

# Clasificador de Noticias usando Autoencoders

Gonzalo Farias, Sebastián Vergara, Ernesto Fabregas, Gabriel Hermosilla, Sebastián Dormido-Canto, and Sebastián Dormido

**Abstract**—This article presents a classification system for news with Deep Learning. With this tool the news are classified in the following categories: Sports, Politics, Economics, Show and Police. Also they receives an scope: Local (Valparaíso), National (Chile) and International (Rest of the World). The classifiers were built using a database with 542 news labeled with the previous criteria. The features were extracted with Autoencoders (AE) to train an Artificial Neural Network (ANN) of multiple classes *Softmax* (*Softmax* ANNs). Both stages were stacked following the concept of Deep Learning. The results with the data test (156 news) reach a success rate of 92.3% for the category classifier and 87.2% for the scope classifier. The general success rate for both, category and scope was 83.75%.

**Keywords**—News Classifier, Autoencoder

## I. INTRODUCTION

La clasificación automática se puede resumir como el ejercicio de separar algún conjunto de elementos específicos a través de un sistema artificial (computador), con el fin de otorgarle una categoría o clase. Los elementos a clasificar pueden ser una gran variedad de tipos, tales como imágenes, videos, sonidos, señales físicas (temperatura, voltaje, corriente, presión, etc.) y texto.

Respecto a este último tipo de dato se han desarrollado una gran cantidad de algoritmos que hacen la extracción de características y clasificación bajo diferentes criterios [1]. Lo anterior se debe a que los textos generan diariamente grandes volúmenes de información, los cuales se almacenan en bases de datos de páginas web, redes sociales y de empresas, por lo que discriminarla puede ser muy provechoso.

Existen algunos ejemplos de estudios realizados para la discriminación de textos, como por ejemplo [2], donde se muestra un sistema clasificador de sentimientos de una base de datos de *tweets* mediante distintos algoritmos de aprendizaje automático: Máquina de Vectores Soporte (SVM) [3], Clasificador Bayesiano [4] y Clasificador de Máxima Entropía [5]. Por otro lado, en [6] se implementa un sistema capaz de discriminar cuando estamos en presencia de un *tweet* de carácter ofensivo mediante Redes Neuronales Recurrentes (RNNs). Con respecto a la extracción de características, en [7] se realizó un estudio comparativo empírico sobre diferentes técnicas que abordan esta problemática.

En el área de los textos noticiosos, se han realizado una serie de investigaciones con la finalidad de discriminarlos bajo diferentes categorías, en [13] se construye un sistema basado en Word2Vec; el cual es un modelo de redes neuronales que

busca otorgar a cada palabra un vector en un espacio generalmente de cientos de dimensiones y teniendo en cuenta además, que las palabras que compartan contextos comunes se encuentran a distancia menores. El método de clasificación en este artículo se basa en algoritmos de Deep Learning. En [14] se presenta un sistema para el análisis de temas de artículos periodísticos, utilizando un clasificador binario basado en Máquina de Vectores Soporte y Clasificador Bayesiano, el algoritmo es entrenado a partir de un vector generado a través de la cuenta de palabras claves extraídas de manera automática. También [15] se realiza un trabajo similar utilizando el método de clasificación *Naive Bayes*.

Particularmente en el presente artículo se busca implementar un sistema de clasificación de noticias según las categorías Deporte, Política, Economía, Espectáculo y Policial, y de acuerdo al ámbito que pertenecen, Local (Región Valparaíso), Nacional (Chile) e Internacional. El clasificador diseñado se basa en Autoencoders (AE) y una Red Neuronal Artificial con capa de salida *Softmax* (ANNs *Softmax*).

Un sistema clasificador de noticias puede ser una herramienta atractiva para usuarios que busquen informarse solamente de tópicos específicos, ya que de esta manera el sistema se encarga de filtrar toda la información que no es de interés del usuario y así tener un acceso a las noticias de forma más amigable.

La estructura de este trabajo inicia con un marco teórico, que busca introducir a los algoritmos utilizados, a continuación, se detalla el funcionamiento del clasificador propuesto y finalmente se exponen los resultados obtenidos.

## II. MARCO TEÓRICO

El *deep learning* (aprendizaje profundo) ha surgido como uno de los enfoques más utilizados en los últimos años ya que es capaz de resolver problemas en muchas áreas, tanto científicas como cotidianas donde otras técnicas de *machine learning* tienen limitaciones. En este apartado se explican de manera general los algoritmos de *deep learning* utilizados en el sistema clasificador propuesto [8].

### A. Red Neuronal Artificial Como multclasificador

El aprendizaje se define como el proceso de adquisición de conocimiento a través de la experiencia, estudio o práctica. Una vez terminado el proceso, el sujeto es capaz de distinguir de acuerdo a lo aprendido. Esta capacidad está asociada directamente con los humanos, pero gracias a los avances tecnológicos, se han logrado encontrar algoritmos capaces de realizar un proceso análogo, pero mediante computadores (máquinas). En otras palabras, el aprendizaje automático busca la capacidad de clasificar de la forma más general y robusta posible diferentes clases y a través de variados tipos de algoritmos como SVM,

G. Farias, S. Vergara and G. Hermosilla are with the Escuela de Ingeniería Eléctrica. Pontificia Universidad Católica de Valparaíso. Valparaíso, Chile e-mail: gfarias,gabriel.hermosilla@pucv.cl, sebastianiovi@gmail.com

E. Fabregas, S. Dormido-Canto and S. Dormido are with Departamento de Informática y Automática. Universidad Nacional de Educación a Distancia. Madrid, Spain email: efabregas@bec.uned.es, sebas,sdormido@dia.uned.es

árboles de decisión, reglas de asociación, redes bayesianas, algoritmos de agrupamiento, ANNs, entre otros.

Las redes neuronales artificiales en particular, son modelos matemáticos que intentan reproducir el comportamiento del sistema nervioso humano. Se conforman a partir de la interconexión de unidades más básicas denominadas perceptrón (o neurona artificial), las cuales son una simplificación matemática de una neurona biológica. Cuando un perceptrón recibe un dato de entrada, este lo pondera (atenúa, amplifica o inhibe) a través de un parámetro llamado peso, asociado a la sinapsis biológica y a su salida se evalúa una función de activación que determina si la neurona está excitada o no [9]. La figura 1 muestra la arquitectura de una red neuronal de 4 capas: la capa de entrada formada por 4 neuronas, las dos capas ocultas formadas por 3 neuronas cada una y la capa de salida formada por una única neurona.

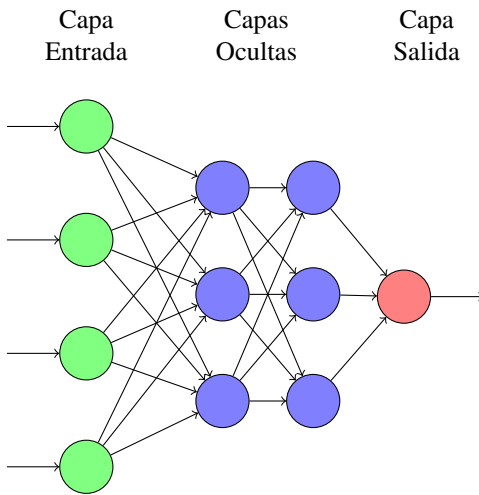


Fig. 1. Red Neuronal Artificial (ANN).

El apilamiento de capas de perceptrones interconectados se entrenan a partir de datos conocidos, es decir, se ingresa un conjunto de datos y se obtiene un resultado a la salida de la red, dicho resultado se compara con el deseado y según el error obtenido se modifican los pesos de cada perceptrón. El algoritmo utilizado para el entrenamiento de la red se denomina *Backpropagation*, el cual funciona en dos etapas. Primeramente, al excitar las neuronas de entrada con un conjunto de datos, estos se propagan a través de la red hasta generar una salida, lo obtenido se compara con la salida deseada y se calcula el error de forma individual. En una segunda etapa el error obtenido se propaga hacia atrás, partiendo por las capas superiores hasta las iniciales, modificando los pesos de la manera más equilibrada posible, con la finalidad de que cada neurona aprenda una característica específica de los datos de entrada. Al realizar el proceso anterior de manera repetitiva y con diferentes muestras, la red es capaz de aprender a clasificar los datos. Cabe destacar el aumento de las capas ocultas en una red neuronal le dan el carácter de profundo (deep).

### B. Clasificador

Existen dos enfoques respecto a la cantidad de clases que se desean clasificar, por un lado, un clasificador busca

separar en dos categorías y un multclasificador en tres o más. Este último puede ser más complicado de entrenar ya que necesita diferenciar más tipos de clases (problema complejo) en comparación al primero. Es importante destacar que a partir de varios clasificadores se puede obtener los resultados de uno múltiple. En particular para el sistema propuesto se utiliza una capa de salida tipo *Softmax* multiclase.

*Softmax*, es una función exponencial normalizada utilizada para otorgar una probabilidad sobre un número determinado de salidas que se emplean como capa final en clasificadores basados en redes neuronales. A cada salida de las neuronas de la última capa se le asignan valores dentro del intervalo  $[0,1]$  y según la neurona que entrega el valor máximo se le otorga una clase a los datos de entrada.

### C. Autoencoder

Las características de un elemento aluden a las cualidades (información) que lo definen. Mientras que el agrupamiento de estas conforman el vector de características y según el tamaño de este último se define la dimensionalidad del problema. Una excesiva dimensionalidad puede generar problemas de sobreajuste y un perjuicio en la velocidad de procesamiento, por lo que reducirla es una etapa fundamental en el proceso de clasificación [10].

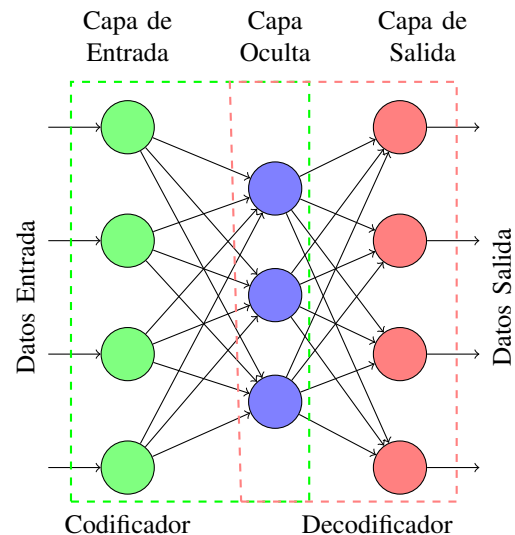


Fig. 2. Esquema de un Autoencoder.

Un *Autoencoder* (AE) tiene como objetivo extraer características de manera automática y sin supervisión mediante una red neuronal artificial capaz de aprender una representación codificada para un conjunto de datos. Se conforma típicamente de tres capas de neuronas artificiales, una de entrada, una de salida y una capa oculta que las conecta. Se buscan los pesos de la red con la finalidad de minimizar la diferencia entre la salida y entrada del AE [11]. La figura 2 muestra gráficamente la arquitectura de un AE, con sus capas y las etapas de codificación y decodificación.

## III. CLASIFICADOR PROPUESTO

En esta sección se detalla el funcionamiento del sistema clasificador de noticias de acuerdo a sus diferentes etapas.

### A. Matriz de Frecuencia

Para el entrenamiento del sistema propuesto como primera etapa se construye la matriz de frecuencia. Esta matriz se genera a partir de un listado que contiene palabras clave respecto a cada **Categoría** y **mbito** de una noticia, las longitudes de las listas son de 1351 y 894 palabras respectivamente y se componen como se indica en la tabla I.

TABLE I  
CANTIDAD DE PALABRAS POR CATEGORIA Y MBITO\*

Categoría o mbito	Cantidad de Palabras
Deporte	290
Política	314
Economía	247
Policía	199
Espectáculo	301
Local*	135
Nacional*	414
Internacional*	345

Cabe destacar que la elección de las palabras claves se realiza de manera manual, es decir, que se inspecciona un conjunto de noticias y se seleccionan aquellas palabras que tengan relación con el tema, por ejemplo para la categoría de Política, palabras como presidente, ministro y diputado forman parte de la lista. Conjuntamente se agregan palabras directamente sin la necesidad de ser extraídas de noticias.

En este punto finalmente se cuenta la cantidad de veces que aparece una palabra específica en una noticia. El resultado debe ser que para cada tipo de noticia, la frecuencia de palabras de igual clase, debe ser mayor. En la tabla II se muestra de manera ilustrativa la matriz de frecuencia para el listado de **Categoría** para noticias de diferentes clases. En nuestro caso, la matriz tiene una dimensionalidad de 542x1351. La matriz de frecuencia para el listado de **mbito** es análoga.

TABLE II  
EJEMPLO DE MATRIZ DE FRECUENCIA

Noticia	D	P	E	Po	Es
Deporte(D)	25	2	1	3	0
Economía(E)	2	0	32	12	5
Policía(Po)	0	1	1	15	3
Espectáculo (Es)	5	3	2	7	24
Política (P)	1	33	15	2	0

Cabe destacar que el proceso de elección de palabras es sensible al sistema y debe funcionar de la manera más robusta posible. Además su determinación debe ser actualizado empíricamente según los resultados obtenidos en la clasificación. Además es fundamental contar con una base de datos de noticias lo más amplia posible. Siguiendo estas premisas, se construyeron dos conjuntos de noticias conformados por 542 para el clasificador de categorías y 395 para el de mbito, subdivididas de acuerdo a lo que se indica en la tabla III.

Una vez construida la matriz de frecuencia con los datos de entrenamiento, se continúa con el proceso de extracción de características.

### B. Extracción de Características

Como se comentó anteriormente, el proceso de extracción de características busca disminuir la dimensionalidad del prob-

TABLE III  
CANTIDAD DE NOTICIAS POR CATEGORIA Y MBITO\*

Categoría o mbito	Cantidad de Noticias
Deporte	142
Política	109
Economía	93
Policía	99
Espectáculo	99
Local*	107
Nacional*	139
Internacional*	149

lema para generar una clasificación más robusta y con menor exigencia en el procesamiento. La dimensionalidad del vector de características de cada noticia respecto a la categoría es de 1351. Mediante la implementación del autoencoder se redujo a 100 características, disminuyendo la dimensionalidad en un factor de 13.51. Con respecto al mbito de la noticia, el número de características es de 894 y a través de otro AE se redujo la dimensionalidad a 45, por lo que el factor de reducción es de 19.86.

### C. Clasificadores Implementados

Los clasificadores implementados se conforman a partir de una única capa tipo *Softmax* de 5 neuronas para el clasificador de categorías (una neurona por categoría) y 3 para el de mbito (una neurona para cada mbito), sus entrenamientos se realizaron a partir de las características extraídas en el proceso anterior.

Como última etapa se apilan los autoencoders de categoría y mbito con la capa *Softmax* respectiva y se realiza un último entrenamiento de refinamiento. En la siguiente figura se ilustra el sistema de clasificación construido según la categoría (análogo para el clasificador de mbito). La figura 3 muestra la implementación del sistema clasificador con el *Neural Network Toolbox* de MATLAB.

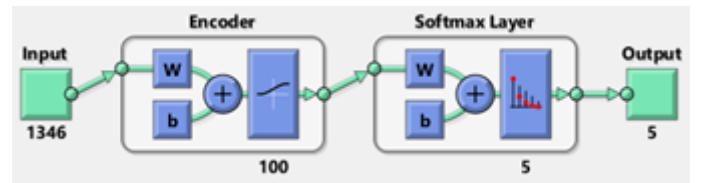


Fig. 3. ANNs diseñado con el *Neural Network Toolbox* de MATLAB.

Es fundamental señalar que la entrada al sistema clasificador es el código HTML que posee la noticia en su totalidad, en la figura 4 se muestra el diagrama de bloques.

## IV. RESULTADOS

Para evaluar el sistema construido se realizaron tres pruebas, en primer lugar, se evaluó el clasificador de Categoría, seguido del de mbito para finalmente realizar un test general del clasificador de noticias. Para representar los resultados obtenidos de cada clasificador se utiliza la matriz de confusión [12]. En la tabla IV se muestra la cantidad de noticias utilizadas en el test. Las filas representan las categorías y las columnas el mbito.

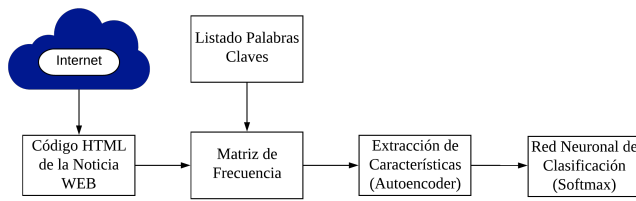


Fig. 4. Diagrama de bloque Sistema Clasificador de Noticias.

TABLE IV  
CONJUNTO DE NOTICIAS DE TEST POR CATEGORIA Y MBITO

Categoría	mbito			
	Local	Nacional	Internacional	Total
Deporte	11	11	11	33
Política	10	10	12	32
Economía	7	10	11	28
Policial	12	10	10	32
Espectáculo	11	10	10	31
Total	51	51	54	156

#### A. Matriz de confusin

La matriz de confusin es muy utilizada en sistemas de clasificacin, en ella se muestra ms detalladamente cmo se predijeron las muestras. La figura 5 muestra un ejemplo de una matriz de confusin.

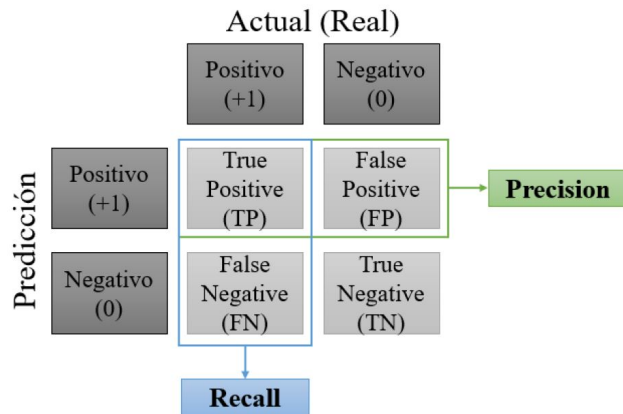


Fig. 5. Matriz de confusin, en la seccin superior o inferior los valores reales de las noticias. En la seccin izquierda, las predicciones realizadas por el algoritmo de clasificacin.

Dicha matriz ilustra cmo se proporciona informacin relevante bajo los siguientes cuatro conceptos:

- **TP: True Positive** (Verdadero Positivo), cuando la clase predicha y la real son iguales.
- **FP: False Positive** (Falso Positivo), cuando se predice como +1 y corresponde a 0.
- **FN: False Negative** (Falso Negativo), se predice como clase 0 pero pertenece a +1.
- **TN: True Negative** (Verdadero Negativo), cuando el clasificador predice como clase 0 y efectivamente pertenece a ella.

De esta matriz, se puede extraer informacin realizando los siguientes dos anlisis:

- 1) Horizontal: Cuntas muestras predichas en la clase  $i$  (fila) pertenece a la clase  $j$  (columna).
- 2) Vertical: Cuntas muestras que efectivamente pertenecen a la clase  $j$  (columna) fueron catalogadas como  $i$  (fila).

En las figuras 6 y 7 se muestran las matrices de confusin obtenidas a partir del clasificador de categoria y mbito respectivamente.

Confusion Matrix					
Output Class	1	2	3	4	5
	32	0	0	0	0
	0	27	2	2	0
	0	2	25	0	0
	0	3	1	30	1
	1	0	0	0	30
Target Class					
1	2	3	4	5	
97.0%	84.4%	89.3%	93.8%	96.8%	92.3%
3.0%	15.6%	10.7%	6.3%	3.2%	7.7%

Fig. 6. Matriz de Confusin clasificador de categoria (1 Deporte, 2 Poltica, 3 Economa, 4 Policial, 5 Espectculo).

De la figura anterior se puede extraer que la tasa de acierto para este clasificador de Categoria alcanz un porcentaje del 92.3%. Los errores generados en la clasificacin se deben a diferentes fenmenos, como por ejemplo:

- Ambigüedad en la clasificacin: esto se ve reflejado en los errores cometidos en noticias polticas y econmicas, ya que stas suelen estar muy relacionadas entre s.
- Ruido en la informacin: las noticias en sitios web muestran generalmente la informacin principal y en su alrededor otros tipos de noticias, esto genera problemas al ahora de clasificar.
- Ruido en el cdigo fuente de la pgina: hay palabras, generalmente con tilde, que al buscarlas en la noticias no se detectan debido que en el cdigo de la pgina estn escritas de diferente manera (poltica y pol<i>tica</i>).

Para combatir los tres puntos anteriores, realizar un preprocesamiento con la finalidad de obtener el cuerpo de la noticia sin otra informacin sera una buena manera de perfeccionar el sistema. En trabajos posteriores se evaluar la inclusin de una etapa de preprocesamiento.

De acuerdo a la matriz de confusin tambin se pueden extraer otros resultados. Por ejemplo, de las 5 categoras clasificadas, con las que se obtiene un mejor desempeo son: Deporte y Espectculo (clases 1 y 5), con un 100% y 96.8% respectivamente. Lo anterior concuerda con lo esperado ya

que este tipo de noticias suelen ser muy diferentes a las otras clases, en cambio para Política, Economía y Policial existe un grado de cercanía, por ejemplo: *son muy comunes noticias donde un Presidente comenta sobre proyectos de ley de carácter policial o económico.*

Confusion Matrix				
Output Class	1	2	3	
	40	1	0	97.6% 2.4%
	11	48	6	73.8% 26.2%
	0	2	48	96.0% 4.0%
	1	2	3	
Target Class				
	78.4% 21.6%	94.1% 5.9%	88.9% 11.1%	87.2% 12.8%

Fig. 7. Matriz de Confusión clasificador de ámbito (1 Local, 2 Nacional, 3 Internacional).

Para el clasificador de ámbito se obtuvo una tasa de acierto de 87.2%, logrando resultados inferiores, aunque cercanos al clasificador anterior y los fenómenos que generan el error son básicamente los mismos.

En este caso el mayor error en la clasificación se genera en las noticias de ámbito nacional, obteniendo una tasa de Falsos Positivos de 26.2%, esto se puede deber a que la gran mayoría de las noticias con que se construye el clasificador son de sitios web nacionales y por ende aparecen constantemente palabras claves relacionadas con este ámbito. Por otro lado, para los ámbitos Local e Internacional se alcanzan aciertos en torno al 96%, ambas tienen palabras muy características, como son el nombre de ciudades locales y de países o ciudades extranjeras.

El comportamiento del sistema general alcanza un porcentaje de acierto de 83.75%, es decir, de las 156 noticias de test, 130 fueron bien predichas tanto en su categoría como ámbito. Cabe destacar que la tasa de acierto total del sistema no corresponde a la multiplicación entre las tasas de acierto de cada clasificador, ya que hay casos donde el error se produce en una misma noticia, de manera contraria se estará contando dos veces un mismo error.

## V. CONCLUSIÓN

Diariamente se generan una gran cantidad de datos de diferente índole, por lo que categorizarlos y extraer información de ellos puede ser muy provechoso. En particular, la clasificación de noticias puede funcionar para generar aplicaciones que las

recopile de diferentes sitios, las ordene en categorías. De esta forma el usuario puede tener un mejor acceso a la actualidad local, nacional o internacional.

En este trabajo se presenta la implementación de dos clasificadores basados en *deep learning*, para determinar la categoría y el ámbito de una noticia extraída de Internet. El proceso de clasificación se genera a través de la construcción de matrices de frecuencias mediante de una lista de palabras clave. Luego se extraen características a través de *Autoencoders* para finalmente realizar la clasificación con una capa final de neuronas artificiales tipo *Softmax*.

Las noticias se clasifican según su Categoría: Deporte, Política, Economía, Espectáculo y Policial; y según el ámbito: Local (Región Valparaíso), Nacional (Chile) o Internacional (resto del mundo). Los resultados obtenidos alcanzaron tasas de acierto del 92.3% para el clasificador de Categoría y 87.2% para el de ámbito. Las noticias bien predichas para ambos casos lograron un porcentaje de acierto de 83.75%.

Para mejorar las tasas del sistema es necesario refinar el listado de palabras claves y ampliar la base de datos de noticias. Además para este último se requiere abarcar la mayor cantidad de sitios web y así generar un sistema más robusto. También se puede filtrar el ruido obtenido al extraer el texto de la noticia, lo cual puede generar una considerable mejora en el comportamiento del sistema.

## ACKNOWLEDGMENT

This work has been funded by the Chilean Ministry of Education under the Project FONDECYT 1161584, and the Spanish Ministry of Economy and Competitiveness under the Project No. ENE2015-64914-C3-2-R.

## REFERENCES

- [1] G.E. Hinton, R.R. Salakhutdinov, "Reducing the dimensionality of data with neural networks," *Science*, vol. 313, no. 5786, pp. 504507, 2006.
- [2] A. Go, R. Bhayani, L. Huang, "Twitter sentiment classification using distant supervision," *CS224N Project Report*, Stanford, vol. 1, no. 12, 2009.
- [3] C. Schudt, I. Laptev, B. Caputo, Recognizing human actions: a local SVM approach. In *Proceedings of the 17th International Conference on Pattern Recognition, ICPR 2004*, vol. 3, pp. 32-36.
- [4] P. Cheeseman, M. J. Kelly, M. Self, J. Stutz, W. Taylor, D. Freeman. "Autoclass: A Bayesian classification system." In *Machine Learning Proceedings*, pp. 54-64. 1988.
- [5] A.L. Berger, J. Vincent, P. Della, A. Stephen. "A maximum entropy approach to natural language processing." *Computational linguistics*, vol. 22, no. 1, pp. 39-71, 1996.
- [6] G. K. Pitsilis, H. Ramampiaro, H. Langseth, "Detecting offensive language in tweets using deep learning," *arXiv preprint arXiv:1801.04433*, 2018.
- [7] G. Forman, "An extensive empirical study of feature selection metrics for text classification," *Journal of machine learning research*, vol. 3, no. Mar, pp. 12891305, 2003.
- [8] F. Chollet. *Deep learning with python*. Manning Publications, 2017.
- [9] S. Haykin. "Neural Networks: A comprehensive foundation," Prentice Hall, 2004.
- [10] I. Guyon, A. Elisseeff. "An introduction to feature extraction." In *Feature extraction*, pp. 1-25. Springer, Berlin, Heidelberg, 2006.
- [11] P. Vincent, H. Larochelle, I. Lajoie, Y. Bengio, P.A. Manzagol. Stacked denoising autoencoders: Learning useful representations in a deep network with a local denoising criterion. *Journal of Machine Learning Research*, pp.3371-3408, 2010.
- [12] T. Fawcett. "An Introduction to ROC Analysis", *Pattern Recognition Letters*. no. 27, vol 8, pp. 861874, 2006.

- [13] Kato, Ryoma, and Hiroyuki Goto. "Categorization of web news documents using word2vec and deep learning." Proceedings of the 2016 International Conference on Industrial Engineering and Operations Management Kuala Lumpur, Malaysia. 2016.
- [14] Bracewell, David B., et al. "Category classification and topic discovery of japanese and english news articles." *Electronic Notes in Theoretical Computer Science* 225 (2009): 51-65.
- [15] Asy'arie, Arni Darliani, and Adi Wahyu Pribadi. "Automatic news articles classification in indonesian language by using naive bayes classifier method." Proceedings of the 11th International Conference on Information Integration and Web-based Applications & Services. ACM, 2009.