

Métodos para la Selección y el Ajuste de Características en el Problema de la Detección de Spam

Carlos M. Lorenzetti^{†§} Rocío L. Cecchini^{†*} Ana G. Maguitman^{†§} Andrés A. Benczúr[‡]

[§]Laboratorio de Inv. y Des. en Inteligencia Artificial

^{*}Laboratorio de Inv. y Des. en Computación Científica

[†]Departamento de Cs. e Ing. de la Computación – Universidad Nacional del Sur

{cml, rlc, agm}@cs.uns.edu.ar

[‡]Data Mining and Web search Research Group – Informatics Laboratory

Computer and Automation Research Institute – Hungarian Academy of Sciences

benczur@ilab.sztaki.hu

1. INTRODUCCIÓN

El correo electrónico es quizás la aplicación que más tráfico genera en la Internet. Es utilizado por millones de personas para comunicarse alrededor del mundo y es una aplicación de misión crítica para muchos negocios. En la última década la avalancha de correo no deseado (Spam) ha sido el mayor problema para los usuarios del correo electrónico, ya que diariamente una cantidad arrolladora de spam entra en las bandejas de los usuarios. En 2004, se estimó que el 62 % de todos los correos que se generaron fueron spam [2]. El spam no solo es frustrante para muchos usuarios, sino que también compromete a la infraestructura tecnológica de las empresas, costando dinero a causa de la pérdida de productividad. En los últimos años, el spam ha evolucionado desde ser una molestia a ser un serio riesgo en la seguridad, llegando a ser el principal medio para el robo de información personal, así como también para la proliferación de software malicioso.

Muchas alternativas se han propuesto para solucionar el problema, desde protocolos de autenticación del remitente a, incluso, cobrarles dinero a los remitentes [10]. Otra alternativa prometedora es el uso de filtros basados en contenido capaces de discriminar automáticamente entre mensajes spam y mensajes legítimos. Los métodos de Aprendizaje Automatizado son atractivos para realizar esta tarea ya que son capaces de adaptarse a las características evolutivas del spam, contándose además con dis-

ponibilidad de datos para entrenar tales modelos. Sin embargo, uno de los aspectos más frustrantes del spam es que cambia continuamente para adaptarse a las nuevas técnicas que intentan detenerlo. Cada vez que se lo ataca de alguna manera, los generadores de spam encuentran una manera de eludir este ataque. Esta carrera ha llevado a una coevolución continua y a un aumento del nivel de sofisticación de ambas partes [10]. Otra diferencia con respecto a muchas tareas en la clasificación de texto consiste en que el costo de un error en la clasificación está fuertemente sesgado: etiquetar un correo legítimo como spam, usualmente llamado *falso positivo*, trae peores consecuencias que el caso inverso.

La detección de spam web puede verse como un problema de clasificación. Para detectar páginas web spam, construimos un clasificador para etiquetar una dada página como spam o como no spam. Centrándonos en el análisis del contenido semántico de los correos y de las páginas, se han estudiado varias técnicas de clasificación de texto basadas en métodos de Aprendizaje Automatizado y Reconocimiento de Patrones, debido principalmente a su mayor capacidad de generalización. Las técnicas de clasificación de texto (ver [25], para una revisión detallada) se aplican básicamente a documentos de texto representados en formato ASCII no estructurado, en formatos estructurados como HTML y también se aplican a mensajes de correo electrónico.

El proceso de clasificación comienza en la

fase de entrenamiento y necesita representar el texto plano que contienen los documentos, por esto el primer paso transforma los documentos a alguna representación interna. Luego se construye un *vocabulario* con todos los términos que se encontraron en los documentos, para luego pasar a una fase de extracción de características en donde, por lo general, se reduce la cardinalidad de las mismas. Esto se lleva a cabo mediante la eliminación de signos de puntuación y de palabras muy frecuentes, y por el stemming (reducción de las palabras a su palabra raíz o stem), con el propósito de descartar términos no discriminantes y de reducir el tamaño del vocabulario (y por lo tanto, de la complejidad computacional). Finalmente se representa el documento como un vector de longitud fija de características, en el cual cada componente (usualmente un número real) está asociado a un término del vocabulario. Los términos usualmente corresponden a palabras individuales, o a frases que se encuentran en los documentos de entrenamiento. Las técnicas de extracción de características más simples están basadas en un método de *bolsa de palabras*, en donde solo se tienen en cuenta la ocurrencia de los términos y se descarta la información de su posición dentro del documento. Las características más comunes son la ocurrencia de la palabra (valor booleano), el número de ocurrencias (valor entero), o su frecuencia relativa a la longitud del documento (valor real). Una característica llamada TFIDF tiene en cuenta el número de apariciones en el documento y en todos los documentos de entrenamiento.

Los clasificadores estadísticos pueden aplicarse a la representación vectorial de características. Las principales técnicas analizadas hasta hoy en este contexto para el filtrado de spam están basadas en el clasificador de texto Bayes Naïve [20] y en los llamados “filtros Bayesianos” [24, 11]. Dado su rendimiento en tareas de clasificación de texto, también se ha investigado el uso de clasificadores Máquina de Vectores de Soporte (SVM, Support Vector Machine [7, 28]).

2. LÍNEA DE INVESTIGACIÓN PROPUESTA

Como se dijo en la sección previa, la identificación de spam puede verse como un problema de clasificación. Por lo tanto proponemos un algoritmo que utiliza un clasificador como uno de sus componentes. Nuestra propuesta no incluye el desarrollo de un clasificador en sí mismo, sino que plantea un ajuste en los datos de entrada del conjunto de entrenamiento del clasificador con el objetivo de mejorar su rendimiento. El esquema general del sistema se muestra en la figura 1 y se describe a continuación.

2.1. Clustering

Dada la heterogeneidad que posee el spam, no puede asumirse que todo spam se asocia a un único tópico. Es por esto que proponemos, como primera etapa, la utilización de un algoritmo de clustering que dividirá a los documentos en subtópicos más pequeños esperando con esto una mejora en el rendimiento global del algoritmo. Una lista detallada de los algoritmos disponibles para este propósito puede encontrarse en [4].

2.2. Descriptores y Discriminadores

Una vez que los datos de entrada se encuentran agrupados en subtópicos más específicos, tomamos cada uno de ellos y calculamos los pesos de los términos en ellos como descriptores y discriminadores de estos subtópicos. En [19] proponemos estudiar el poder descriptivo y discriminante de un término en base a su distribución a través de los tópicos de las páginas recuperadas por un motor de búsqueda.

Para distinguir entre descriptores y discriminadores de tópicos argumentamos que *buenos descriptores de tópicos* pueden encontrarse buscando aquellos términos que aparecen con frecuencia en documentos relacionados con el tópico deseado. Por otro lado, *buenos discriminadores de tópicos* pueden hallarse buscando términos que aparecen solo en documentos relacionados con el tópico deseado. Ambos tipos

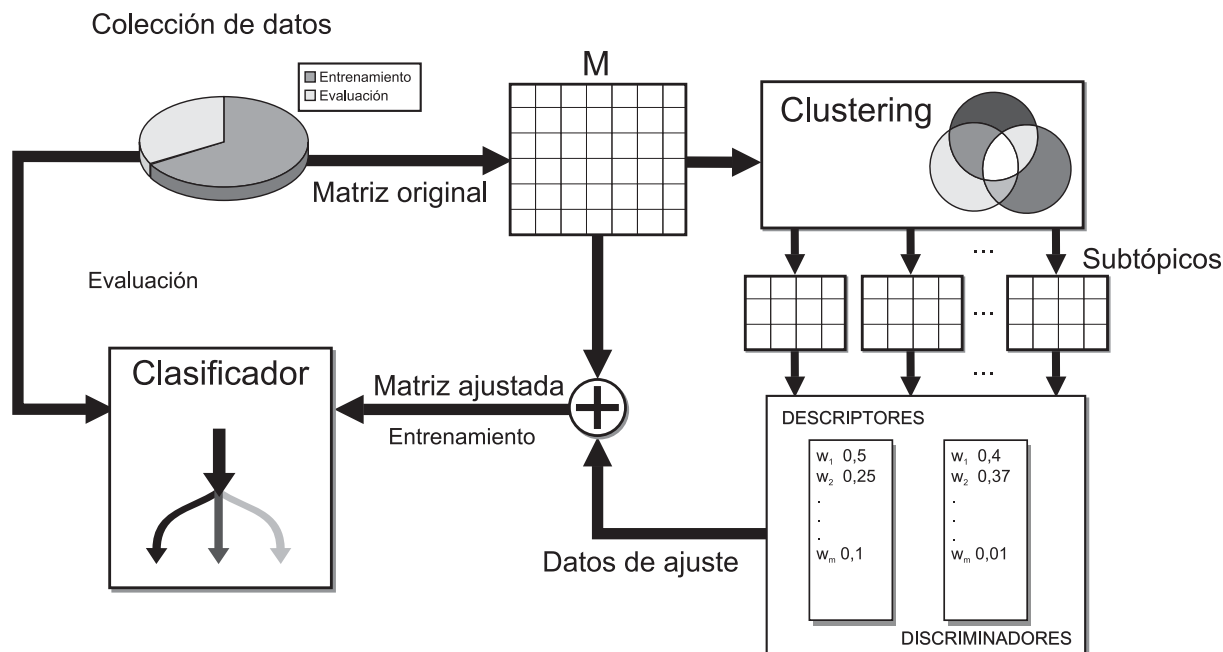


Figura 1: Diagrama esquemático de la propuesta

de términos son importantes a la hora de generar consultas. Utilizar términos descriptores del tópico mejora el problema de los resultados falso negativo porque aparecen frecuentemente en páginas relevantes. De la misma manera, los buenos discriminadores de tópicos ayudan a reducir el problema de los falsos positivos, ya que aparecen principalmente en páginas relevantes.

Esta etapa da como resultado listas de términos con información asociada a la importancia de los mismos como descriptores y discriminadores. Dicha información se utilizará para ajustar la matriz de datos de entrenamiento para reflejar de forma más fidedigna los pesos de los términos en los documentos.

2.3. Clasificador

Los clasificadores son implementados a partir de un conjunto de instancias o ejemplos previamente etiquetados, en donde cada ejemplo tiene un vector de atributos o características. En general, en los conjuntos de datos que utilizaremos en nuestras evaluaciones (descritos en la sección 3) las etiquetas fueron determinadas por personas.

La clasificación involucra la creación de un

modelo durante la etapa de entrenamiento que predecirá la etiqueta de cada instancia del conjunto de testeo usando los valores del vector de características. Para construir el clasificador, primero lo entrenamos sobre un número de ejemplos del conjunto etiquetado de entrenamiento y determinamos los parámetros de nuestro clasificador. Durante la etapa de testeo, el clasificador examina el vector de características de forma conjunta para determinar si una página pertenece a una dada categoría o no. La evaluación del clasificador se realiza en la etapa de testeo comparando, para cada instancia, la etiqueta calculada por el clasificador con la asignada a esa instancia.

Una lista detallada de los algoritmos disponibles para este propósito puede encontrarse en [16, 23]. Se prevé utilizar el entorno Weka [13] en esta etapa.

3. EVALUACIÓN

Para la evaluación de nuestra propuesta utilizaremos distintos conjuntos de datos disponibles, como por ejemplo el conjunto de datos UK-2007 del workshop internacional Air-Web [8], el conjunto de datos de la conferen-

cia internacional ECML PKDD [15], el conjunto de datos del track de Spam de la conferencia TREC [5] y el corpus de correos electrónicos SpamAssassin [26]. Para analizar la eficacia del método propuesto evaluaremos el rendimiento del clasificador. Para ello utilizaremos las métricas estándares de evaluación, como precisión, cobertura, F-score, Media Geométrica, área bajo la curva ROC, área bajo la curva Precisión-Cobertura y estadísticas de Kolmogorov-Smirnov.

4. CONCLUSIONES

La técnica propuesta en este trabajo ataca uno de los problemas más grandes a los que se deben enfrentar los usuarios de los sistemas de información actuales. Mejorar la representación de los documentos mediante el uso de vocabularios más representativos, así como el ajuste de los datos realizado a través de la detección de buenos descriptores y discriminadores ha mostrado ser efectivo en otras áreas de recuperación de información [18, 17]. Anticipamos que aplicar estos métodos será ventajoso para abordar diversos problemas de clasificación, en particular en el ámbito de la detección de spam.

Nuestro trabajo está relacionado con muchos estudios previos sobre clasificación de páginas web spam basados en características. Desde los comienzos de la World Wide Web ha existido una necesidad de calificar a las páginas de acuerdo a su relevancia con una dada consulta. Sin embargo se ha puesto un nuevo énfasis al problema dadas las grandes ganancias que genera la publicidad a través de Internet. La clasificación de spam web es uno de los desafíos más importantes de los motores de búsqueda [14], en particular debido a la degradación de la calidad de sus resultados. Un método prometedor para la identificación del spam web es la utilización de la información de los enlaces web que contienen las páginas [6, 1, 3, 12, 27]. Por otro lado recientemente se ha estudiado la clasificación de spam web basándose en el contenido de la página [22, 9]. La detección de spam en blogs se estudió en [21].

Todos estos ejemplos son solo los primeros pasos en el combate contra el spam: la naturaleza necesariamente adversaria de la tarea conlleva a un problema que evoluciona rápidamente, y esta característica (de tener que buscar técnicas que sean exitosas a la luz de la adaptación del enemigo) es algo nuevo en la comunidad de Aprendizaje Automatizado y trae consigo numerosos desafíos y oportunidades de investigación.

REFERENCIAS

- [1] E. Amitay, D. Carmel, A. Darlow, R. Lempel, and A. Soffer. The connectivity sonar: detecting site functionality by structural patterns. In *HYPertext '03: Proceedings of the fourteenth ACM conference on Hypertext and hypermedia*, pages 38–47, New York, NY, USA, 2003. ACM.
- [2] I. Androutsopoulos, J. Koutsias, K. Chandrinos, G. Paliouras, and C. D. Spyropoulos. An evaluation of Naive Bayesian anti-spam filtering. In G. Potamias, V. Moustakis, and M. van Someren, editors, *Proceedings of the Workshop on Machine Learning in the New Information Age, 11th European Conference on Machine Learning (ECML 2000)*, pages 9–17, Barcelona, Spain, 2000.
- [3] L. Becchetti, C. Castillo, D. Donato, S. Leonard, and R. Baeza-Yates. Using rank propagation and probabilistic counting for link-based spam detection. In *Proceedings of the Workshop on Web Mining and Web Usage Analysis (WebKDD)*, Pennsylvania, USA, August 2006. ACM Press.
- [4] P. Berkhin. A survey of clustering data mining techniques. In *Grouping Multidimensional Data*, pages 25–71. Springer Berlin Heidelberg, 2006.
- [5] G. V. Cormack. TREC 2007 Spam Track Overview. In *TREC*, 2007.
- [6] B. D. Davison. Recognizing nepotistic links on the Web. In *Artificial Intelligence for Web Search*, pages 23–28, Austin, Texas, USA, July 2000. AAAI Press.
- [7] H. Drucker, D. Wu, and V. Vapnik. Support vector machines for spam categorization. *IEEE Transactions on Neural Networks*,

- 10(5):1048–1054, 1999.
- [8] D. Fetterly and Z. Gyöngyi, editors. *AIR-Web '09: Proceedings of the 5th International Workshop on Adversarial Information Retrieval on the Web*, New York, NY, USA, 2009. ACM.
- [9] D. Fetterly, M. Manasse, and M. Najork. Spam, damn spam, and statistics: using statistical analysis to locate spam web pages. In *WebDB '04: Proceedings of the 7th International Workshop on the Web and Databases*, pages 1–6, New York, NY, USA, 2004. ACM.
- [10] J. Goodman, D. Heckerman, and R. Rounthwaite. Stopping spam. *Scientific American*, 292(4):42–49, April 2005.
- [11] P. Graham. A Plan for Spam, August 2002.
- [12] Z. Gyöngyi and H. Garcia-Molina. Link spam alliances. In K. Böhm, C. S. Jensen, L. M. Haas, M. L. Kersten, P.-Å. Larson, and B. C. Ooi, editors, *Proceedings of the 31st International Conference on Very Large Data Bases (VLDB)*, pages 517–528. ACM, 2005.
- [13] M. Hall, E. Frank, G. Holmes, B. Pfahringer, P. Reutemann, and I. H. Witten. The weka data mining software: an update. *ACM SIGKDD Explorations Newsletter*, 11(1):10–18, 2009.
- [14] M. R. Henzinger, R. Motwani, and C. Silverstein. Challenges in web search engines. In *IJCAI'03: Proceedings of the 18th International Joint Conference on Artificial Intelligence*, pages 1573–1579, San Francisco, CA, USA, 2003. Morgan Kaufmann Publishers Inc.
- [15] A. Hotho, D. Benz, R. Jäschke, and B. Krause, editors. *ECML PKDD Discovery Challenge 2008 (RSDC'08)*. Workshop at 18th European Conference on Machine Learning (ECML'08) / 11th European Conference on Principles and Practice of Knowledge Discovery in Databases (PKDD'08), 2008.
- [16] A. K. Jain, M. N. Murty, and P. J. Flynn. Data clustering: a review. *ACM Computing Surveys*, 31(3):264–323, 1999.
- [17] C. M. Lorenzetti and A. G. Maguitman. Tuning Topical Queries through Context Vocabulary Enrichment: A Corpus-based approach. In R. Meersman, Z. Tari, and P. Herrero, editors, *On the Move to Meaningful Internet Systems: OTM 2008 Workshops*, volume 5333 of *LNCS*, pages 646–655. Springer, 2008.
- [18] C. M. Lorenzetti and A. G. Maguitman. A semi-supervised incremental algorithm to automatically formulate topical queries. *Information Sciences*, 179(12):1881–1892, 2009. Including Special Issue on Web Search.
- [19] A. Maguitman, D. Leake, T. Reichherzer, and F. Menczer. Dynamic Extraction of Topic Descriptors and Discriminators: Towards automatic context-based topic search. In *Proceedings of the Thirteenth Conference on Information and Knowledge Management (CIKM)*, pages 463–472, Washington, DC, November 2004. ACM Press.
- [20] A. McCallum and K. Nigam. A Comparison of Event Models for Naive Bayes Text Classification. In *Learning for Text Categorization: Papers from the 1998 AAAI Workshop*, volume 752, pages 41–48, 1998.
- [21] G. Mishne, D. Carmel, and R. Lempel. Blocking Blog Spam with Language Model Disagreement. In *Proceedings of the First International Workshop on Adversarial Information Retrieval on the Web (AIRWeb)*, Chiba, Japan, May 2005.
- [22] A. Ntoulas, M. Najork, M. Manasse, and D. Fetterly. Detecting spam web pages through content analysis. In *WWW '06: Proceedings of the 15th International Conference on World Wide Web*, pages 83–92, New York, NY, USA, 2006. ACM.
- [23] X. Qi and B. D. Davison. Web page classification: Features and algorithms. *ACM Computing Surveys (CSUR)*, 41(2):1–31, 2009.
- [24] M. Sahami, S. Dumais, D. Heckerman, and E. Horvitz. A bayesian approach to filtering junk E-mail. In *Learning for Text Categorization: Papers from the 1998 Workshop*, Madison, Wisconsin, 1998. AAAI Technical Report WS-98-05.
- [25] F. Sebastiani. Machine learning in automated text categorization. *ACM Computing Surveys*, 34(1):1–47, 2002.
- [26] SpamAssassin. SpamAssassin public corpus, <http://spamassassin.apache.org/publiccorpus/>.
- [27] B. Wu and B. D. Davison. Identifying link farm spam pages. In *WWW '05: Special interest tracks and posters of the 14th International Conference on World Wide Web*, pages 820–829, New York, NY, USA, 2005. ACM Press.
- [28] L. Zhang, J. Zhu, and T. Yao. An evaluation of statistical spam filtering techniques. *ACM Transactions on Asian Language Information Processing (TALIP)*, 3(4):243–269, 2004.