



DEPARTAMENTO
DE COMPUTACION

Facultad de Ciencias Exactas y Naturales - UBA

Métodos Numéricos

Trabajo Práctico 2

Reconocimiento de dígitos.

Resumen

En este trabajo pondremos en práctica distintos algoritmos para reconocer una cierta cantidad de dígitos manuscritos. Se trabajará con una base train, a la cual le realizaremos particiones para entrenar los algoritmos. Se implementará el método de kNN (k -Nearest Neighbors). Debido a que este es sensible a la dimensión de los objetos a considerar, además se implementará el método de Análisis de Componentes Principales para reducir el tamaño de dichos objetos. Se llevarán a cabo experimentaciones para poder determinar los parámetros óptimos para cada método. Al final del trabajo, se llegarán a conclusiones sobre lo descubierto.

Integrante	LU	Correo electrónico
Armagno, Julián	377/12	julian.armagno@gmail.com
More, Ángel	931/12	angel_21_fer@hotmail.com
Pinzón, Germán	475/13	pinzon.german.94@gmail.com
Porto, Jorge	376/11	cuanto.p.p@gmail.com

Palabras claves:

Learning Machine. kNN . Análisis de Componentes Principales. Auto-Valores. Auto-Vectores. Método de las Potencias. K-fold cross validation

Índice

1. Introducción Teórica	3
1.1. k Nearest Neighbor y Principal Component Analysis	3
1.2. Método de la potencia y deflación	4
2. Desarrollo	5
2.1. k-Nearest Neighbors (kNN):	5
2.2. Principal Component Analysis (PCA):	6
3. Resultados	9
4. Discusión	10
5. Conclusiones	11
6. Apéndices	12

1. Introducción Teórica

1.1. k Nearest Neighbor y Principal Component Analysis

En el presente trabajo utilizaremos los métodos de kNN y PCA para llevar a cabo el reconocimiento de dígitos manuscritos. Comenzaremos explicando resumidamente los métodos utilizados, de una manera más general. Dado un vector v y un conjunto de vectores U , queremos saber a que clase pertenece v . Nosotros sabemos a que clase pertenece cada vector $u_i \in U$, entonces una manera sencilla de "estimar" a que clase pertenece v en base a la información que poseemos, es comparándolo con cada vector de U . Una manera de efectuar esta comparación es tomando la distancia (norma 2) entre v y cada vector $u_i \in U$ y quedarnos con la clase del u_m que minimice esa distancia sobre todos los demás. Lo que hace el algoritmo kNN es generalizar un poco esta idea. En lugar de quedarnos con la clase del u_m que minimice la distancia a v , kNN toma las clases de los k vectores u_1, u_2, \dots, u_k cuyas distancias sean mínimas y elige entre ellas la clase que más aparezca. Para conocer el parámetro k es necesario realizar experimentos probando distintos valores y tomar una decisión en base a la calidad de los resultados.

El algoritmo kNN es conceptualmente fácil de entender e implementar, lo cual constituye una ventaja, pero su desventaja principal es el tiempo de cómputo que requiere. Esto se debe a que, como los vectores representan imágenes, la dimensión de estos vectores suele ser grande. Dado que la idea detras del reconocimiento de imágenes en este trabajo radica en utilizar una base de datos relativamente grande (de ahora en más dbTrain) para determinar que tipo de imagen es la que estamos tratando de reconocer, el costo de computar estas distancias se multiplica por la cantidad de imágenes que tenemos en nuestra base de datos. Vemos que este método, así como está planteado, no escala.

Como ya dijimos, por cuestiones de performance no es conveniente utilizar kNN de manera directa con los datos que tenemos. Lo que vamos a intentar hacer entonces, es realizar una transformación de nuestros datos de manera tal que luego, cuando queramos aplicar kNN no nos resulte tan costoso. Recordemos que nuestros datos son vectores en un espacio \mathbb{R}^n , entonces lo que vamos a querer hacer es transformarlos a vectores en otro espacio \mathbb{R}^α tal que $\alpha < n$. Pero también vamos a querer que esta transformación no "cambie por completo" a nuestros vectores, en el sentido de que sigan representando las imágenes que representaban o al menos conserven las "partes relevantes" de ellas. Para poder llevar a cabo esta transformación realizaremos los siguientes pasos:

1. Tomamos los vectores de dbTrain que usaremos para comparar la imagen que queremos reconocer en forma matricial (cada vector una fila de la matriz) y hallar la matriz de covarianza $M \in \mathbb{R}^{n \times n}$.
2. Hallar una matriz $P \in \mathbb{R}^{n \times \alpha}$ que nos permita disminuir la covarianza de los datos de dbTrain. Esto equivale a buscar las variables que tengan la mayor varianza entre sí y covarianza 0. El objetivo de esto es disminuir la redundancia en nuestros datos.
3. Sea P' la matriz P con las primeras α columnas (este parámetro se fija mediante experimentación), es decir $P' \in \mathbb{R}^{n \times \alpha}$, y x_i la i -ésima imagen de dbTrain, realizamos el producto $x'_i = P'^t \hat{x}_i$, ahora x'_i es nuestra nueva i -ésima imagen de train y está en el espacio $\mathbb{R}^{n \times \alpha}$.
4. Sea x una imagen vectorizada cuya clase queremos reconocer, realizamos el producto $x' = P'^t \hat{x}$ donde $\hat{x} = (x - \mu) / (\sqrt{n - 1})$, μ es la media de las imágenes de dbTrain y n la cantidad de imágenes de dbTrain. Como ahora $\hat{x} \in \mathbb{R}^\alpha$ y los vectores de dbTrain están en el mismo espacio, podemos aplicar kNN con x' y los x'_i .

Para obtener P lo que hacemos es hallar la base ortonormal de autovectores de M . Sabemos que dicha base existe por ser M simétrica. Entonces la columna i de P va a ser el vector v_i , el i -ésimo autovector de M . Esta base lo que nos permite hacer es "observar" a nuestros datos desde otro lugar, o sea ahora nuestros ejes de coordenadas van a estar en las direcciones donde más varianza existe entre los datos.

Una vez completados estos tres pasos, cuando queramos reconocer a que clase pertenece un vector v , trabajamos con $v' = P^t v$ y, dado que ahora v' es un vector de \mathbb{R}^α al igual que los vectores de dbTrain podemos aplicar kNN .

1.2. Método de la potencia y deflación

En la sección anterior explicamos los métodos que nos permiten reducir la dimensión de las imágenes con las cuales vamos a trabajar y, en base a cierta información que tenemos, efectuar comparaciones para predecir a que clase pertenece una imagen. Vimos que, uno de los pasos de PCA es obtener una matriz P cuyas columnas son los autovectores de otra matriz M . El método de la potencia, junto con el método de deflación, nos permite encontrar los autovectores y autovalores asociados a la matriz M .

Dada una matriz $A \in \mathbb{R}^{n \times n}$ cuyos autovalores $\lambda_1, \dots, \lambda_n$ cumplen $|\lambda_1| > |\lambda_2| \geq |\lambda_3| \geq \dots \geq |\lambda_n|$ y v_1, \dots, v_n los autovectores asociados, el Método de la potencia estimará v_1 y λ_1 . Decimos estimará porque para tener el valor exacto deberíamos calcular un límite, entonces vamos a "simular" este límite computacionalmente con lo cual el resultado puede tener cierto error. Este método necesita de un vector inicial x_0 y un valor de k relativamente grande y lo que va a hacer es calcular $x_{i+1} = \frac{Ax_i}{\|Ax_i\|_2}$ y $\hat{\lambda}_1 = \frac{\Phi(Ax_{i+1})}{\Phi(Ax_i)}$ desde $i = 1$ hasta k . El resultado será $x_{k+1} = \hat{v}_1 \approx v_1$ y $\hat{\lambda}_1 \approx \lambda_1$. Para que esto se cumpla es importante aclarar que la función $\Phi : \mathbb{R}^n \rightarrow \mathbb{R}$ que usamos debe cumplir las siguientes propiedades:

- Debe ser continua
- Si $v \in \mathbb{R}^n$ y $v \neq 0$ entonces $\Phi(v) \neq 0$
- Si $v \in \mathbb{R}^n$ y $\alpha \in \mathbb{R}$ entonces $\Phi(\alpha v) = \alpha \Phi(v)$

El Método de la potencia también pide que el vector inicial x_0 no sea perpendicular a v_1 , sin embargo a la hora de implementarlo en una computadora esto puede no ser tenido en cuenta y no traerá grandes consecuencias. Esto se debe a que nuestro resultado va a estar arrastrando cierto error a medida que el método itere y este error va a estar cambiando la dirección de x_i (aunque sea mínimamente) y en ese caso dejaría de ser perpendicular a v_1 .

Vimos como funciona el Método de la potencia y que nos devuelve el autovalor de mayor módulo junto con su autovector asociado, sin embargo nosotros queremos obtener todos los autovalores y autovectores de A . Esto se soluciona definiendo una nueva matriz \hat{A} y aplicando nuevamente el Método de la potencia pero ahora sobre \hat{A} . Supongamos que ya tenemos los valores $\hat{v}_1 \approx v_1$ y $\hat{\lambda}_1 \approx \lambda_1$ obtenidos al aplicar el Método de la potencia sobre A , entonces para obtener \hat{v}_2 y $\hat{\lambda}_2$ aplicamos nuevamente el método pero ahora sobre $\hat{A} = A - \hat{\lambda}_1 \hat{v}_1 \hat{v}_1^t$. Para el caso general, si queremos obtener \hat{v}_{i+1} y $\hat{\lambda}_{i+1}$, tenemos \hat{v}_i y $\hat{\lambda}_i$ y nuestra matriz es \hat{A} , definimos $\hat{\hat{A}} = \hat{A} - \hat{\lambda}_i \hat{v}_i \hat{v}_i^t$ y aplicamos el método sobre $\hat{\hat{A}}$. A esto se le llama *deflación*. Los autovalores y autovectores de la nueva matriz que definimos van a ser los mismos que la matriz original (la primera o la que definimos en el paso anterior) excepto por el autovalor de mayor módulo y su autovector asociado.

Si bien vimos que para aplicar el método de la potencia en una matriz A se tiene que cumplir que A tenga un autovalor mayor estricto (de multiplicidad 1) en módulo a todos los demás, para poder hacer deflación iterativamente y obtener todos los autovalores y autovectores de A es necesario que valgan las desigualdades de manera estricta: $|\lambda_1| > |\lambda_2| > |\lambda_3| > \dots > |\lambda_n|$. Esto debe ser así para que cada nueva matriz que definimos cumpla las hipótesis que requiere el Método de la potencia.

2. Desarrollo

2.1. k-Nearest Neighbors (kNN):

El primer método utilizado es kNN. Como se dijo en la introducción teórica, este método es muy intuitivo, se basa en tomar a cada imagen como un punto en un espacio, y se hace una votación de la moda de los k vecinos cercanos. Todas las imágenes de dbTrain (28x28 píxeles), se encuentran guardadas vectorizadas en una matriz de $\mathbb{R}^{42000 \times 785}$. La primera columna establece el label de cada imagen, las columnas restantes son los píxeles.

Mediante el K que determine la partición, vamos a dividir esa matriz, en dos, una de train, con $(42000/K) * (K - 1)$ imágenes y otra de test, con $42000/K$ imágenes. A su vez, vamos a eliminar la primera columna de cada matriz, y nos las vamos a guardar en 2 vectores diferentes. Así podremos comparar cada imagen de la matriz de test, con todas las imágenes de la matriz de train y al mismo tiempo tener guardados los labels correspondientes a cada imagen. Finalmente la idea algorítmica se detalla en el siguiente pseudocódigo:

Algorithm 1 kNN(int k , matriz Test, matriz Train, vector digTest, vector digTrain)

```

1: reconocidos = 0
2: for z = 0; z < Test.filas() ; z++ do
3:   vector<pair<int, int> > normas2AlCuadrado
4:   for m = 0; m < Train.filas(); m++ do
5:     \\ calculamos las distancias
6:     distanciaAlCuadrado = 0
7:     for i = 0; i < Test.columnas() ; i++ do
8:       distanciaAlCuadrado += (Test[z][i] - Train[m][i]) *(Test[z][i] - Train[m][i])
9:     end for
10:    if normas2AlCuadrado.size() < k then
11:      \\ colocamos las primeras k normas
12:      pair<unsigned int, int> a
13:      a.first = digTrain[m]
14:      a.second = distanciaAlCuadrado
15:      normas2AlCuadrado.pushback(a)
16:    else
17:      if haymayor(normas2AlCuadrado, distanciaAlCuadrado) then
18:        \\ si ya tengo k voy sacando las mayores distancias
19:        int posmayor = dondemayor(normas2AlCuadrado)
20:        normas2AlCuadrado[posmayor].first = digTrain[m]
21:        normas2AlCuadrado[posmayor].second = distanciaAlCuadrado
22:      end if
23:    end if
24:  end for
25: end for
26: \\ Ahora hago la votación entre los k vecinos
27: digitos[10] = 0, 0, 0, 0, 0, 0, 0, 0, 0, 0
28: for int t = 0; t < k; t++ do
29:   digitos[normas2AlCuadrado[t].first]++
30: end for
31: ganador = digitos[0]
32: for int x = 0; x < 10; x++ do
33:   if digitos[x] > digitos[ganador] then
34:     ganador = x
35:   end if
36: end for
37: \\ Si en la votación gano el verdadero label de la imagen de test, sumamos reconocidos
38: if ganador == digTest[z] then
39:   reconocidos++
40: end if
41: tasaDeReconocimiento = reconocidos/Test.filas()

```

2.2. Principal Component Analysis (PCA):

El segundo método utilizado es PCA. Dado que las imágenes están representadas por vectores y las dimensiones de los mismos pueden ser muy grandes, a la hora de analizar las imágenes se requiere una gran cantidad de tiempo de cómputo así como de memoria. PCA tiene como objetivo redimensionar dichos vectores, siempre que se obtenga un determinado compromiso entre la nueva cantidad de variables y la calidad de las imágenes que ahora se representan. Para esto se construye un nuevo sistema de ecuaciones donde en el eje i se representa la i -ésima varianza de mayor valor para el conjunto original de imágenes dadas. De esta manera cada variable está correlacionada, siendo las primeras las de mayor importancia, por lo que se pueden obviar las últimas componentes (ya que serían las que menos importancia tienen) reduciendo el número de variables. Para

poder llevar a cabo esta transformacion lineal es necesario construir la matriz de covarianza para el conjunto de valores que representan a las imagenes. Para esto, sean n observaciones (imagenes) y dadas dos coordenadas x_i, x_k su covarianza se obtiene como:

$$\sigma_{x_j x_k} = (1/(n-1)) \sum_{i=1}^n (x_j^{(i)} - \mu_j)(x_k^{(i)} - \mu_k) = (1/(n-1))(x_k^{(i)} - \mu_k v)^t (x_j^{(i)} - \mu_j v)$$

con μ_j, μ_k la media de las coordenadas j y k respectivamente; $x_j^{(i)}, x_k^{(i)}$ las coordenadas j y k de la i -esima imagen vectorizada y $v^t = (1, \dots, 1)$.

De esta forma la covarianza entre dos muestras puede ser calculado como el producto entre dos vectores. Por lo que si $A \in \mathbb{R}^{n \times 784}$ representa n imagenes de train vectorizadas de 784 variables cada una, la covarianza pueden ser calculadas como:

$$M_x = A^t * A'; \quad \text{con } A' = A/\sqrt{n-1} \text{ y } M_x \in \mathbb{R}^{784 \times 784}$$

obteniendo en la diagonal la varianza de cada coordenada. Con el fin de disminuir las redundancias (covarianza), se plantea un cambio de base. Para esto se diagonaliza la matriz M_x , de forma tal que:

$$M_x = P * D * P^t; \quad D \text{ diagonal y } P \text{ matriz ortogonal con los autovectores de } A \text{ como columnas}^1$$

Para obtener los autovectores de M_x se empleó el método de la potencia². Dado que la misma es una técnica iterativa que converge al autovector asociado, se estableció que la misma se detenga luego de alcanzar un maximo de 3000 iteraciones o cuando la diferencia entre cada coordenada del vector obtenido en la iteración k respecto de la $k+1$ sea muy poca (en este caso optamos por 0.0000009). Una vez obtenido el primer autovector procedimos a calcular el segundo, aplicando el método de deflación³. Ambas técnicas serán aplicadas un número α de veces. En la sección experimentación se evaluará como se comportan los resultados y el tiempo al variar este valor. Una vez obtenidos los α autovectores se procedió a llevar a cabo el cambio de base, tanto para las imagenes utilizadas como train como para las de test, por lo que se multiplicó cada autovector obtenido por las imagenes vectorizadas en A' :

$$v_i^t * a^{(i)}; \quad v_i \text{ el } i\text{-esimo autovector y } a^{(i)} \text{ la } i\text{-esima imagen vectorizada}$$

lo que matricialmente se puede obtener al hacer: $A' * P$ (llamaremos a la matriz resultante $TcTrain$). De esta forma la matriz $TcTrain \in \mathbb{R}^{n * \alpha}$. Para poder comparar las imagenes de test (matriz B) con $TcTrain$, es necesario llevar a las imagenes de test a las mismas dimensiones, por lo que a cada una se le restara la media que corresponda y dividirá por $\sqrt{n-1}$, obteniendo B' . Una vez realizado esto procedemos a cambiar la base tal como se hizo para las de train:
 $TcTest = B' * P$

Finalizado los cambios de bases obtenemos imagenes en un mismo espacio vectorial cuyas dimensiones son menores a las originales. Pudiendo aplicar distintas técnicas para llevar a cabo el reconocimiento de los dígitos. Optaremos por aplicar KNN para así poder comparar las diferencias entre tiempo de computo y análisis cuando se busca reconocer dígitos utilizando PCA + KNN y solo KNN.

A continuación se muestra un pseudocódigo de los procedimientos descriptos anteriormente:

Algorithm 2 PCA(matriz Train, matriz Test, int α)

```

1: n = cantidad de imagenes de train
2: \\ Se procede a calcular el promedio de las imagenes de train
3: vectorPromedio  $\leftarrow$  crear(784)
4: for i = 0; i < 784; i++ do
5:     sum = 0
6:     for j = 0; j < n; j++ do
7:         sum = sum + Train[j][i]
8:     end for
9:     promedio[i] = sum/n
10: end for
11: \\ calculamos la matriz de covarianza:
12: for i = 0; i < n; i++ do
13:     for j = 0; j < 784; j++ do
14:         Train[i][j] = Train[i][j] - promedio[j]
15:     end for
16:     Train[i][j] = Train[i][j] /  $\sqrt{n-1}$ 
17: end for
18: matriz covarianza  $\leftarrow$  Traint * Train
19: \\ obtenemos la matriz con los  $\alpha$  autovectores de la matriz de covarianza como columna (P):
20: P  $\leftarrow$  metodo de la potencia y de deflacion (covarianza,  $\alpha$ )
21: for i = 0; i < cantidad de imagenes de test; i++ do
22:     for j = 0; j < 784; j++ do
23:         test[i][j] = test[i][j] - promedio[j]
24:     end for
25:     test[i][j] = test[i][j] /  $\sqrt{n-1}$ 
26: end for
27: \\ realizamos la transformación característica para train (TcTrain) y para test (TcTest):
28: TcTrain  $\leftarrow$  Train * P
29: TcTest  $\leftarrow$  Test * P
30: \\ finalmente procedemos a reconocer los dígitos:
31: Knn(TcTrain, TcTest)

```

3. Resultados

4. Discusión

5. Conclusiones

6. Apéndices

Referencias

- [1] Faires, J. D. and Burden, R. L. Análisis Numérico, 7rd ed. Teorema 9.10 y Corolario 9.11. Pág 555, año 2011.
- [2] Faires, J. D. and Burden, R. L. Análisis Numérico, 7rd ed. Pág 560, año 2011.
- [3] Faires, J. D. and Burden, R. L. Análisis Numérico, 7rd ed. Teorema 9.15. Pág. 570, año 2011.