

Capítulo 2

Representación esparsiva y diccionarios de aprendizaje.

2.1 Reconstrucción de señales

A finales del siglo XIX da surgimiento una teoría matemática que sería de gran uso en el procesamiento de señales. El matemático francés Joseph Fourier estableció que una señal podría ser representada como la suma de series de senos y cosenos, que hoy en día tal postulado se ha dado a conocer como la transformada de Fourier. La transformada de Fourier ha sido de gran utilidad en la resolución de problemas científicos e ingenieriles en diferentes campos, tales como física cuántica, óptica, electrónica y muchos otros [8].

Desde un punto de vista ingenieril la transformada de Fourier se describe como un fenómeno físico más que una herramienta matemática. Las señales pueden ser interpretadas como una combinación lineal de ondas armónicas por lo que se observa que la señal en un instante de tiempo es reemplazada por la suma de varios armónicos.

Otra herramienta útil para el análisis de señales es la transformada Wavelet. Inicialmente un geofísico francés llamado Jean Morlet trabajaba sobre un método para modelar la propagación del sonido a través de la corteza terrestre. Morlet utilizó un sistema basado en una función prototipo, que cumpliendo ciertos requerimientos matemáticos y mediante dos procesos denominados dilatación y traslación se formaban un set de bases que representaban las señales de propagación con la misma robustez y versatilidad que la transformada de Fourier [8].

Las características propias de la transformada Wavelet nos otorgan la posibilidad de representar señales en diferentes niveles de resolución, analizar señales no estacionarias permitiéndonos saber el contenido en frecuencia de una señal [8].

2.2 Representación esparsiva.

La representación esparsiva de señales ha sido de gran interés de estudio en los años recientes. El problema a resolver de la representación esparsiva está en la búsqueda de compactar una señal en términos de una combinación lineal de átomos en un “diccionario sobre-completo”. Comparando los métodos basados en transformaciones orto-normales o procesamiento en el dominio del tiempo, la representación esparsiva ofrece un mejor rendimiento en el modelado de señales [9].

Las señales discretas x son vectores en N -dimensiones denotados en un espacio euclidiano por R^N . El conjunto de señales atómicas $D = \{d_i\}_{i=1}^k$ es llamado base para el conjunto $x = \{x_i\}_{i=0}^M$ si las señales atómicas, también llamadas átomos, son linealmente independientes.

Esto implica que cada señal en x puede ser únicamente reconstruida por una sola combinación lineal de átomos en D .

$$x = \sum_{i=1}^k \alpha_i d_i = D\alpha, \quad (1)$$

donde $D \in R^{N \times K}$ y es una matriz compuesta de K columnas $D = [d_1, d_2 \dots d_k]$ y las entradas del vector $\alpha = [\alpha_1, \alpha_2 \dots \alpha_k]^T$ son los coeficientes de esta combinación lineal.

Si los elementos de la base D son mutuamente ortonormales entonces $d_i^T d_j = 0$ si $i \neq j$ y cada átomo radicara en la hiper- esfera $d_i^T d_i = 1 (D^T D = I)$. Después de esto, la base está completa como es el caso de la base de Fourier.

Permitiendo el uso de átomos linealmente dependientes el diccionario respectivo de la matriz D llega a ser redundante. Existen múltiples opciones del vector α para la construcción de una señal x acorde a $x = D\alpha$.

Decimos que la señal x admite una representación esparsiva sobre la base D cuando es concisamente reconstruida con pocos átomos y la señal es caracterizada por un factor esparsivo L cuando su representación esparsiva α tiene L entradas diferente de cero ($l^0 - norm$).

$$\|\alpha\|_0 \leq L \quad (2)$$

El resultado de la representación, es un poderoso modelo para reducción de almacenamiento y transmisión. Este modelo sugiere diccionarios redundantes óptimos para diferentes clases de señales.

2.3 Algoritmo K-SVD (K Descomposiciones de Valores Singulares)

El algoritmo K-SVD (K Descomposiciones de Valores Singulares) trabaja en conjunto con cualquier otro algoritmo de búsqueda. Se considera un algoritmo simple y diseñado para ser una generalización directa del algoritmo K-means. El K-SVD es muy eficiente debido a la alta codificación esparsiva además de tener un método de actualización de diccionarios acelerado.

2.3.1 K-means algoritmo para Vector de Cuantización (VQ).

Un código que incluye K palabras de código es usado para representar una amplia familia de vectores (señales) $Y = \{y_i\}_{i=1}^N$ ($N \gg K$) mediante la búsqueda del vecino más cercano. Esto lleva a la comprensión o la descripción de esas señales como las agrupaciones en R^N rodeando las palabras de código elegido.

K-means puede ampliarse para sugerir una asignación difusa y una matriz de covarianza por cada grupo, por lo que los datos son modelados por una combinación de Gaussianas.

Denotamos la matriz de código $C = [c_1, c_2, \dots, c_k]$, donde el codewords vienen siendo las columnas.

Cuando \mathbf{C} es dada, cada señal es representada por su palabra de código más cercano, por lo que nosotros podemos escribir $\mathbf{Y}_i = \mathbf{C}\mathbf{x}_i$, donde $\mathbf{x}_i = \mathbf{e}_j$ y es un vector de base trivial con todas las entradas de cero excepto el que está en la posición j -th.

La j indexada es seleccionada como:

$$\forall_{k \neq j} \|\mathbf{y}_i - \mathbf{C}\mathbf{e}_j\|_2^2 \leq \|\mathbf{y}_i - \mathbf{C}\mathbf{e}_k\|_2^2 \quad (3)$$

Esto es considerado como un caso extremo de codificación esparsiva en la que solamente un átomo es permitido para participar en la construcción de \mathbf{y}_i y el coeficiente es forzado a ser 1. La representación MSE (Error Cuadrático Medio) por \mathbf{y}_i es definida como $e_i^2 = \|\mathbf{y}_i - \mathbf{C}\mathbf{x}_i\|_2^2$ y el MSE total es:

$$E = \sum_{i=1}^K e_i^2 = \|\mathbf{Y} - \mathbf{C}\mathbf{X}\|_F^2 \quad (4)$$

El problema de entrenamiento del VQ (Cuantización de Vectores) es encontrar un código \mathbf{C} que minimice el error E , sujeto al límite de la estructura de \mathbf{X} , cuyas columnas se deben de tomar desde la base trivial,

$$\min_{\mathbf{C}, \mathbf{X}} \{\|\mathbf{Y} - \mathbf{C}\mathbf{X}\|_F^2\} \text{ sujeto a } \forall_i, \mathbf{x}_i = \mathbf{e}_k \quad k \neq i \quad (5)$$

El algoritmo K-means es un método iterativo usado para el diseño de un óptimo código para VQ [13]. En cada iteración hay dos etapas, la primera es para la codificación esparsiva que evalúa esencialmente \mathbf{X} y la segunda etapa es para la actualización del código.

La etapa de la codificación esparsiva asume un código conocido \mathbf{C}^{j-1} y calcula una \mathbf{X} factible que minimiza el valor de (5). Es evidente que en cada iteración se garantiza ya sea una reducción o que no haya ningún cambio en el MSE (Error Cuadrático Medio).

El problema de la representación esparsiva puede ser vista como una generalización de VQ (5), donde nosotros permitimos cada señal de entrada para

ser representada por una combinación lineal de palabras de código que llamaremos diccionarios.

Por lo tanto el vector de coeficientes no se le permite más de una entrada diferente de cero, por lo que estos pueden tener valores arbitrarios. Para este caso, la minimización corresponde a la ecuación (5) y de esa forma obtener el mejor diccionario posible de la representación esparsiva de este ejemplo \mathbf{Y} .

$$\min_{D, X} \{ \|\mathbf{Y} - \mathbf{DX}\|_F^2 \text{ sujeto a } \forall_i \|\mathbf{x}_i\|_0 \leq T_o \} \quad (6)$$

donde T_o es el factor de esparsividad.

Un objetivo principal, alternativamente, podría cumplirse considerando:

$$\min_{D, X} \sum_i \|\mathbf{x}_i\|_0 \text{ sujeto a } \|\mathbf{Y} - \mathbf{DX}\|_F^2 \leq \varepsilon, \quad (7)$$

Por un valor fijo ε .

donde ε es un error.

Se minimiza la ecuación (6) de forma iterativa. Primero se fija \mathbf{D} donde el objetivo es encontrar el mejor coeficiente de la matriz \mathbf{X} . Cualquier algoritmo puede ser utilizado para el cálculo de coeficientes siempre y cuando se puede suministrar una solución con un predeterminado número de elementos distintos de cero, T_o .

Una vez hecha la codificación esparsiva, el segundo paso sólo desempeña la búsqueda de un mejor diccionario. Este proceso actualiza una columna en el tiempo, fijando todas las columnas en \mathbf{D} , excepto una, \mathbf{d}_k , y encontrando una nueva columna \mathbf{d}_k y nuevos valores para sus coeficientes logrando reducir el MSE.

El proceso de actualización de una columna de \mathbf{D} en el tiempo es un problema que tiene una sencilla solución basado en el “Descomposiciones de Valores Singulares” (SVD) permitiendo un cambio en los coeficientes mientras actualiza acelerando la convergencia de columnas, ya que las columnas posteriores basaran sus actualizaciones en los coeficientes más relevantes.

2.3.2 Descripción detallada del K-SVD.

Se muestra el K-SVD a detalle. Se recuerda que la función objetivo es

$$\min_{\mathbf{D}, \mathbf{X}} \{ \|\mathbf{Y} - \mathbf{DX}\|_F^2 \} \text{ sujeto a } \forall_i, \|\mathbf{x}_i\|_0 \leq T_0 \quad (8)$$

Se considera en primer lugar la fase de codificación dispersa, donde asumimos que \mathbf{D} es fija, y consideramos el problema de optimización anterior como una búsqueda de representación dispersa con coeficientes resumidos en la matriz \mathbf{X} [10].

El término de penalización puede ser reescrito como:

$$\|\mathbf{Y} - \mathbf{DX}\|_F^2 = \sum_{i=1}^N \|\mathbf{y}_i - \mathbf{D}\mathbf{x}_i\|_2^2$$

donde $\|\mathbf{x}\|_2$ se refiere a la norma 2, que es distancia euclidiana.

Además el problema planteado en (8) puede ser desacoplado para N distintos problemas de la forma:

$$i = 1, 2, \dots, N, \min_{\mathbf{x}_i} \{ \|\mathbf{y}_i - \mathbf{D}\mathbf{x}_i\|_2^2 \} \text{ sujeto a } \|\mathbf{x}_i\|_0 \leq T_0 \quad (9)$$

Pasamos ahora a la segunda parte, al proceso de actualización de los diccionarios con los coeficientes diferente de cero. Asumimos que \mathbf{X} y \mathbf{D} son fijas poniendo en cuestión una columna en el diccionario, \mathbf{d}_k , y a los coeficientes que le corresponden de la i th fila en \mathbf{X} , denotado como \mathbf{x}_T^i (este no es un vector que está en la i th columna de \mathbf{X}).

Regresando al objetivo de la función (5) el término de penalización pueden ser descritas como [10]:

$$\|\mathbf{Y} - \mathbf{DX}\|_F^2 = \|\mathbf{Y} - \sum_{j=1}^K \mathbf{d}_j \mathbf{x}_j\|_F^2 = \|(\mathbf{Y} - \sum_{j \neq k} \mathbf{d}_j \mathbf{x}_j^i) - \mathbf{d}_k \mathbf{x}_T^k\|_F^2 = \|\mathbf{E}_k - \mathbf{d}_k \mathbf{x}_T^k\|_F^2 \quad (10)$$

Se ha descompuesto la multiplicación \mathbf{DX} para la suma de \mathbf{K} matrices de rango 1. La matriz \mathbf{E}_k representa el error para todos los ejemplos \mathbf{N} cuando se retira el átomo de k -ésimo. Para este caso estaría bien el uso del SVD para encontrar \mathbf{d}_k

y \mathbf{x}_T^k . El SVD encuentra el rango más cerca que se aproxime a \mathbf{E}_k , y esto reducirá el error el error definido en (10). Sin embargo habrá un error porque el nuevo vector \mathbf{x}_T^k será ocupado por una nueva actualización de d_k no logrando cumplir la restricción de escasez.

Un remedio para el problema es simple y bastante intuitivo. Se define ω_i como el grupo de índices apuntados para ejemplos $\{y_i\}$ que usa el átomo \mathbf{d}_k donde $\mathbf{x}_T^k(i)$ es diferente de cero [10].

Así que,

$$\omega_k = \{i | 1 \leq i \leq K, \mathbf{x}_T^k(i) \neq 0\}, \quad (11)$$

Donde se define a ω_k tiene una matriz de tamaño $N \times |\omega_k|$, con unos en la $(\omega_k(i), i)$ y ceros en lo demás.

Multiplicando $\mathbf{x}_R^k = \mathbf{x}_T^k \omega_k$, este se contrae a la fila vector \mathbf{x}_T^k descartando las entradas de cero, resultando con la columna vector \mathbf{x}_R^k de longitud $|\omega_k|$ que incluye un subconjunto de ejemplos que actualmente usan al átomo \mathbf{d}_k .

El mismo efecto pasa con $\mathbf{E}_k^R = \mathbf{E}_k \omega_k$ implicando una selección de columnas de error que corresponde a ejemplos que usan el átomo \mathbf{d}_k .

Con esta notación se regresa a (10) y se sugiere minimizar con respecto a \mathbf{d}_k y \mathbf{x}_T^k , pero la solución forzada de $\tilde{\mathbf{x}}_T^k$ para tener el mismo soporte como la original \mathbf{x}_T^k .

Esto es equivalente a la minimización de [10]:

$$\|\mathbf{E}_k \omega_k - \mathbf{d}_k \mathbf{x}_T^k \omega_k\|_F^2 = \|\mathbf{E}_k^R - \mathbf{d}_k \mathbf{x}_R^k\|_F^2 \quad (12)$$

Tomando la matriz restringida \mathbf{E}_k^R , el SVD la descompone a $\mathbf{E}_k^R = \mathbf{U} \Delta \mathbf{V}^T$. Se define la solución $\tilde{\mathbf{d}}_k$ como la primera columna de \mathbf{U} , y el vector de coeficiente \mathbf{x}_R^k como la primera columna de \mathbf{V} multiplicada por $\Delta(1,1)$. Para esta solución fue necesario [10]:

- Las columnas de \mathbf{D} deben permanecer normalizadas
- Cualquier soporte de todas las representaciones se queda igual o presenta una pequeña anulación de los términos

Se llama a este algoritmo K-SVD. Mientras K-Means aplica K cálculos de principales actualizaciones del código, el K-SVD obtiene una actualización del diccionario por el cálculo K-SVD, que determina una columna.

A continuación se muestra el algoritmo del K-SVD [10],

Encuentra el mejor diccionario para representar los datos muestreados $\{y_i\}_{i=1}^N$ como composición esparsiva resolviendo,

$$\min_{D, X} \{ \|\mathbf{Y} - \mathbf{D}\mathbf{X}\|_F^2 \} \text{ sujeto a } \forall_i \|x_i\|_0 \leq T_0.$$

Inicialización: Establecer la matriz del diccionario $D^0 \in R^{n \times k}$ con l^2 columnas normalizadas.

Se establece $J=1$.

Se repite hasta converger.

- Etapa del código esparsivo: usa cualquier algoritmo de búsqueda para calcular la representación de vectores x_i por cada ejemplo y_i para aproximar la solución de,

$$i = 1, 2, \dots, N, \min_{x_i} \{ \|y_i - \mathbf{D}x_i\|_2^2 \} \text{ Sujeto a } \|x_k\|_0 \leq T_0.$$

- Actualización del código: por cada columna $k = 1, 2, \dots, K$ en $D^{(J-1)}$, Define el grupo de ejemplos que usa este átomo, $\omega_k = \{i | 1 \leq i \leq N, x_i^k \neq 0\}$.

Calcula la representación total de la matriz de error, E_k , por,

$$\mathbf{E}_k = \mathbf{Y} - \sum_{j \neq k} \mathbf{d}_j \mathbf{x}_T^j$$

Restringe E_k al tomar solamente las columnas correspondientes a ω_k y obtener E_k^R .

- Aplica la descomposición SVD $E_k^R = U\Delta V^T$. Toma la columna \tilde{d}_k del diccionario actualizado para ser la primera columna de U . Actualiza el vector de coeficiente X_R^k para ser la primera columna de V multiplicado por $\Delta(1,1)$.
- Establece $J = J + 1$.

2.4 Aprendizaje de diccionarios redundantes.

En lugar de usar diccionarios predefinidos, tal como las wavelets para la reconstrucción esparsiva de señales, los diccionarios pueden ser adaptados para un conjunto de señales de entrenamiento.

Dejando las columnas de la matriz $X = [x_1, x_2, \dots, x_j] \in \mathbf{R}^{N \times j}$ como el conjunto de señales de entrenamiento de una clase particular para ser reconstruida esparsivamente a través de un óptimo diccionario redundante D . El aprendizaje de diccionario se describe con el siguiente problema de minimización:

$$[D, A] = \arg \min_{DA} \|X - DA\|_F^2 \text{ subject to } \|\alpha_i\|_0; \forall i = 1, 2, \dots, j \quad (13)$$

Las columnas de la matriz $A = [\alpha_1, \alpha_2, \dots, \alpha_j] \in \mathbf{R}^{K \times j}$ son la representación esparsiva del conjunto de señales en X tomando en cuenta que L es el factor esparsivo.

Diferentes enfoques para el aprendizaje de diccionarios han sido desarrollados y basados en dos pasos. El primer paso consiste en encontrar una representación esparsiva A del entrenamiento de la señal X basado en un diccionario D a través de un algoritmo de búsqueda tal como el *Matching Pursuit* (MP) u *Orthogonal Matching Pursuit* (OMP).

Durante el segundo paso, las señales atómicas son actualizadas asumiendo coeficientes de reconstrucción fija. Un algoritmo adecuado para adaptar

diccionarios para la representación esparsiva es el método K-SVD (K-Descomposición de Valores Singulares).

La idea principal del algoritmo K-SVD consiste en expresar la reconstrucción total de la función de error.

$$\|\mathbf{X} - \mathbf{D} \mathbf{A}\|_F^2 = \|\mathbf{X} - \sum_{i=1}^K \mathbf{d}_i \alpha_i^T\|_F^2 = \left\| \mathbf{X} - \sum_{\substack{i=1 \\ i \neq j}}^K \mathbf{d}_i \alpha_i^T - \mathbf{d}_j \alpha_j^T \right\|_F^2 = \|\mathbf{E}_j - \mathbf{d}_j \alpha_j^T\|_F^2, \quad (14)$$

Donde el átomo \mathbf{d}_i es el i th columna de \mathbf{D} y α_i^T es la i th columna de \mathbf{A} . El algoritmo K-SVD es un proceso iterativo de dos pasos para formar un diccionario, en cada iteración el primer paso consiste en estimar la representación esparsiva de \mathbf{A} de acuerdo a [10]:

$$\mathbf{A} = \arg \min_{\mathbf{A}} \|\mathbf{X} - \mathbf{D} \mathbf{A}\|_F^2 \text{ subject to } \|\alpha_i\|_0 < L; i = 1, 2, \dots, J. \quad (15)$$

donde J es el número de señales de prueba.

En el segundo paso, cada átomo \mathbf{d}_i y correspondiendo a la columna α_i^T son encontrados al usar el método SVD conforme a [10]:

$$[\mathbf{d}_j, \alpha_j^T] = \arg \min_{\mathbf{d}_j, \alpha_j^T} \|\mathbf{E}_j \Omega_j - \mathbf{d}_j \alpha_j^T\|_F^2, \quad (16)$$

Donde la matriz Ω_j reduce \mathbf{E}_j manteniendo sólo aquellas columnas que tienen coeficientes de reconstrucción diferentes de cero en α_j^T . Usando el SVD, la matriz $\mathbf{E}_j \Omega_j$ puede ser aproximada por una matriz rank-1.

$$\mathbf{E}_j \Omega_j = \sum_k \sigma_k \mathbf{u}_k \mathbf{v}_k^T \approx \sigma_1 \mathbf{u}_1 \mathbf{v}_1^T; \sigma_1 > \sigma_2 > \sigma_3 > \dots, \quad (17)$$

donde $\mathbf{d}_j = \mathbf{u}_1$, y $\alpha_j^T = \sigma_1 \mathbf{v}_1^T$

2.5 Clasificación de señales basadas en aprendizaje de diccionarios

El aprendizaje de diccionarios se ha aplicado con éxito a diferentes tareas de aplicación. Dejando el diccionario $\mathbf{D}_i \in R^{N \times K}$ al ser entrenado esparsivamente reconstruye un conjunto de señales $X_i \in R^{N \times J}$ perteneciente a la i th clase C_i [10].

Una prueba de la señal $X \in R^N$ es asignada para una de las P clases para la primera estimación del conjunto de residuos $\{r_i(\mathbf{x}, \mathbf{D}_i)\}_{i=1}^P$ donde cada residuo corresponde a la reconstrucción esparsiva de la señal bajo un diccionario.

$$r_i(\mathbf{x}, \mathbf{D}_i) = \min_{\alpha} [\|\mathbf{x} - \mathbf{D}_i \alpha\|_2^2 + \lambda \|\alpha\|_0]; i = 1, \dots, P; \quad (18)$$

Después se asigna una clase para encontrar el residuo más pequeño.

$$class = \arg \min_{i \in \{1, \dots, p\}} r_i(\mathbf{x}, \mathbf{D}_i) \quad (19)$$