-8 1e-7 1e-6 1e-5 1e-4 1e-3 1e-2 1e-1 2e-1 5e-1 9e-1 1e1]" 90 10
```

obtiene como resultado:

```
alpha dv-err
-----
1.0e-08 19.550
1.0e-07 18.933
1.0e-06 14.083
1.0e-05 6.317
1.0e-04 4.267
1.0e-03 6.383
1.0e-02 10.000
1.0e-01 11.967
2.0e-01 12.200
5.0e-01 13.550
9.0e-01 18.683
1.0e+01 93.867
```

Como se puede observar, la mejor tasa de error se consigue con $\alpha = 10^{-4}$. La evaluación final con el valor de α óptimo entrenando con todo el conjunto de entrenamiento y evaluando en el conjunto de test se puede realizar con el script `gaussian-eva.m`:

```

1  #!/usr/bin/octave -qf
2
3  if (nargin!=5)
4  printf("Usage: gaussian-eva.m <trdata> <trlabs> <tedata> <telabs>...
5  exit(1);
6  end;
7
8  arg_list=argv();
9  trdata=arg_list{1};
10 trlabs=arg_list{2};
11 tedata=arg_list{3};
12 telabs=arg_list{4};
13 alpha=str2num(arg_list{5});
14
15 % Loading data
16 load(trdata);
17 load(trlabs);
18 load(tedata);
19 load(telabs);
20
21 [ete] = gaussian(X,xl,Y,yl,alpha);
22
23 printf("\n  alpha te-err");
24 printf("\n----- \n");
25 printf("%.1e %6.3f\n",alpha,ete);

```

Si ejecutamos este script desde la línea de comandos del terminal:

```

./gaussian-eva.m train-images-idx3-ubyte.mat.gz train-labels-idx1-ubyte.mat.gz
t10k-images-idx3-ubyte.mat.gz t10k-labels-idx1-ubyte.mat.gz 1e-4

```

obtenemos la tasa de error en el conjunto de test:

```

  alpha te-err
-----
1.0e-04  4.180

```

Esta última tasa de error es la que se puede comparar con las reportadas en la web de MNIST⁹.

⁹<http://yann.lecun.com/exdb/mnist>

D. Función *plot* en Octave

La función `plot`¹⁰ permite la representación gráfica bidimensional de secuencias de puntos definidas por sus coordenadas *x* e *y* como vectores:

```
plot(x, y, format, ...)
```

donde *format* son los argumentos de formato que se define por una cadena con cuatro partes opcionales `<linestyle><marker><color><;displayname>` que indican el tipo de línea ('-' continua, '-' discontinua, ':' punteada, etc.), el tipo de punto ('+' cruz, 'o' círculo, 's' cuadrado, etc.), el color ('r' rojo, 'g' verde, 'b' azul, etc.) y la etiqueta de la leyenda, respectivamente. Por ejemplo:

```
plot (tr(:,1),tr(:,2),"sr")
```

que representa gráficamente el conjunto de entrenamiento como cuadrados rojos. También cabe la posibilidad de añadir a la función `plot` modificadores de propiedades:

```
plot(x, y, format, property, value, ...)
```

```
plot(x, y, property, value, ...)
```

donde algunos valores útiles en esta práctica para *property* son `markersize`, `linewidth`, `markerfacecolor`, `color`, etc., y cuyos valores (*value*) se pueden consultar en el manual de la función `plot`.

Asimismo, será necesario añadir anotaciones¹¹ sobre la gráfica representada. Algunas funciones que te serán de utilidad para realizar estas anotaciones son `text`, `title` y `grid`. Además, la configuración de los ejes¹² se puede realizar mediante la función `axis`. Finalmente, se puede guardar la gráfica a fichero para incluirla posteriormente en la memoria mediante la función `print`¹³.

¹⁰https://octave.org/doc/v4.2.1/Two_002dDimensional-Plots.html

¹¹<https://octave.org/doc/v4.2.1/Plot-Annotations.html>

¹²<https://octave.org/doc/v4.2.1/Axis-Configuration.html>

¹³<https://octave.org/doc/v4.2.1/Printing-and-Saving-Plots.html>

E. Recordatorio de teoría de SVM

Sea $\phi(\mathbf{x}; \boldsymbol{\theta}, \theta_0)$, una FDL obtenida mediante el método SVM, donde $\boldsymbol{\theta}$ es el vector de pesos óptimo que se calcula como:

$$\boldsymbol{\theta} = \sum_{m \in \mathcal{V}} c_m \alpha_m \mathbf{x}_m$$

siendo c_m , la etiqueta de clase, α_m , el multiplicador de Lagrange óptimo, y \mathcal{V} , el conjunto de vectores soporte. El peso umbral θ_0 , por las condiciones de KKT se calcula para cualquier vector soporte $m \in \mathcal{V}$ correcto ($\alpha_m < \mathcal{C}$) como:

$$\theta_0 = c_m - \boldsymbol{\theta}^t \mathbf{x}_m$$

El margen es $\frac{2}{\|\boldsymbol{\theta}\|}$ y la tolerancia de margen ζ_m para un vector soporte $m \in \mathcal{V}$ se calcula como:

$$\zeta_m = 1 - c_m(\boldsymbol{\theta}^t \mathbf{x}_m + \theta_0)$$

Para el caso bidimensional donde $\mathbf{x}, \boldsymbol{\theta} \in \mathbb{R}^2$, tenemos que la ecuación de la recta de separación asociada a ϕ se calcula como:

$$\phi(\mathbf{x}; \boldsymbol{\theta}, \theta_0) = 0 \Rightarrow \theta_1 x_1 + \theta_2 x_2 + \theta_0 = 0 \Rightarrow x_2 = -\frac{\theta_1}{\theta_2} x_1 - \frac{\theta_0}{\theta_2},$$

mientras que las ecuaciones de las rectas que definen las fronteras del margen son:

$$\phi(\mathbf{x}; \boldsymbol{\theta}, \theta_0) = +1 \Rightarrow \theta_1 x_1 + \theta_2 x_2 + \theta_0 = +1 \Rightarrow x_2 = -\frac{\theta_1}{\theta_2} x_1 - \frac{\theta_0 - 1}{\theta_2}$$

$$\phi(\mathbf{x}; \boldsymbol{\theta}, \theta_0) = -1 \Rightarrow \theta_1 x_1 + \theta_2 x_2 + \theta_0 = -1 \Rightarrow x_2 = -\frac{\theta_1}{\theta_2} x_1 - \frac{\theta_0 + 1}{\theta_2}$$