

ESCUELA COLOMBIANA DE INGENIERÍA
JULIO GARAVITO

PROYECTO DE GRADO

Advanced Natural Language Processing Techniques to Profile Cybercriminals

Autor:

Alejandro ANZOLA ÁVILA

Director:

Dr. Daniel Orlando DÍAZ LÓPEZ

Programa de Ingeniería de Sistemas

BOGOTÁ, COLOMBIA



31 de marzo de 2019

Índice general

1. Introducción	1
1.1. Objetivo general	1
1.2. Objetivos específicos	1
2. Cronograma	2
3. Marco teórico	4
3.1. Análisis de vínculos (Link Analysis)	4
3.2. Agentes de software (Software Agents)	5
3.3. Minería de datos (Data Mining)	5
3.3.1. Minería de texto (Text Mining)	5
3.3.2. Clasificación	5
Metodologías de clasificación	6
3.3.3. Clustering	6
3.4. Sistemas Basados en Conocimiento (Knowledge Based Systems)	6
3.4.1. Fuzzy Knowledge Based Systems	7
3.5. Redes Neuronales Artificiales (Artificial Neural Network)	7
3.5.1. Sistemas de Detección de Anomalías (Anomaly Detection Systems)	7
3.5.2. Mapa autoorganizado (Self-organizing Maps)	7
3.6. Maquina de soporte vectorial (Support Vector Machine)	10
3.6.1. Maximum Margin Hyperplanes	10
3.6.2. Función Kernel	12
3.7. Clasificadores Bayesianos	12
3.7.1. Teorema de Bayes	12
Usando el teorema de Bayes para clasificación	12
3.7.2. Clasificador Naïve Bayes	13
Como funciona el clasificador Naïve Bayes	13
3.8. Aprendizaje de maquina (Machine Learning)	13
4. Estado del arte	14
5. Propuesta	16
5.1. Análisis	16
5.2. Diseño	17
5.3. Resultados	18
Glosario	21
Bibliografía	24

Índice de figuras

2.1. Diagrama Gantt de actividades de 1 ^{er} periodo	2
2.2. Diagrama Gantt de actividades de 2 ^{do} periodo	2
3.1. Arquitectura KBS	6
3.2. Proceso de adaptación de SOM	9
3.3. Ejemplo de salida de SOM	9
3.4. Ejemplo de uso de SOM en aplicaciones de perfilado	10
3.5. Maximum Margin Hyperplanes	11
3.6. Transformación de espacios en Support Vector Machine	11

Índice de cuadros

2.1. Detalle de cronograma de actividades	3
---	---

Capítulo 1

Introducción

1.1. Objetivo general

El objetivo de este proyecto es generar herramientas y estrategias para el perfilado de cibercriminales con ayuda de metodologías de *Natural Language Processing (NLP)* aplicado a datos recolectados de comunicaciones y redes sociales.

1.2. Objetivos específicos

- Diseñar e implementar una solución de lenguaje natural para realizar el perfilado de sospechosos.
- Identificar el estado del arte en sistemas que usan NLP para apoyar agencias de seguridad del Estado.
- Implementación de artefactos para la construcción de *Datasets* con información recolectada de medios privados como de fuentes abiertas.
- Validar la solución desarrollada frente a un escenario real.
- Modelado de diferentes metodologías, heurísticas y meta–heurísticas para NLP.

Capítulo 2

Cronograma

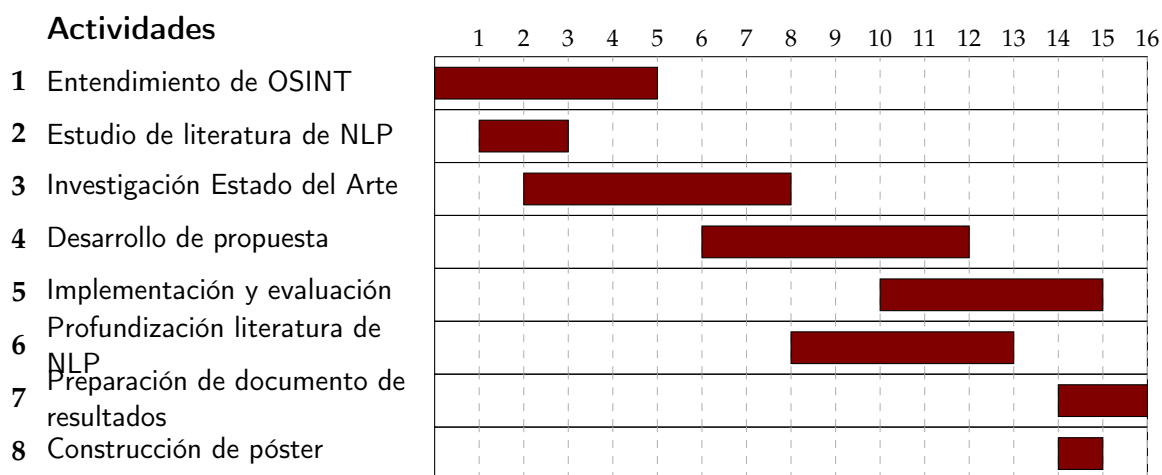


FIGURA 2.1: Diagrama Gantt de actividades de 1^{er} periodo

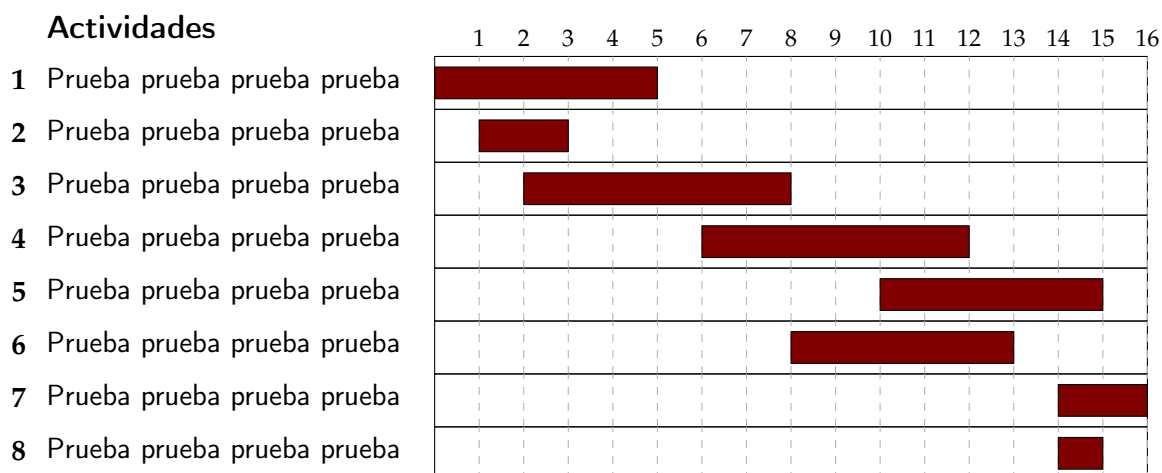


FIGURA 2.2: Diagrama Gantt de actividades de 2^{do} periodo

	Detalle
1	Entendimiento de NLP y proyecto <i>Open Source Intelligence (OSINT)</i>
2	Estudio de la literatura
3	Investigación del Estado del Arte
4	Desarrollo de la propuesta de investigación
5	Desarrollo de la implementación y pruebas
6	Realización de curso de NLP de National Research University Higher School of Economics de <i>Coursera</i>
7	Preparación final de documento de libro de proyectos y artículo de investigación
8	Construcción de póster para presentación en la Vitrina Académica

CUADRO 2.1: Detalle de cronograma de actividades

Capítulo 3

Marco teórico

La Web contiene una gran cantidad de opiniones respecto a productos, políticos, y mucho mas, expresado en forma de noticias, sitios de opinión, reseñas en tiendas online, redes sociales. Como resultado, el problema de “Minería de opinión” ha obtenido una atención creciente en las ultimas dos décadas y es un factor decisivo para las nuevas organizaciones (como es mencionado en [5]). De esto mismo partimos que el análisis de textos para extraer el significado y demás componentes extraíbles del texto componen un factor que debe considerarse al momento de realizar decisiones, de manera que los avances hechos hasta ahora tienen como meta una aplicación practica de lo que se conoce como Natural Language Processing.

Luego de los ataques terroristas del 11 de Septiembre de 2001 en Estados Unidos, se realizaron fuertes criticas respecto a la inteligencia, donde el director del FBI *Robert S. Mueller* indico que el principal problema que la agencia tuvo fue que se enfocaba demasiado en lidiar con el crimen luego de que fue cometido y ponía muy poco énfasis en prevenirlo (adaptado de [3]). Es por esto que el uso de NLP para temas de seguridad como también de metodologías de Machine Learning y Deep Learning han sido ampliamente utilizadas en ámbito de seguridad luego de estos eventos.

Para obtener una mejor inteligencia se necesito de mejores tecnologías a las que se tenían entonces (véase [3, pág 2]):

- Integración de datos (ó *Data Integration (DI)* en inglés)
- Análisis de vínculos (ó *Link Analysis (LA)* en inglés)
- Agentes de software (ó *Software Agents (SA)* en inglés)
- Minería de texto (ó *Text Mining (TM)* en inglés)
- Redes neuronales (ó *Artificial Neural Network (ANN)* en inglés)
- Algoritmos de Machine Learning (ó *Machine Learning Algorithms (MLA)* en inglés)

3.1. Análisis de vínculos (Link Analysis)

Es la visualización de asociaciones entre entidades y eventos, por lo general involucran una visualización por medio de una gráfica o un mapa que muestre las relaciones entre sospechosos y ubicaciones, sea por medio físico o por comunicaciones en la red.

3.2. Agentes de software (Software Agents)

Es el software que realiza tareas asignadas por el usuario de manera autónoma, donde sus habilidades básicas son:

- **Realización de tareas:** Hacen obtención de información, filtrado, monitoreo y reporte.
- **Conocimiento:** Pueden usar reglas programadas, o pueden aprender reglas nuevas (véase 3.4).
- **Habilidades de comunicación:** Reportar a humanos e interactuar con otros agentes.

3.3. Minería de datos (Data Mining)

Según [7], la minería de datos se define como el proceso de descubrir información útil en repositorios grandes de datos. Las técnicas de minería de datos son desplegadas para limpiar grandes bases de datos para encontrar patrones nuevos y útiles que de lo contrario podrían permanecer desconocidos. También ofrecen capacidades para predecir la salida de observaciones futuras, tales como predecir si un cliente nuevo gastara más de \$100 en una tienda.

No todas las tareas de descubrimiento de información son considerados como *Data Mining (DM)*. Por ejemplo, realizar una consulta de campos individuales usando un sistema de base de datos o encontrar una página web por medio de una búsqueda en Internet son tareas relacionadas con *adquisición de información*.

3.3.1. Minería de texto (Text Mining)

Es un subcampo de Inteligencia Artificial conocida como Natural Language Processing, en donde las herramientas de minería de datos pueden capturar rasgos críticos del contenido de un documento basado en el análisis de sus características lingüísticas.

La mayoría de los crímenes son electrónicos por naturaleza, por lo que se dejan rastros textuales que investigadores pueden seguir y analizar. Estas se enfocan en el descubrimiento de relaciones en texto no-estructurado y pueden ser aplicados al problema de *búsqueda y localización de palabras clave*.

3.3.2. Clasificación

Clasificación es la tarea de asignarle una de varias categorías predefinidas a objetos, y es una tarea que tiene una variedad extensa de aplicaciones. Ejemplos de esto se encuentran la detección de correos no deseados en mensajes de e-mails basándose del encabezado o el cuerpo del mensaje, categorización de células benignas de malignas basándose en los resultados de escaneados MRI o incluso la clasificación de galaxias basado en su forma.

Definido formalmente, clasificación es la tarea de aprender una función objetivo f que mapea cada conjunto de atributos x a una clase predefinida de etiquetas y .

La función objetivo también se define informalmente como un *modelo de clasificación*.

Metodologías de clasificación

Existen muchos métodos para la clasificación de datos no-estructurados, entre los descritos aquí están:

- Clasificador basado en reglas (véase 3.4)
- Redes neuronales artificiales (véase 3.5)
- Maquinas de soporte vectorial (véase 3.6)
- Clasificador de Naïve Bayes (véase 3.7.2)

3.3.3. Clustering

El análisis de clusters agrupa objetos de datos basándose únicamente en la información encontrada en los datos que describen los objetos y sus relaciones. El objetivo es que objetos dentro de un grupo sean similares (o relacionados) el uno al otro, y que sean diferentes (o sin relación) a objetos en otros grupos. Entre mayor sea la similitud dentro de un grupo y entre mayor sea la diferencia entre grupos, sera mejor o mas distintivo el clustering.

Los métodos de clustering se hacen referencia comúnmente en *Machine Learning* (ML) como métodos no-supervisados, los cuales se describirán en 3.8. Un método de estos se describe en 3.5.2 conocidos como mapas autoorganizados.

3.4. Sistemas Basados en Conocimiento (Knowledge Based Systems)

Según [6], los *Knowledge Based Systems* (KBS) son uno de los mayores miembros de la familia de *Artificial Intelligence* (AI). El KBS consiste de una *Knowledge Base* (KB) y un programa de búsqueda llamado *Inference Engine* (IE) representado en la figura 3.1. La KB puede ser usado como un repositorio de conocimiento de varias formas.

Existen 5 tipos de KBS, donde uno de ellos es conocido como *Expert Systems* (ES), usados como *Rule-based Systems* (RBS), donde su KB esta dado como reglas y el IE esta dado por algo llamado *Working Memory* (WM), que representa los hechos que se conocen inicialmente del sistema junto con los hechos que se van dando como inferencia de las reglas.

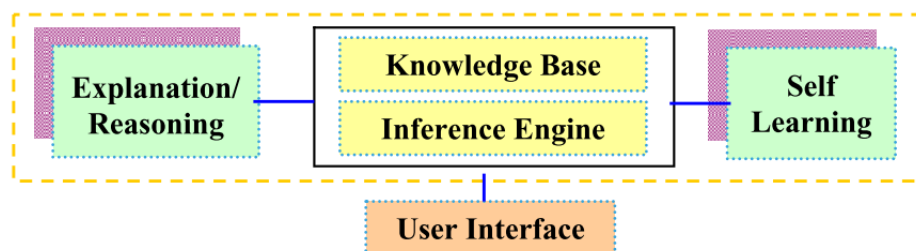


FIGURA 3.1: Arquitectura KBS. Tomado de [6]

Estas reglas pueden resumirse como una colección de condicionales de la forma **IF/ELSE** que se componen de un *antecedente* y un *consecuente*.

Existen dos tipos de RBS, definidos como *Deductive Systems (DS)* y *Reactive Systems (RS)*, donde el DS tiene como objetivo realizar una conclusión en base a los hechos iniciales en la WM, por el otro lado se tienen los RS, los cuales de igual manera a los DS, toman los hechos de la WM y realizan sea una acción interactiva con su entorno o bien una modificación de los hechos que se encuentran en la WM tal como la adición o eliminación de hechos. Tómese el ejemplo de la ecuación 3.1 tomada de [4], donde x es la temperatura y AC es aire acondicionado.

$$\begin{cases} \text{IF } x \text{ es moderado,} & \text{THEN } y = \text{ajustar AC a bajo} \\ \text{IF } x \text{ es alto,} & \text{THEN } y = \text{ajustar AC a moderado a alto} \\ \text{IF } x \text{ es muy alto,} & \text{THEN } y = \text{ajustar AC a alto} \end{cases} \quad (3.1)$$

3.4.1. Fuzzy Knowledge Based Systems

Pendiente

3.5. Redes Neuronales Artificiales (Artificial Neural Network)

El estudio de redes neuronales artificiales fue inspirado por los intentos de simular los sistemas biológicos de neuronas. El cerebro humano se compone principalmente de células nerviosas llamadas *neuronas*, enlazadas con otras neuronas por medio de hebras de fibra conocidas como *axones*. Los axones son usados para transmitir impulsos nerviosos de una neurona a otra cada vez que las neuronas son estimuladas. Una neurona esta conectada a axones de otras neuronas por medio de *dendritas*, las cuales son extensiones desde el cuerpo de la neurona. El punto de contacto entre una dendrita y un axón se conoce como *sinapsis*. Los neurólogos han descubierto que el cerebro humano aprende por medio de cambiar la fuerza de la conexión sináptica entre las neuronas a través de estimulación repetitiva por el mismo impulso.

De manera análoga a la estructura del cerebro humano, una ANN se compone de una estructura interconectada de nodos y vínculos directos.

3.5.1. Sistemas de Detección de Anomalías (Anomaly Detection Systems)

Pendiente

3.5.2. Mapa autoorganizado (Self-organizing Maps)

El objetivo principal de los *Self-organizing Maps (SOM)* es de transformar una patrón de entrada m -dimensional en un mapa discreto uni- o bi-dimensional, donde sus principales características es que es un algoritmo que se basa en *Unsupervised Learning*, es *Feedforward*, tiene una sola capa de neuronas donde su propósito es realizar *Clustering* y una reducción de dimensionalidad sobre los datos de una forma topologicamente ordenada.

Los SOM tienen tres características distintivas:

- **Competencia:** por cada patrón de entrada, las neuronas en la red competirán entre ellas para determinar un ganador.
- **Cooperación:** la neurona ganadora determina la ubicación espacial (vecinos) alrededor de donde otras vecinas también se verán estimuladas.
- **Adaptación:** la neurona ganadora como también sus vecinas tendrán sus pesos asociados actualizados, y se tiene que los vecinos entre mas cerca estén del ganador, mayor es el grado de adaptación.

El algoritmo de aprendizaje de SOM parte de primero inicializar los pesos de las o neuronas con pesos aleatorios pequeños de una distribución de probabilidad aleatoria o uniforme, donde cada vector de entrada se define como $x = [x_1, \dots, x_m]^T \in \mathbb{R}^m$ y la entrada general de N patrones como $\mathbf{X}^{m \times N}$, el vector de pesos de la neurona i es $\mathbf{w}_i = [w_{i1}, \dots, w_{im}] \in \mathbb{R}^{1 \times m}$, con la matriz de pesos $\mathbf{W}^{o \times m}$.

Para alcanzar el objetivo de *competencia*, se realiza por cada patrón de entrada x_i una comparación con cada uno de los pesos de las o neuronas y se establece la de menor distancia respecto x_i (típicamente la distancia Euclidiana), dejando un ganador *winner*, tal como en la ecuación 3.2.

$$winner = \operatorname{argmin}_j \|x_i - w_j\|; j = 1, \dots, o \quad (3.2)$$

Luego de establecer la neurona ganadora, se realiza el paso para alcanzar la *cooperación*, que consiste en que por medio de una función kernel h (típicamente una distribución gaussiana), que permite establecer un área de afectación de las otras neuronas según su ubicación física en el mapa, definidos como r_{winner} y r_j que son la ubicación de la neurona ganadora y la neurona vecina j , en el cual el grado de afectación de la neurona vecina depende de la distancia de la que esta de la neurona ganadora, definido en la ecuación 3.3.

$$h_{j,winner}(t) = \exp\left(\frac{-\|r_j - r_{winner}\|^2}{2\sigma(t)^2}\right) \quad (3.3)$$

Parte importante del proceso de convergencia del SOM es que a medida que avanzan las iteraciones t del algoritmo el área de afectación se va reduciendo como parte del proceso de adaptación, por lo que definimos $\sigma(t) = \sigma_0 \exp(-t/\tau_1)$, donde τ_1 es una constante heurística y σ_0 la dimensión del mapa SOM.

Finalmente para alcanzar la *adaptación* se realiza una actualización de los pesos de la matriz \mathbf{W} en base a la influencia de área $\sigma(t)$ y de una tasa de aprendizaje $\alpha(t) = \alpha_0 \exp(-t/\tau_2)$, donde τ_2 es otra constante heurística y α_0 es una constante de aprendizaje inicial, que debe ser $0 \leq \alpha_0 \leq 1$, la actualización se describe por la ecuación 3.4 y el proceso puede ser visto gráficamente en la figura 3.2, tanto de forma uni- como bi-dimensional.

$$w_j(t+1) = w_j(t) + \alpha(t)h_{j,winner}(t)[x_i - w_j(t)] \quad (3.4)$$

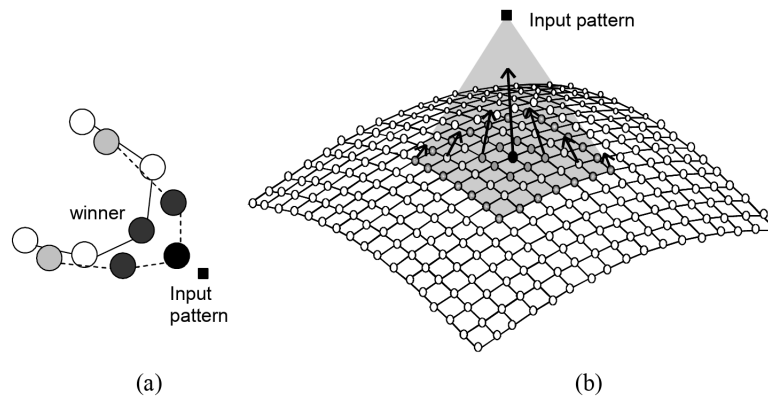


FIGURA 3.2: Proceso de adaptación de SOM, (a) uni-dimensional, (b) bi-dimensional. Tomado de [1]

Luego de que el algoritmo de aprendizaje termina de realizar las iteraciones, la salida de este es la matriz de pesos W , en la figura 3.3 se puede apreciar una aproximación del algoritmo con un mapa uni-dimensional tratando de aproximar una función sinusoidal con ruido adicionado en un gráfico 2D.

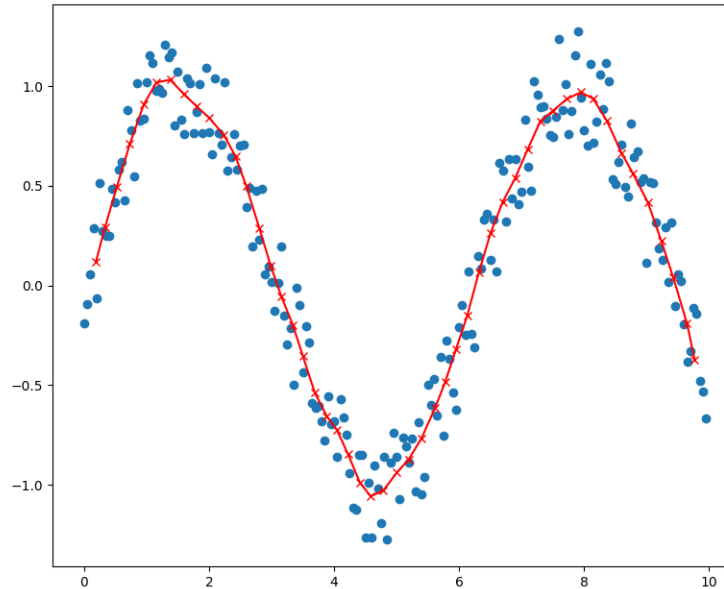


FIGURA 3.3: Ejemplo de salida de SOM. Implementación propia.

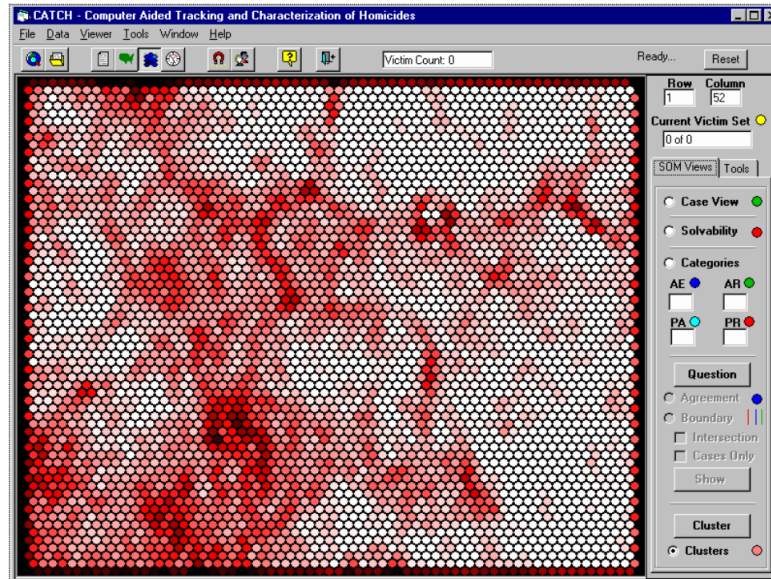


FIGURA 3.4: Ejemplo de uso de SOM en aplicaciones de perfilado. Tomado de [3]

3.6. Máquina de soporte vectorial (Support Vector Machine)

glssvm es una técnica de clasificación que tiene sus raíces en la teoría de aprendizaje estadístico que ha mostrado resultados empíricos prometedores en muchas aplicaciones prácticas, desde reconocimiento de dígitos escritos a mano a categorización de texto. *Support Vector Machine (SVM)* también funciona muy bien con datos de alta dimensionalidad. Otro aspecto destacable de esta aproximación es que representa la frontera de decisión usando un subconjunto de las muestras de entrenamiento, conocidos como los *support vectors*.

3.6.1. Maximum Margin Hyperplanes

Se puede entender a los *Maximum Margin Hyperplanes* como hiper-planos que ayudan a separar datos en un hiper-espacio y que poseen un margen de decisión entre los datos, como ejemplo tómese la figura 3.5, donde el hiper-plano B_1 tiene un margen de decisión mas grande que el hiper-plano B_2 .

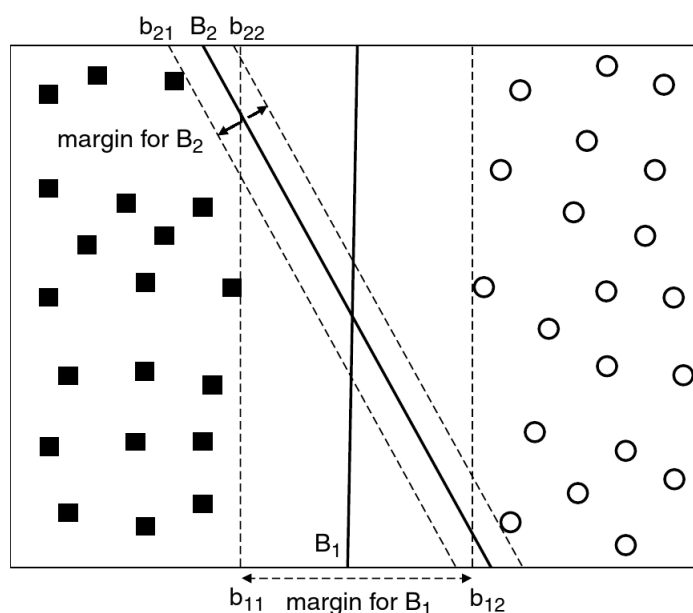


FIGURA 3.5: Maximum Margin Hyperplanes. Tomado de [7]

Finalmente, el objetivo final de los SVM es la búsqueda de un hiper-plano con el mayor margen de decisión. Existen dos tipos de SVM, el lineal y el no-lineal. El lineal realiza la separación de los datos con su hiper-plano a partir de los datos de entrada en su espacio vectorial original, mientras que el no-lineal consta de realizar una transformación de los espacios de los datos de entrada a uno en que sea linealmente separable (véase como ejemplo la figura 3.6), sin embargo al realizar la transformación, el algoritmo de SVM se ve afectado por la dimensionalidad de la entrada, por lo que existe lo que se conoce como la función *kernel* para remediarlo.

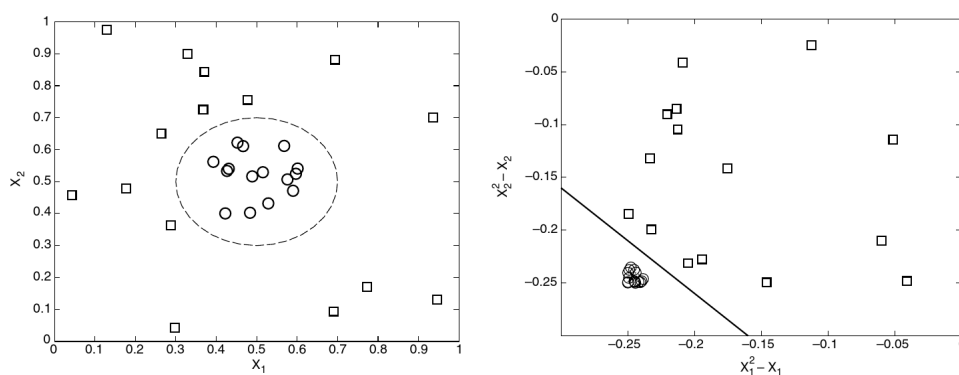


FIGURA 3.6: Transformación de espacios en Support Vector Machine. Tomado de [7]

3.6.2. Función Kernel

La función polinomial de similaridad, K , la cual es calculada en el espacio original de los datos de entrada, se le conoce como la **función Kernel**. En principio se asegura que la función kernel puede ser expresada siempre como el producto punto entre dos vectores de entrada en algún espacio de alta dimensionalidad, la función de kernel también tiene la particularidad de que el computo de los productos punto con la función toman considerablemente menos tiempo que realizar la transformación de espacios, dejando de lado la transformación, acelerando la tarea de clasificación.

3.7. Clasificadores Bayesianos

En muchas aplicaciones de relaciones entre el conjunto de atributos y la etiqueta es no-determinante. Es decir, la etiqueta de clase de un dato de un conjunto de prueba no puede ser determinado con certeza a pesar de ser un atributo idéntico a los atributos de entrenamiento. Esto puede ser producto de que los datos poseen ruido o la presencia de ciertos factores que afectan la clasificación pero no son incluidos en el análisis. Para esto es crucial el teorema de Bayes, el cual es un principio estadístico que combina el conocimiento previo de las clases con la nueva evidencia que se obtiene de los datos.

3.7.1. Teorema de Bayes

El teorema de Bayes dice que para un par de variables aleatorias X y Y y que $P(X = x|Y = y)$ la probabilidad de que la variable X tome el valor x dado que el valor de la variable Y es y . Se tiene entonces la ecuación 3.5.

$$P(Y|X) = \frac{P(X|Y)P(Y)}{P(X)} \quad (3.5)$$

Usando el teorema de Bayes para clasificación

Para denotar el problema de clasificación desde una perspectiva estadística se define a X como el conjunto de atributos y Y como el conjunto de etiquetas de clase. Si la etiqueta de clase tiene una relación no-determinante con los atributos, entonces se pueden tomar a X y a Y como variables aleatorias y capturar su relación probabilística con $P(Y|X)$, conocida como la probabilidad posterior para Y , dada su probabilidad previa $P(Y)$.

Durante la fase de entrenamiento, es necesario aprender las probabilidades posteriores $P(Y|X)$ para cualquier combinación de X y Y basándose en la información recolectada de los datos de entrenamiento.

Dado que lo que se quiere realizar es una clasificación que represente la probabilidad de que dado un valor de $X = x$ este relacionado con que $Y = y$, se puede reconocer primero que X se mantiene constante para lo que son los datos de entrenamiento, y que lo desconocido sea la clasificación $Y = y$ con probabilidad $P(Y|X)$, al conocer esta probabilidad, un valor de prueba X' puede ser clasificado por medio de encontrar la clase Y' que maximice la probabilidad posterior $P(Y'|X')$.

3.7.2. Clasificador Naïve Bayes

Un clasificador de Naïve Bayes estima la probabilidad condicional de las clases por medio de suponer que los atributos son condicionalmente independientes, dado la etiqueta de clasificación y . La suposición de independencia condicional se puede dar por la ecuación 3.6.

$$P(X|Y = y) = \prod_{i=1}^d P(X_i|Y = y) \quad (3.6)$$

Donde cada conjunto de atributos $X = \{X_1, \dots, X_d\}$ que consiste de d atributos.

Como funciona el clasificador Naïve Bayes

Con la suposición de independencia condicional, en vez de computar la probabilidad condicional de clases para cada combinación de X , solo se debe realizar para establecer la probabilidad condicional de cada X_i , dado Y .

Para clasificar un dato de prueba, el clasificador computa la probabilidad posterior para cada clase Y como se muestra en la ecuación 3.7.

$$P(Y|X) = \frac{P(Y) \prod_{i=1}^d P(X_i|Y)}{P(X)} \quad (3.7)$$

3.8. Aprendizaje de maquina (Machine Learning)

Pendiente

Capítulo 4

Estado del arte

Lorem ipsum dolor sit amet, consectetur adipiscing elit. Ut purus elit, vestibulum ut, placerat ac, adipiscing vitae, felis. Curabitur dictum gravida mauris. Nam arcu libero, nonummy eget, consectetur id, vulputate a, magna. Donec vehicula augue eu neque. Pellentesque habitant morbi tristique senectus et netus et malesuada fames ac turpis egestas. Mauris ut leo. Cras viverra metus rhoncus sem. Nulla et lectus vestibulum urna fringilla ultrices. Phasellus eu tellus sit amet tortor gravida placerat. Integer sapien est, iaculis in, pretium quis, viverra ac, nunc. Praesent eget sem vel leo ultrices bibendum. Aenean faucibus. Morbi dolor nulla, malesuada eu, pulvinar at, mollis ac, nulla. Curabitur auctor semper nulla. Donec varius orci eget risus. Duis nibh mi, congue eu, accumsan eleifend, sagittis quis, diam. Duis eget orci sit amet orci dignissim rutrum.

Nam dui ligula, fringilla a, euismod sodales, sollicitudin vel, wisi. Morbi auctor lorem non justo. Nam lacus libero, pretium at, lobortis vitae, ultricies et, tellus. Donec aliquet, tortor sed accumsan bibendum, erat ligula aliquet magna, vitae ornare odio metus a mi. Morbi ac orci et nisl hendrerit mollis. Suspendisse ut massa. Cras nec ante. Pellentesque a nulla. Cum sociis natoque penatibus et magnis dis parturient montes, nascetur ridiculus mus. Aliquam tincidunt urna. Nulla ullamcorper vestibulum turpis. Pellentesque cursus luctus mauris.

Nulla malesuada porttitor diam. Donec felis erat, congue non, volutpat at, tincidunt tristique, libero. Vivamus viverra fermentum felis. Donec nonummy pellentesque ante. Phasellus adipiscing semper elit. Proin fermentum massa ac quam. Sed diam turpis, molestie vitae, placerat a, molestie nec, leo. Maecenas lacinia. Nam ipsum ligula, eleifend at, accumsan nec, suscipit a, ipsum. Morbi blandit ligula feugiat magna. Nunc eleifend consequat lorem. Sed lacinia nulla vitae enim. Pellentesque tincidunt purus vel magna. Integer non enim. Praesent euismod nunc eu purus. Donec bibendum quam in tellus. Nullam cursus pulvinar lectus. Donec et mi. Nam vulputate metus eu enim. Vestibulum pellentesque felis eu massa.

Quisque ullamcorper placerat ipsum. Cras nibh. Morbi vel justo vitae lacus tincidunt ultrices. Lorem ipsum dolor sit amet, consectetur adipiscing elit. In hac habitasse platea dictumst. Integer tempus convallis augue. Etiam facilisis. Nunc elementum fermentum wisi. Aenean placerat. Ut imperdiet, enim sed gravida sollicitudin, felis odio placerat quam, ac pulvinar elit purus eget enim. Nunc vitae tortor. Proin tempus nibh sit amet nisl. Vivamus quis tortor vitae risus porta vehicula.

Fusce mauris. Vestibulum luctus nibh at lectus. Sed bibendum, nulla a faucibus semper, leo velit ultricies tellus, ac venenatis arcu wisi vel nisl. Vestibulum diam. Aliquam pellentesque, augue quis sagittis posuere, turpis lacus congue quam, in hendrerit risus

eros eget felis. Maecenas eget erat in sapien mattis porttitor. Vestibulum porttitor. Nulla facilisi. Sed a turpis eu lacus commodo facilisis. Morbi fringilla, wisi in dignissim interdum, justo lectus sagittis dui, et vehicula libero dui cursus dui. Mauris tempor ligula sed lacus. Duis cursus enim ut augue. Cras ac magna. Cras nulla. Nulla egestas. Curabitur a leo. Quisque egestas wisi eget nunc. Nam feugiat lacus vel est. Curabitur consectetur.

Suspendisse vel felis. Ut lorem lorem, interdum eu, tincidunt sit amet, laoreet vitae, arcu. Aenean faucibus pede eu ante. Praesent enim elit, rutrum at, molestie non, nonummy vel, nisl. Ut lectus eros, malesuada sit amet, fermentum eu, sodales cursus, magna. Donec eu purus. Quisque vehicula, urna sed ultricies auctor, pede lorem egestas dui, et convallis elit erat sed nulla. Donec luctus. Curabitur et nunc. Aliquam dolor odio, commodo pretium, ultricies non, pharetra in, velit. Integer arcu est, nonummy in, fermentum faucibus, egestas vel, odio.

Sed commodo posuere pede. Mauris ut est. Ut quis purus. Sed ac odio. Sed vehicula hendrerit sem. Duis non odio. Morbi ut dui. Sed accumsan risus eget odio. In hac habitasse platea dictumst. Pellentesque non elit. Fusce sed justo eu urna porta tincidunt. Mauris felis odio, sollicitudin sed, volutpat a, ornare ac, erat. Morbi quis dolor. Donec pellentesque, erat ac sagittis semper, nunc dui lobortis purus, quis congue purus metus ultricies tellus. Proin et quam. Class aptent taciti sociosqu ad litora torquent per conubia nostra, per inceptos hymenaeos. Praesent sapien turpis, fermentum vel, eleifend faucibus, vehicula eu, lacus.

Capítulo 5

Propuesta

5.1. Análisis

Lorem ipsum dolor sit amet, consectetur adipiscing elit. Ut purus elit, vestibulum ut, placerat ac, adipiscing vitae, felis. Curabitur dictum gravida mauris. Nam arcu libero, nonummy eget, consectetur id, vulputate a, magna. Donec vehicula augue eu neque. Pellentesque habitant morbi tristique senectus et netus et malesuada fames ac turpis egestas. Mauris ut leo. Cras viverra metus rhoncus sem. Nulla et lectus vestibulum urna fringilla ultrices. Phasellus eu tellus sit amet tortor gravida placerat. Integer sapien est, iaculis in, pretium quis, viverra ac, nunc. Praesent eget sem vel leo ultrices bibendum. Aenean faucibus. Morbi dolor nulla, malesuada eu, pulvinar at, mollis ac, nulla. Curabitur auctor semper nulla. Donec varius orci eget risus. Duis nibh mi, congue eu, accumsan eleifend, sagittis quis, diam. Duis eget orci sit amet orci dignissim rutrum.

Nam dui ligula, fringilla a, euismod sodales, sollicitudin vel, wisi. Morbi auctor lorem non justo. Nam lacus libero, pretium at, lobortis vitae, ultricies et, tellus. Donec aliquet, tortor sed accumsan bibendum, erat ligula aliquet magna, vitae ornare odio metus a mi. Morbi ac orci et nisl hendrerit mollis. Suspendisse ut massa. Cras nec ante. Pellentesque a nulla. Cum sociis natoque penatibus et magnis dis parturient montes, nascetur ridiculus mus. Aliquam tincidunt urna. Nulla ullamcorper vestibulum turpis. Pellentesque cursus luctus mauris.

Nulla malesuada porttitor diam. Donec felis erat, congue non, volutpat at, tincidunt tristique, libero. Vivamus viverra fermentum felis. Donec nonummy pellentesque ante. Phasellus adipiscing semper elit. Proin fermentum massa ac quam. Sed diam turpis, molestie vitae, placerat a, molestie nec, leo. Maecenas lacinia. Nam ipsum ligula, eleifend at, accumsan nec, suscipit a, ipsum. Morbi blandit ligula feugiat magna. Nunc eleifend consequat lorem. Sed lacinia nulla vitae enim. Pellentesque tincidunt purus vel magna. Integer non enim. Praesent euismod nunc eu purus. Donec bibendum quam in tellus. Nullam cursus pulvinar lectus. Donec et mi. Nam vulputate metus eu enim. Vestibulum pellentesque felis eu massa.

Quisque ullamcorper placerat ipsum. Cras nibh. Morbi vel justo vitae lacus tincidunt ultrices. Lorem ipsum dolor sit amet, consectetur adipiscing elit. In hac habitasse platea dictumst. Integer tempus convallis augue. Etiam facilisis. Nunc elementum fermentum wisi. Aenean placerat. Ut imperdiet, enim sed gravida sollicitudin, felis odio placerat quam, ac pulvinar elit purus eget enim. Nunc vitae tortor. Proin tempus nibh sit amet nisl. Vivamus quis tortor vitae risus porta vehicula.

Fusce mauris. Vestibulum luctus nibh at lectus. Sed bibendum, nulla a faucibus semper, leo velit ultricies tellus, ac venenatis arcu wisi vel nisl. Vestibulum diam. Aliquam pellentesque, augue quis sagittis posuere, turpis lacus congue quam, in hendrerit risus eros eget felis. Maecenas eget erat in sapien mattis porttitor. Vestibulum porttitor. Nulla facilisi. Sed a turpis eu lacus commodo facilisis. Morbi fringilla, wisi in dignissim interdum, justo lectus sagittis dui, et vehicula libero dui cursus dui. Mauris tempor ligula sed lacus. Duis cursus enim ut augue. Cras ac magna. Cras nulla. Nulla egestas. Curabitur a leo. Quisque egestas wisi eget nunc. Nam feugiat lacus vel est. Curabitur consectetur.

Suspendisse vel felis. Ut lorem lorem, interdum eu, tincidunt sit amet, laoreet vitae, arcu. Aenean faucibus pede eu ante. Praesent enim elit, rutrum at, molestie non, nonummy vel, nisl. Ut lectus eros, malesuada sit amet, fermentum eu, sodales cursus, magna. Donec eu purus. Quisque vehicula, urna sed ultricies auctor, pede lorem egestas dui, et convallis elit erat sed nulla. Donec luctus. Curabitur et nunc. Aliquam dolor odio, commodo pretium, ultricies non, pharetra in, velit. Integer arcu est, nonummy in, fermentum faucibus, egestas vel, odio.

Sed commodo posuere pede. Mauris ut est. Ut quis purus. Sed ac odio. Sed vehicula hendrerit sem. Duis non odio. Morbi ut dui. Sed accumsan risus eget odio. In hac habitasse platea dictumst. Pellentesque non elit. Fusce sed justo eu urna porta tincidunt. Mauris felis odio, sollicitudin sed, volutpat a, ornare ac, erat. Morbi quis dolor. Donec pellentesque, erat ac sagittis semper, nunc dui lobortis purus, quis congue purus metus ultricies tellus. Proin et quam. Class aptent taciti sociosqu ad litora torquent per conubia nostra, per inceptos hymenaeos. Praesent sapien turpis, fermentum vel, eleifend faucibus, vehicula eu, lacus.

5.2. Diseño

Lorem ipsum dolor sit amet, consectetur adipiscing elit. Ut purus elit, vestibulum ut, placerat ac, adipiscing vitae, felis. Curabitur dictum gravida mauris. Nam arcu libero, nonummy eget, consectetur id, vulputate a, magna. Donec vehicula augue eu neque. Pellentesque habitant morbi tristique senectus et netus et malesuada fames ac turpis egestas. Mauris ut leo. Cras viverra metus rhoncus sem. Nulla et lectus vestibulum urna fringilla ultrices. Phasellus eu tellus sit amet tortor gravida placerat. Integer sapien est, iaculis in, pretium quis, viverra ac, nunc. Praesent eget sem vel leo ultrices bibendum. Aenean faucibus. Morbi dolor nulla, malesuada eu, pulvinar at, mollis ac, nulla. Curabitur auctor semper nulla. Donec varius orci eget risus. Duis nibh mi, congue eu, accumsan eleifend, sagittis quis, diam. Duis eget orci sit amet orci dignissim rutrum.

Nam dui ligula, fringilla a, euismod sodales, sollicitudin vel, wisi. Morbi auctor lorem non justo. Nam lacus libero, pretium at, lobortis vitae, ultricies et, tellus. Donec aliquet, tortor sed accumsan bibendum, erat ligula aliquet magna, vitae ornare odio metus a mi. Morbi ac orci et nisl hendrerit mollis. Suspendisse ut massa. Cras nec ante. Pellentesque a nulla. Cum sociis natoque penatibus et magnis dis parturient montes, nascetur ridiculus mus. Aliquam tincidunt urna. Nulla ullamcorper vestibulum turpis. Pellentesque cursus luctus mauris.

Nulla malesuada porttitor diam. Donec felis erat, congue non, volutpat at, tincidunt tristique, libero. Vivamus viverra fermentum felis. Donec nonummy pellentesque ante.

Phasellus adipiscing semper elit. Proin fermentum massa ac quam. Sed diam turpis, molestie vitae, placerat a, molestie nec, leo. Maecenas lacinia. Nam ipsum ligula, eleifend at, accumsan nec, suscipit a, ipsum. Morbi blandit ligula feugiat magna. Nunc eleifend consequat lorem. Sed lacinia nulla vitae enim. Pellentesque tincidunt purus vel magna. Integer non enim. Praesent euismod nunc eu purus. Donec bibendum quam in tellus. Nullam cursus pulvinar lectus. Donec et mi. Nam vulputate metus eu enim. Vestibulum pellentesque felis eu massa.

Quisque ullamcorper placerat ipsum. Cras nibh. Morbi vel justo vitae lacus tincidunt ultrices. Lorem ipsum dolor sit amet, consectetur adipiscing elit. In hac habitasse platea dictumst. Integer tempus convallis augue. Etiam facilisis. Nunc elementum fermentum wisi. Aenean placerat. Ut imperdiet, enim sed gravida sollicitudin, felis odio placerat quam, ac pulvinar elit purus eget enim. Nunc vitae tortor. Proin tempus nibh sit amet nisl. Vivamus quis tortor vitae risus porta vehicula.

Fusce mauris. Vestibulum luctus nibh at lectus. Sed bibendum, nulla a faucibus semper, leo velit ultricies tellus, ac venenatis arcu wisi vel nisl. Vestibulum diam. Aliquam pellentesque, augue quis sagittis posuere, turpis lacus congue quam, in hendrerit risus eros eget felis. Maecenas eget erat in sapien mattis porttitor. Vestibulum porttitor. Nulla facilisi. Sed a turpis eu lacus commodo facilisis. Morbi fringilla, wisi in dignissim interdum, justo lectus sagittis dui, et vehicula libero dui cursus dui. Mauris tempor ligula sed lacus. Duis cursus enim ut augue. Cras ac magna. Cras nulla. Nulla egestas. Curabitur a leo. Quisque egestas wisi eget nunc. Nam feugiat lacus vel est. Curabitur consectetur.

Suspendisse vel felis. Ut lorem lorem, interdum eu, tincidunt sit amet, laoreet vitae, arcu. Aenean faucibus pede eu ante. Praesent enim elit, rutrum at, molestie non, nonummy vel, nisl. Ut lectus eros, malesuada sit amet, fermentum eu, sodales cursus, magna. Donec eu purus. Quisque vehicula, urna sed ultricies auctor, pede lorem egestas dui, et convallis elit erat sed nulla. Donec luctus. Curabitur et nunc. Aliquam dolor odio, commodo pretium, ultricies non, pharetra in, velit. Integer arcu est, nonummy in, fermentum faucibus, egestas vel, odio.

Sed commodo posuere pede. Mauris ut est. Ut quis purus. Sed ac odio. Sed vehicula hendrerit sem. Duis non odio. Morbi ut dui. Sed accumsan risus eget odio. In hac habitasse platea dictumst. Pellentesque non elit. Fusce sed justo eu urna porta tincidunt. Mauris felis odio, sollicitudin sed, volutpat a, ornare ac, erat. Morbi quis dolor. Donec pellentesque, erat ac sagittis semper, nunc dui lobortis purus, quis congue purus metus ultricies tellus. Proin et quam. Class aptent taciti sociosqu ad litora torquent per conubia nostra, per inceptos hymenaeos. Praesent sapien turpis, fermentum vel, eleifend faucibus, vehicula eu, lacus.

5.3. Resultados

Lorem ipsum dolor sit amet, consectetur adipiscing elit. Ut purus elit, vestibulum ut, placerat ac, adipiscing vitae, felis. Curabitur dictum gravida mauris. Nam arcu libero, nonummy eget, consectetur id, vulputate a, magna. Donec vehicula augue eu neque. Pellentesque habitant morbi tristique senectus et netus et malesuada fames ac turpis egestas. Mauris ut leo. Cras viverra metus rhoncus sem. Nulla et lectus vestibulum urna fringilla ultrices. Phasellus eu tellus sit amet tortor gravida placerat. Integer sapien est, iaculis in, pretium quis, viverra ac, nunc. Praesent eget sem vel leo ultrices bibendum. Aenean faucibus. Morbi dolor nulla, malesuada eu, pulvinar at, mollis ac,

nulla. Curabitur auctor semper nulla. Donec varius orci eget risus. Duis nibh mi, congue eu, accumsan eleifend, sagittis quis, diam. Duis eget orci sit amet orci dignissim rutrum.

Nam dui ligula, fringilla a, euismod sodales, sollicitudin vel, wisi. Morbi auctor lorem non justo. Nam lacus libero, pretium at, lobortis vitae, ultricies et, tellus. Donec aliquet, tortor sed accumsan bibendum, erat ligula aliquet magna, vitae ornare odio metus a mi. Morbi ac orci et nisl hendrerit mollis. Suspendisse ut massa. Cras nec ante. Pellentesque a nulla. Cum sociis natoque penatibus et magnis dis parturient montes, nascetur ridiculus mus. Aliquam tincidunt urna. Nulla ullamcorper vestibulum turpis. Pellentesque cursus luctus mauris.

Nulla malesuada porttitor diam. Donec felis erat, congue non, volutpat at, tincidunt tristique, libero. Vivamus viverra fermentum felis. Donec nonummy pellentesque ante. Phasellus adipiscing semper elit. Proin fermentum massa ac quam. Sed diam turpis, molestie vitae, placerat a, molestie nec, leo. Maecenas lacinia. Nam ipsum ligula, eleifend at, accumsan nec, suscipit a, ipsum. Morbi blandit ligula feugiat magna. Nunc eleifend consequat lorem. Sed lacinia nulla vitae enim. Pellentesque tincidunt purus vel magna. Integer non enim. Praesent euismod nunc eu purus. Donec bibendum quam in tellus. Nullam cursus pulvinar lectus. Donec et mi. Nam vulputate metus eu enim. Vestibulum pellentesque felis eu massa.

Quisque ullamcorper placerat ipsum. Cras nibh. Morbi vel justo vitae lacus tincidunt ultrices. Lorem ipsum dolor sit amet, consectetur adipiscing elit. In hac habitasse platea dictumst. Integer tempus convallis augue. Etiam facilisis. Nunc elementum fermentum wisi. Aenean placerat. Ut imperdiet, enim sed gravida sollicitudin, felis odio placerat quam, ac pulvinar elit purus eget enim. Nunc vitae tortor. Proin tempus nibh sit amet nisl. Vivamus quis tortor vitae risus porta vehicula.

Fusce mauris. Vestibulum luctus nibh at lectus. Sed bibendum, nulla a faucibus semper, leo velit ultricies tellus, ac venenatis arcu wisi vel nisl. Vestibulum diam. Aliquam pellentesque, augue quis sagittis posuere, turpis lacus congue quam, in hendrerit risus eros eget felis. Maecenas eget erat in sapien mattis porttitor. Vestibulum porttitor. Nulla facilisi. Sed a turpis eu lacus commodo facilisis. Morbi fringilla, wisi in dignissim interdum, justo lectus sagittis dui, et vehicula libero dui cursus dui. Mauris tempor ligula sed lacus. Duis cursus enim ut augue. Cras ac magna. Cras nulla. Nulla egestas. Curabitur a leo. Quisque egestas wisi eget nunc. Nam feugiat lacus vel est. Curabitur consectetur.

Suspendisse vel felis. Ut lorem lorem, interdum eu, tincidunt sit amet, laoreet vitae, arcu. Aenean faucibus pede eu ante. Praesent enim elit, rutrum at, molestie non, nonummy vel, nisl. Ut lectus eros, malesuada sit amet, fermentum eu, sodales cursus, magna. Donec eu purus. Quisque vehicula, urna sed ultricies auctor, pede lorem egestas dui, et convallis elit erat sed nulla. Donec luctus. Curabitur et nunc. Aliquam dolor odio, commodo pretium, ultricies non, pharetra in, velit. Integer arcu est, nonummy in, fermentum faucibus, egestas vel, odio.

Sed commodo posuere pede. Mauris ut est. Ut quis purus. Sed ac odio. Sed vehicula hendrerit sem. Duis non odio. Morbi ut dui. Sed accumsan risus eget odio. In hac habitasse platea dictumst. Pellentesque non elit. Fusce sed justo eu urna porta tincidunt. Mauris felis odio, sollicitudin sed, volutpat a, ornare ac, erat. Morbi quis dolor. Donec pellentesque, erat ac sagittis semper, nunc dui lobortis purus, quis congue purus metus ultricies tellus. Proin et quam. Class aptent taciti sociosqu ad litora torquent per conubia

nostra, per inceptos hymenaeos. Praesent sapien turpis, fermentum vel, eleifend faucibus, vehicula eu, lacus.

Glosario

Símbolos

Clustering

Pendiente. 7

Coursera

Sitio web de cursos de aprendizaje en <https://www.coursera.org>. 3

Dataset

Pendiente. 1

Feedforward

Pendiente. 7

Maximum Margin Hyperplanes

Hiperplanos que permiten separar datos en espacios de alta dimensionalidad con un margen asociado para separarlos.. 10

Unsupervised Learning

Pendiente. 7

Artificial Intelligence (AI)

Pendiente. 6

Data Integration (DI)

Para acceder a múltiples y diversas fuentes de información. 4

Data Mining (DM)

Proceso de descubrir automáticamente información útil en repositorios grandes de datos. 5

Deductive Systems (DS)

Pendiente. 7

Expert Systems (ES)

Pendiente. 6

Inference Engine (IE)

Pendiente. 6

*Knowledge Based Systems (KBS)***Pendiente.** 6*Knowledge Base (KB)***Pendiente.** 6*Link Analysis (LA)*

Para visualizar asociaciones y relaciones criminales y terroristas. 4

Machine Learning (ML)

Informalmente ha sido definido como “El campo de estudio que le da a computadores la habilidad de aprender sin ser explícitamente programados”, este tiene tres tipos de algoritmos de aprendizaje: aprendizaje supervisado, aprendizaje no-supervisado, y aprendizaje por refuerzo. 6

Machine Learning Algorithms (MLA)

Para extraer perfiles de perpetradores y mapas gráficos de crímenes. 4

Artificial Neural Network (ANN)

Para predecir la probabilidad de crímenes y nuevos ataques terroristas. 4, 7

Natural Language Processing (NLP)

Rama de la inteligencia artificial que lidia con la interacción entre computadores y humanos usando el lenguaje natural. 1–4

Open Source Intelligence (OSINT)

Disciplina responsable de la adquisición, procesamiento y posterior transformación en inteligencia de información obtenida de fuentes públicas como prensa, radio, televisión, internet, informes de diferentes sectores y, en general, cualquier recurso de acceso público (Tomado de [2]). 3

*Reactive Systems (RS)***Pendiente.** 7*Rule-based Systems (RBS)***Pendiente.** 6, 7*Software Agents (SA)*

Para el monitoreo, obtención, análisis y actuación sobre la información. 4

*Self-organizing Maps (SOM)***Pendiente.** 7, 8*Support Vector Machine (SVM)***Pendiente.** 10, 11

Text Mining (TM)

Búsqueda sobre terabytes de información en documentos, paginas web y correos electrónicos. 4

Working Memory (WM)

Pendiente. 6, 7

Bibliografía

- [1] Leandro Nunes De Castro. *Fundamentals of natural computing: basic concepts, algorithms, and applications*. Chapman y Hall/CRC, 2006.
- [2] Martín José Hernández Medina, Ricardo Andrés Pinto Rico y Cristian Camilo Pinzón Hernández. *Inteligencia de fuentes abiertas para el contexto colombiano*. Inf. téc. Escuela Colombiana de Ingeniería Julio Garavito, 2018.
- [3] J. Mena. *Investigative Data Mining for Security and Criminal Detection*. Elsevier Science, 2003. ISBN: 9780080509389. URL: <https://books.google.com.co/books?id=3mDlrtJuZv4C>.
- [4] Jerry M Mendel. *Uncertain Rule-Based Fuzzy Systems*. ISBN: 9783319513690. DOI: [10.1007/978-3-319-51370-6](https://doi.org/10.1007/978-3-319-51370-6).
- [5] Ana Maria Popescu y Orena Etzioni. «Extracting product features and opinions from reviews». En: *Natural Language Processing and Text Mining* (2007), págs. 9-28. ISSN: 08247935. DOI: [10.1007/978-1-84628-754-1_2](https://doi.org/10.1007/978-1-84628-754-1_2). arXiv: [0309034](https://arxiv.org/abs/0309034) [cs].
- [6] Priti Srinivas Sajja y Rajendra Akerkar. «Knowledge-based systems for development». En: *Advanced Knowledge Based Systems: Model, Applications & Research 1* (2010), págs. 1-11.
- [7] Pang-Ning Tan, Michael Steinbach y Vipin Kumar. *Introduction to Data Mining*. US ed. Addison Wesley, mayo de 2005. ISBN: 0321321367. URL: <http://www.amazon.com/exec/obidos/redirect?tag=citeulike07-20&path=ASIN/0321321367>.