



Clasificación automática de correos electrónicos

Juan M. Fernandez¹, Nicolás Cavasin¹, Agustín Rodríguez¹, Marcelo Errecalde²
 {jfernandez, ncavasin, arodriguez}@unlu.edu.ar, merreca@unsl.edu.ar



¹Departamento de Ciencias Básicas, Universidad Nacional de Luján

²LIDIC, Universidad Nacional de San Luis

Contexto

Este Proyecto de Investigación, aprobado por Disposición CDD-CB N° 086/2020 del Departamento de Ciencias Básicas de la Universidad Nacional de Luján, tiene como objetivo estudiar y analizar el conocimiento existente sobre técnicas de aprendizaje automático aplicadas a la clasificación automática de textos, particularmente de correos electrónicos, y generar modelos que aborden problemas concretos. Asimismo, se está trabajando en conjunto con el Laboratorio de Investigación y Desarrollo en Inteligencia Computacional (LIDIC) de la Universidad Nacional de San Luis a efectos de integrar los resultados encontrados en este proyecto con desarrollos de ese Laboratorio en torno a la clasificación de preguntas para sistemas conversacionales (chatbots).

LÍNEAS DE I+D

ESTRATEGIAS DE ETIQUETADO DE DOCUMENTOS

La estrategia tradicional para el etiquetado de documentos consiste en que esta tarea sea realizada por un humano, lo cual, con el crecimiento de la cantidad de datos, resulta poco conveniente. En el marco de este proyecto, se indagan las **estrategias de etiquetado semi-automático y automático existentes** y se exploran nuevas soluciones que resulten escalables y efectivas para la clasificación de correos electrónicos.

REPRESENTACIÓN DE DOCUMENTOS

Para construir un sistema de clasificación eficiente previamente se necesita tener un modelo de representación de documentos efectivo. En este proyecto, se estudian y aplican diferentes estrategias de representación de documentos usualmente utilizadas como **bag of words, topic modeling, embeddings y BERT**. El objetivo de esta línea de investigación es evaluar las diferentes técnicas aplicadas a la clasificación de correos electrónicos.

EVALUACIÓN DE ALGORITMOS DE CLASIFICACIÓN

En el marco del proyecto se están utilizando **Naive Bayes, la regresión logística, las máquinas de vector soporte, las redes neuronales recurrentes LSTM y XGBoost**.

El objetivo de esta línea es evaluar las técnicas y encontrar las características de los problemas que hacen más acorde la utilización de un algoritmo de aprendizaje automático por encima del resto.

MODELOS APLICABLES A PROBLEMAS CONCRETOS

La última línea de I+D consiste en el desarrollo de modelos de clasificación que resuelvan problemas concretos; la transferencia de estos conocimientos y herramientas a la sociedad para resolver problemas del mundo real resulta de alto interés.

Actualmente se está desarrollando un **clasificador automático para la clasificación de consultas académicas realizadas por los estudiantes** a la administración académica de la Universidad Nacional de Luján.

FORMULARIO DE CONTACTO PARA ESTUDIANTES
 Enviado: 03/24/2021-13:22:06

Boleto Universitario

Nombre y Apellido: Claudia [REDACTED]

Cambio Carrera

Legajo: 101000

Documento: [REDACTED]

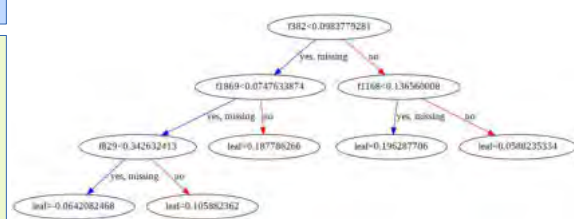
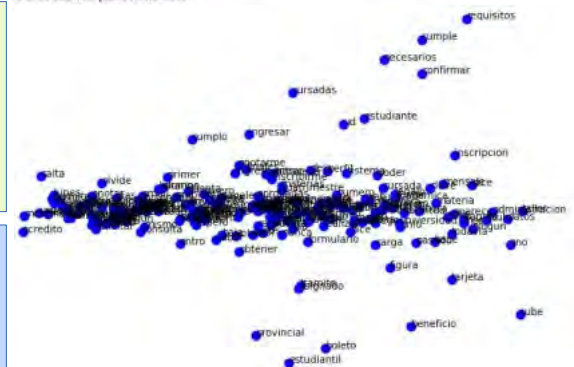
Carrera: LICENCIATURA EN CS. DE LA EDUCACION(4) ...

Teléfono: [REDACTED]

E-Mail: [REDACTED]

Ingreso UNLu

Mensaje / Consulta: Hola queria saber donde tengo que pedir la conform dos materias para recibirme.



Formación de Recursos Humanos

Se espera que este proyecto contribuya a consolidar un grupo de investigación en la temática y brindar un marco adecuado para la formación de recursos humanos en la Universidad Nacional de Luján a partir de la incorporación de saberes y competencias provenientes de la participación en actividades de investigación. Concretamente, se ha incluido como integrantes del proyecto a dos docentes auxiliares del Departamento de Ciencias Básicas y un estudiante de la Carrera de Licenciatura en Sistemas de Información. A su vez, el director trabaja con estos temas en su tesis de Maestría, esperando culminarla en el marco del proyecto.

Referencias

- Alghoul, A., Al Ajrami, S., Al Jarousha, G., Harb, G., and Abu-Naser, S. S. Email classification using artificial neural network. ACM (2018).
- Chen, Z. Tao, R., Wu, X., Wei, Z., and Luo, X. Active learning for spam email classification. In Proceedings of the 2019 2nd International Conference on Algorithms, Computing and Artificial Intelligence (2019).
- Harish, B. S., Guru, D. S., and Manjunath, S. Representation and classification of text documents: A brief review. IJCA, Special Issue on RTIPPR (2) (2010), 110-119.
- Henkel, M., Perjons, E., and Sneider. Examining the potential of language technologies in public organizations by means of a business. International Journal of Information Management (2017).
- Saha, S., Das Gupta, S., and Das, S. K. Spam mail detection using data mining: A comparative analysis. In Smart Intelligent Computing and Applications. Springer, 2019, pp. 571-580.
- Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., ... & Polosukhin, I. (2017). Attention is all you need. arXiv preprint arXiv:1706.03762.