



UNIVERSIDAD DE CHILE
FACULTAD DE CIENCIAS FÍSICAS Y MATEMÁTICAS
DEPARTAMENTO DE INGENIERÍA INDUSTRIAL

**MODELAMIENTO Y SEGUIMIENTO DE TÓPICOS PARA DETECCIÓN DE
MODUS OPERANDI EN ROBO DE VEHÍCULOS**

TESIS PARA OPTAR AL GRADO DE MAGÍSTER EN GESTIÓN DE OPERACIONES

MEMORIA PARA OPTAR AL TÍTULO DE INGENIERO CIVIL INDUSTRIAL

DIEGO GARRIDO

PROFESOR GUÍA:
RICHARD WEBER

MIEMBROS DE LA COMISIÓN:
PROFESOR 2
PROFESOR 3

Este trabajo ha sido parcialmente financiado por:
NOMBRE INSTITUCIÓN

SANTIAGO, CHILE
2020

*Una frase de dedicatoria,
pueden ser dos líneas.*

Saludos

Agradecimientos

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.

Tabla de Contenidos

| | |
|---|-----------|
| 1. Introducción | 1 |
| 2. Definición del problema y Objetivos | 2 |
| 2.1. Problema | 2 |
| 2.2. Objetivo, Resultados esperados y Alcances | 3 |
| 2.3. Metodología de trabajo | 3 |
| 3. Revisión bibliográfica | 4 |
| 4. Marco teórico | 6 |
| 4.1. Modelos de tópicos | 6 |
| 4.1.1. Latent Dirichlet Allocation | 6 |
| 4.1.1.1. Distribución Dirichlet | 6 |
| 4.1.1.2. LDA | 7 |
| 4.1.2. Hierarchical Dirichlet Process | 8 |
| 4.1.2.1. Proceso Dirichlet | 8 |
| 4.1.2.1.1. Stick breaking construction | 9 |
| 4.1.2.2. HDP | 9 |
| 4.1.2.3. LDA versus HDP | 11 |
| 4.2. Modelamiento de la evolución de los tópicos en el tiempo | 11 |
| 4.2.1. Gráfo de similitud temporal | 11 |
| 4.2.2. Medidas de similitud | 13 |
| 5. Experimento | 15 |
| 5.1. Datos | 15 |
| 5.2. Procesamiento | 15 |
| 5.3. Análisis cuantitativo de resultados | 18 |
| 5.4. Análisis cualitativo de resultados | 18 |
| 6. Conclusiones | 19 |
| Bibliografía | 20 |

Índice de Tablas

| | | |
|------|--|----|
| 5.1. | Estadísticas del corpus bajo distintos niveles de procesamientos, raw : sin procesamiento, ch : eliminación de símbolos de puntuación, correos electrónicos y tokens con números, ch+s+l+f : además incluye eliminación de stopwords (s), lematización (l) y eliminación de tokens con baja ocurrencia (f). | 18 |
| 5.2. | Evolución del vocabulario en el tiempo, old_vocabulary : corresponde al vocabulario del período $t - 1$, new_vocabulary : corresponde al vocabulario del período t , %old_tokens : porcentaje de tokens del período $t - 1$ que ya no están en el período t y %new_tokens : porcentaje de tokens del período t que no están en el período $t - 1$ | 18 |

Índice de Ilustraciones

| | | |
|------|--|----|
| 2.1. | (a) Cantidad de robos de vehículos y robos de accesorios de vehículos anuales en Chile (2004-2014). Fuente: Informe anual Carabineros, 2004-2014, INE. (b) Tasa de robos con violencia del total de robo de autos de lujo 2011-2016. . . . | 2 |
| 4.1. | Efecto de los parámetros de una distribución Dirichlet en el muestreo para $K = 3$ | 7 |
| 4.2. | Representación gráfica de LDA: círculos denotan variables aleatorias, círculos abiertos denotan parámetros, círculos sombreados denotan variables observadas y los platos indican replicación. | 8 |
| 4.3. | Ilustración de <i>stick breaking construction</i> . (a) Tenemos una barra de largo 1, el cual se rompe en un punto aleatorio β_1 , el largo de la pieza que conservamos es llamada π_1 , luego recursivamente rompemos la barra restante, así generando π_2, π_3, \dots (b) Muestras de π_k para $\alpha = 2$ y $\alpha = 5$ | 9 |
| 4.4. | Representación gráfica de HDP: círculos denotan variables aleatorias, círculos abiertos denotan parámetros, círculos sombreados denotan variables observadas y los platos indican replicación. | 10 |
| 4.5. | Representación gráfica de la construcción stick-breaking de HDP: círculos denotan variables aleatorias, círculos abiertos denotan parámetros, círculos sombreados denotan variables observadas y los platos indican replicación. | 11 |
| 4.6. | Ilustración conceptual del grafo de similitud que modela la dinámica de los tópicos en el tiempo. Un nodo corresponde a un tópico en una época específica; el ancho de los arcos es proporcional a la similitud entre los tópicos, arcos ausentes fueron eliminados por presentar una similitud menor a un umbral. | 12 |
| 4.7. | Espacio vectorial de los <i>word embeddings</i> de las palabras de dos tópicos con un vocabulario de tamaño 4. | 14 |
| 5.1. | Frecuencia acumulada de los tokens únicos aplicando hasta el primer paso de procesamiento. El eje horizontal es el acumulado de tokens únicos en orden decreciente de ocurrencia. Los puntos corresponden a los cuantiles 60 %, 80 %, 90 %, 95 % y 99 %. | 16 |
| 5.2. | Frecuencia acumulada de los tokens únicos aplicando hasta el cuarto nivel de procesamiento. El eje horizontal es el acumulado de tokens únicos en orden decreciente de ocurrencia. Los puntos corresponden a los cuantiles 60 %, 80 %, 90 %, 95 % y 99 %. | 17 |

RESUMEN DE LA MEMORIA PARA OPTAR
AL TÍTULO DE MAGÍSTER EN CIENCIAS
DE LA INGENIERÍA
POR: **DIEGO GARRIDO**
FECHA: 2020
PROF. GUÍA: RICHARD WEBER

MODELAMIENTO Y SEGUIMIENTO DE TÓPICOS PARA DETECCIÓN DE MODUS OPERANDI EN ROBO DE VEHÍCULOS

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.

Capítulo 1

Introducción

Capítulo 2

Definición del problema y Objetivos

2.1. Problema

El robo de vehículos o accesorios de vehículos es un problema que afecta a toda la sociedad en Chile y en el mundo. Este problema se ha vuelto más relevante el último tiempo debido al crecimiento en el robo de vehículo motorizado y de los robos con violencia (ver Figura 2.1). Este fenómeno trae consigo un montón de costos para la sociedad, como incremento en la percepción de la seguridad, aumentos en la prima de los seguros de los asegurados, aumento en los costos de las aseguradoras ¹ y el incremento de otros tipos de delitos ²

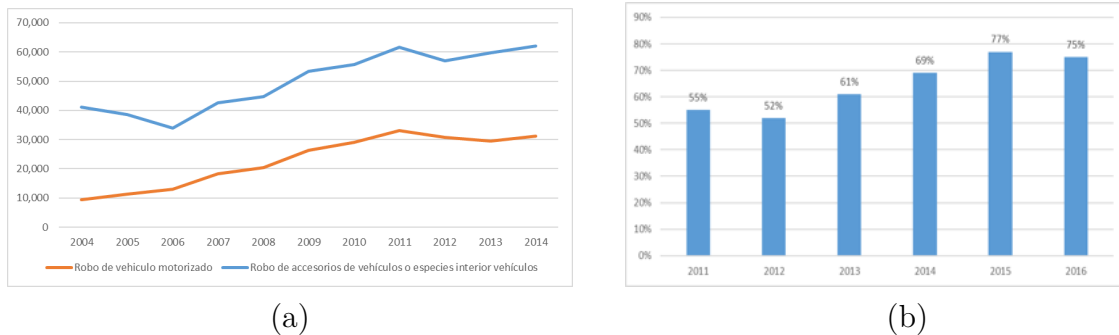


Figura 2.1: (a) Cantidad de robos de vehículos y robos de accesorios de vehículos anuales en Chile (2004-2014). Fuente: Informe anual Carabineros, 2004-2014, INE. (b) Tasa de robos con violencia del total de robo de autos de lujo 2011-2016.

Bajo este contexto la Universidad de Chile junto a la Pontificia Universidad Católica de Chile se adjudicó el 2017 un proyecto Fondef para desarrollar un proyecto que lleva por nombre “Observatorio Digital de Delincuencia en Chile: Un sistema inteligente de apoyo a la industria automotriz chilena, en el robo de vehículos y accesorios” cuyo director es Richard

¹ Considerando que el costo promedio incurrido en un auto asegurado robado y no recuperado es de \$ 5.000.000 de pesos, la pérdida total considerando solo los vehículos no recuperados para el año 2015 es de unos \$15.720 millones de pesos.

² El destino de los vehículos robados es variado, se usan los autos para perpetrar otros delitos y huir, venderlos por piezas en talleres clandestinos o blanquear sus documentos para pasarlos por la frontera y venderlos o cambiarlos por droga en el extranjero”.

Weber Haas y la institución beneficiaria es la Asociación de Aseguradores de Chile (AACH).

Para este problema se cuenta con las fuentes de datos de la AACH, lo que corresponde a relatos de las víctimas del robo de sus vehículos desde el 2011 hasta el 2016, lo cual corresponde a 49.015 relatos. Cabe destacar que se estima que un tercio del parque automotriz se encuentra asegurado, por lo que se trabaja con una muestra del parque automotriz.

2.2. Objetivo, Resultados esperados y Alcances

El objetivo del trabajo de tesis es caracterizar los *modus operandi* de los delincuentes a partir de los relatos de víctimas de robo de vehículo entregados por la AACH.

El resultado esperado es descubrir los *modus operandi* ocultos en los relatos de las víctimas y caracterizarlos a partir de las palabras, como también ver su evolución a través del tiempo, siendo capaz de detectar cuando nacen y mueren, y como cambian en el tiempo.

El presente trabajo tiene un propósito académico, puesto que no cuenta con un cliente particular y tiene por objetivo estudiar técnicas de *clustering* dinámico para detectar patrones en el contexto de robo de vehículos, sin embargo, potenciales beneficiarios del trabajo podrían ser las aseguradoras, los asegurados, carabineros de Chile y la sociedad.

2.3. Metodología de trabajo

El trabajo se realizó bajo la metodología CRISP-DM (Cross Industry Standard Process for Data mining)([Chapman et al., 2000](#)), la cual posee las siguientes seis etapas:

1. Comprensión del negocio:
2. Comprensión de los datos
3. Preparación de los datos
4. Modelamiento
5. Evaluación
6. Implementación

Dentro de los alcances del trabajo no está contemplado la puesta en producción de una solución basada en *machine learning*, el alcance es hasta la evaluación e interpretación de los resultados arrojados por el modelo.

Capítulo 3

Revisión bibliográfica

El problema planteado consiste en un problema de *clustering*, puesto que no se cuenta con una etiqueta del *modus operandi* al que corresponde cada relato, siendo el propósito del trabajo descubrirla. Dentro de los métodos de *clustering* que involucran texto el modelamiento de tópicos es el enfoque más prometedor. El modelamiento de tópicos es una herramienta estadística que busca encontrar los temas (tópicos) presentes en un conjunto de documentos (corpus), permitiendo organizar, buscar, indexar, explorar y comprender grandes colecciones de documentos. Los modelos de tópicos asumen que los documentos pueden ser representados por una mezcla de tópicos, donde los tópicos son distribuciones sobre las palabras, los tópicos son latentes y la inferencia tiene por objetivo descubrir la mezcla de tópicos que originó cada documento y la distribución sobre las palabras de cada tópico. En modelamiento de tópicos las personas son las que le dan una interpretación a los tópicos inferidos a partir de las palabras más relevantes y en base a esa información los etiquetan, por ejemplo, para un tópico, dentro de sus cinco palabras más probables se halla la siguiente secuencia: “llaves”, “domicilio”, “individuos”, “casa” y “portón”, una etiqueta válida para este tópico podría ser “portonazo”.

Algunas de las técnicas de modelamiento de tópicos están basadas en factorización matricial como LSI (Latent Semantic Indexing) (Dumais, 2004) o NMF (Non-negative Matrix Factorization) (Xu et al., 2003), pero en este trabajo se utilizarán técnicas basadas en modelos probabilísticos generativos, como LDA (Latent Dirichlet Allocation) (Blei et al., 2003) o HDP (Hierarchical Dirichlet Process) (Teh et al., 2005). Ambos enfoques tienen sus pros y contras, en este trabajo se prefiere el enfoque probabilísticos ya que es capaz de expresar incertidumbre en la asignación de un tópico a un documento y en la asignación de palabras a los tópicos, además, este enfoque suele aprender tópicos más descriptivos (Stevens et al., 2012).

El presente trabajo busca capturar el dinamismo que puede presentar el fenómeno del robo de vehículos. El aspecto dinámico del problema considera:

1. Nacimiento, muerte, fusión y división de tópicos: En el contexto de robos es natural que en el tiempo aparezcan nuevos *modus operandi* como también que desaparezcan aquellos que ya no parecen tan atractivos.
2. Dinamismo en la mezcla de tópicos: esto permite capturar la popularidad de los tópicos en el tiempo.

3. Evolución de los tópicos: la evolución de los tópicos se refleja en el cambio en la distribución sobre las palabras, esto permite detectar cambios en cómo se comete un mismo tipo de delito, por ejemplo, el “portonazo” en un determinado momento se comete en grupos de 2-3 personas con arma blanca, luego evoluciona de arma blanca a arma de fuego y lo perpetran jóvenes menores de edad.

Dentro de los modelos de tópicos probabilísticos existen modelos estáticos y dinámicos:

1. Dentro de los modelos estáticos destaca LDA y HDP. La diferencia principal en estos dos modelos es que el primero necesita de antemano fijar el número de tópicos a descubrir y el segundo lo infiere a partir del corpus.
2. Dentro de los modelos dinámicos están aquellos que mantienen el número de tópicos fijos durante el tiempo y los que no:
 - a) En el primer grupo destaca Dynamic Topic Modelling (DTM)([Blei and Lafferty, 2006](#)) junto Topic Over Time (TOC)([Wang and McCallum, 2006](#)), la gran desventaja de estos modelos es que si aparece un nuevo tópico este quedará clasificado dentro de un tópico que existía desde el comienzo, por lo que solo es capaz de capturar el punto 2 y 3.
 - b) Dentro de los modelos que no mantienen el número fijo de tópicos en el tiempo existen de dos tipos, aquellos que modelan todo el problema bajo un modelo monolítico, en este grupo destaca Dynamic Hierarchical Dirichlet Process (DHDP)([Ahmed and Xing, 2012](#)), el cual modela el problema de dinamismo de una forma elegante pero a la vez acompañada de una inferencia bastante complicada, de los dinamis-mos mencionados captura los puntos 2, 3 y el 1 parcialmente, ya que no es capaz de capturar fusión y división de tópicos, uno de los principales contras de esta solución es que no se trata de una tecnología madura, puesto que no cuenta con una implementación disponible a diferencia de los otros modelos mencionados, los cuales se encuentran disponibles en múltiples lenguajes de programación y cuentan con una amplia adopción de la comunidad científica. El segundo tipo de modelos que no mantienen fijo el número de tópicos utilizan modelos de tópicos estáticos de forma iterativa, lo que hacen es dividir el corpus en épocas, luego entrenan de forma independiente un modelo de tópico para cada época y luego unen los resultados obtenidos, un ejemplo utilizando LDA en ([Wilson and Robinson, 2011](#)) y con HDP en ([Beykikhoshk et al., 2018](#)).

En este trabajo se utilizarán técnicas de modelado dinámico de tópicos como las presentadas en ([Wilson and Robinson, 2011](#); [Beykikhoshk et al., 2018](#)), debido a que son capaces de modelar los tres puntos mencionados sobre dinamismo y se basan en tecnologías maduras.

Capítulo 4

Marco teórico

4.1. Modelos de tópicos

4.1.1. Latent Dirichlet Allocation

4.1.1.1. Distribución Dirichlet

La distribución Dirichlet es una generalización multivariada de la distribución beta, la cual tiene soporte sobre un símplex, definido por:

$$S_K = \{x : 0 \leq x_k \leq 1, \sum_{k=1}^K x_k = 1\} \quad (4.1)$$

Luego, su función de densidad de probabilidad (pdf):

$$Dir(x|\alpha) = \frac{1}{B(\alpha)} \prod_{k=1}^K x_k^{\alpha_k-1} \quad (4.2)$$

La distribución Dirichlet es una distribución útil para generar distribuciones de probabilidades categóricas. En general se asume simetría en los parámetros de la distribución, es decir, $\alpha_k = \frac{\alpha}{K}$. En la Figura 4.1 se observa el efecto de los parámetros de una Dirichlet en la muestra generada, para $\alpha_k = 1$ se tiene una distribución uniforme en el dominio S_K , α_k controla la *sparsity*, mientras más se acerca a 0 los vectores generados tienen más componentes nulos y se concentra la masa en unas pocas coordenadas, mientras más grande α_k la masa tiende a distribuirse uniformemente en las coordenadas de los vectores generados, por último, cuando α no es simétrico la masa se concentra en aquellas coordenadas cuyo α_k es más grande.

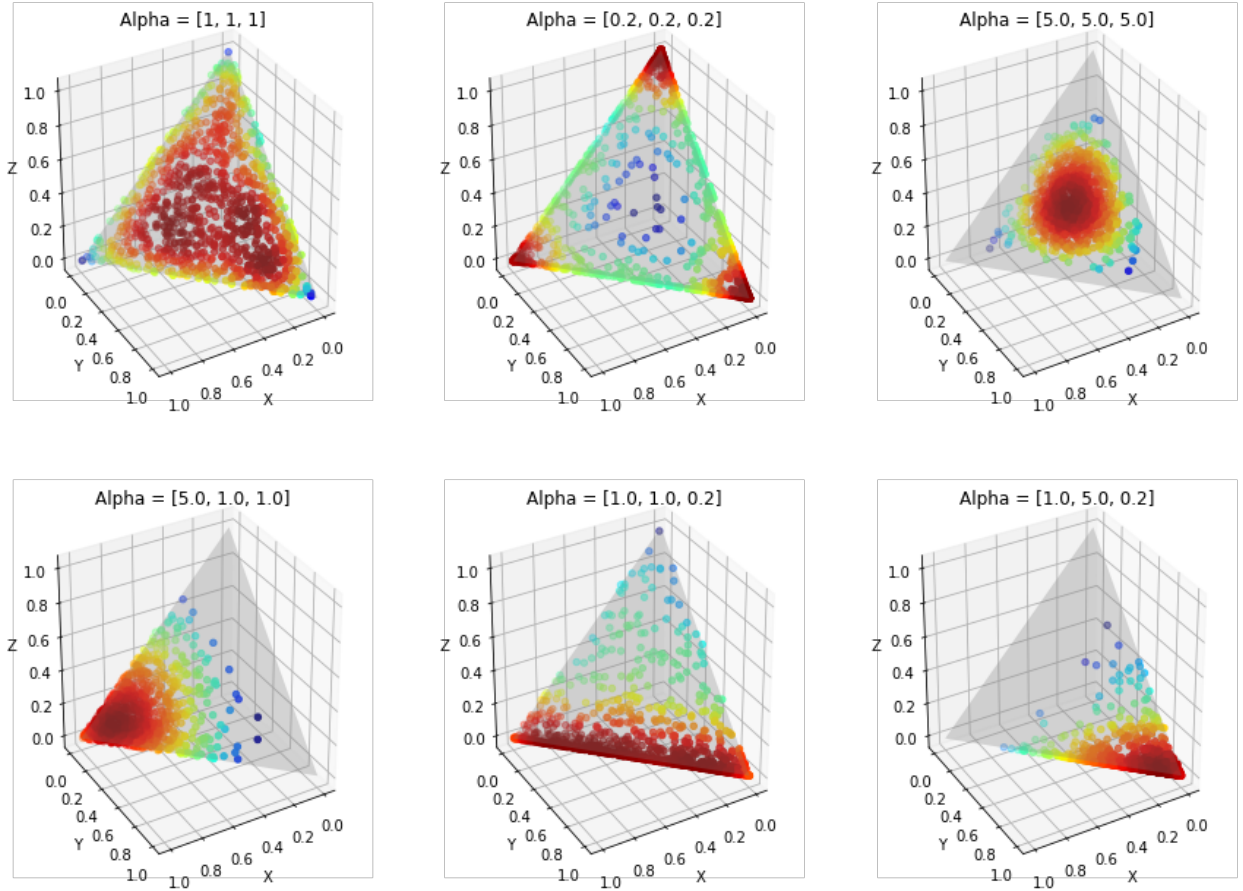


Figura 4.1: Efecto de los parámetros de una distribución Dirichlet en el muestreo para $K = 3$.

4.1.1.2. LDA

A continuación se describe el proceso generativo de Latent Dirichlet Allocation (LDA). Sean K tópicos, $\phi_{1:K}$ distribuciones de probabilidad sobre un vocabulario fijo, dibujadas por una $Dir(\frac{\eta}{|V|}1_{|V|})$. Para cada documento d del corpus D se asume que es dibujado por el siguiente proceso generativo (ver representación gráfica del modelo en la Figura 4.2):

1. Dibujar una mezcla de tópicos $\pi_d \sim Dir(\frac{\alpha}{K}1_K)$
2. Para cada palabra:
 - a) Escoger un tópico $z_{d,n} \sim Mult(\pi_d)$
 - b) Escoger una palabra $w_{d,n} \sim Mult(\phi_{z_{d,n}})$

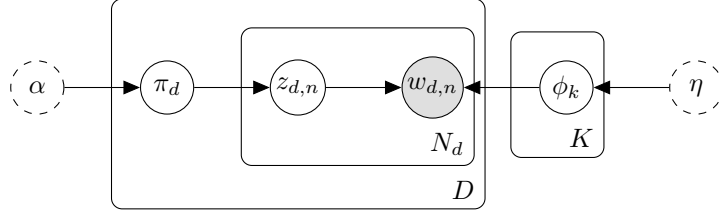


Figura 4.2: Representación gráfica de LDA: círculos denotan variables aleatorias, círculos abiertos denotan parámetros, círculos sombreados denotan variables observadas y los platos indican replicación.

La probabilidad conjunta del modelo:

$$p(\phi, \pi, z, w | \alpha, \eta) = \prod_{k=1}^K p(\phi_k | \eta) \prod_{d=1}^D p(\pi_d | \alpha) \prod_{n=1}^{N_d} p(z_{n,d} | \pi_d) p(w_{d,n} | \phi_{1:K}, z_{d,n}) \quad (4.3)$$

La distribución a posterior:

$$p(\phi, \pi, z | w, \alpha, \eta) = \frac{p(\phi, \pi, z, w | \alpha, \eta)}{p(w | \alpha, \eta)} \quad (4.4)$$

La distribución posterior es computacionalmente intratable para inferencia exacta, debido a que para normalizar la distribución debemos marginalizar sobre todas las variables ocultas y escribir la constante de normalización en términos de los parámetros del modelo. Para poder computar la posterior es necesario utilizar algoritmos de inferencia aproximada, donde el enfoque habitual es Markov Chain Monte Carlo (MCMC), en (Griffiths and Steyvers, 2004) se propone un algoritmo basado en Gibbs Sampling. Luego para estimación puntual de las cantidades relevantes se utiliza Monte Carlo para computar la esperanza.

4.1.2. Hierarchical Dirichlet Process

4.1.2.1. Proceso Dirichlet

En los modelos probabilísticos de *clustering* existen aquellos basados en mezcla finita de componentes que utilizan como *prior* la distribución Dirichlet y aquellos basados en mezcla infinita de componentes que utilizan como *prior* el proceso Dirichlet, un *prior* no paramétrico que no impone una cota en el número de *clusters* a encontrar (K).

Una representación equivalente en LDA sería generar cada palabra de un documento d a partir de una multinomial sobre un tópico dibujado por una distribución G_d , formalmente, $w_{d,n} \sim Mult(\phi_{d,n})$, donde $\phi_{d,n} \sim G_d$ con $\phi_{d,n} \in \{\phi_k\}_{k=1}^K$, y $G_d(\phi) = \sum_{k=1}^K \pi_{d,k} \delta_{\phi_k}(\phi)$, donde $\delta_{\phi_k}(\phi) = \begin{cases} 1 & \text{si } \phi_k = \phi \\ 0 & \text{si no} \end{cases}$.

Un proceso Dirichlet (DP) es una distribución sobre medidas de probabilidad $G : \Theta \rightarrow \mathbf{R}^+$, donde $G(\theta) \geq 0$ y $\int_{\Theta} G(\theta) d\theta = 1$. Un DP se define implícitamente por cumplir que para cualquier partición finita (T_1, \dots, T_K) de Θ , una medida base H y un parámetro de concentración α se tiene que $(G(T_1), \dots, G(T_K)) \sim Dir(\alpha H(T_1), \dots, \alpha H(T_K))$.

4.1.2.1.1. Stick breaking construction

En esta sección describiremos una definición constructiva para el DP conocida como *stick breaking construction*. Sea $\pi = \{\pi_k\}_{k=1}^{\infty}$ una secuencia de mezcla de pesos derivadas a partir del siguiente proceso:

$$\beta_k \sim \text{Beta}(1, \alpha) \quad (4.5)$$

$$\pi_k = \beta_k \prod_{l=1}^{k-1} (1 - \beta_l) = \beta_k (1 - \sum_{l=1}^{k-1} \pi_l) \quad (4.6)$$

Esto se suele denotar como $\pi \sim GEM(\alpha)$, donde GEM representa Griffiths, Engen y McCloskey. Algunos ejemplos de este proceso son mostrados en la Figura 4.3.

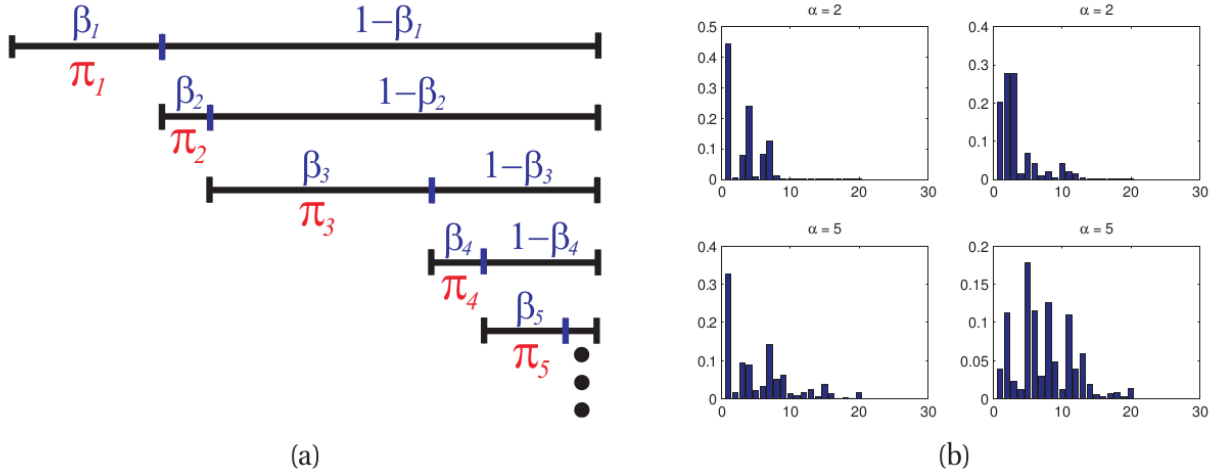


Figura 4.3: Ilustración de *stick breaking construction*. (a) Tenemos una barra de largo 1, el cual se rompe en un punto aleatorio β_1 , el largo de la pieza que conservamos es llamada π_1 , luego recursivamente rompemos la barra restante, así generando π_2, π_3, \dots (b) Muestras de π_k para $\alpha = 2$ y $\alpha = 5$.

Se puede demostrar que este proceso terminará con probabilidad 1 (convergencia casi segura), a pesar que el número de elementos que este genera incrementa con α . Además, el tamaño del componente π_k decrece en promedio. Ahora definamos

$$G(\phi) = \sum_{k=1}^{\infty} \pi_k \delta_{\phi_k}(\phi) \quad (4.7)$$

donde $\pi \sim GEM(\alpha)$ y $\phi_k \sim H$, se puede demostrar que $G \sim DP(\alpha, H)$. Como consecuencia de esta construcción, las muestras de un DP son discretas con probabilidad uno. En otras palabras, al ir muestreando se tendrán mas repeticiones de valores generados previamente, por lo que la mayoría de los datos vendrán de los ϕ_k con π_k mayor.

4.1.2.2. HDP

Hierarchical Dirichlet Process (HDP) es una colección de DP que comparten una distribución base G_0 , la cual además es dibujada a partir de un DP (ver representación gráfica del modelo en la Figura 4.4). Matemáticamente, a nivel corpus se tiene que la distribución base

$H \sim \text{Dir}(\frac{1}{|V|} \mathbf{1}_{|V|})$ y $G_0 \sim \text{DP}(\gamma, H)$, luego, para cada documento d del corpus D se asume que es dibujado por el siguiente proceso generativo:

1. Dibujar un DP $G_d \sim \text{DP}(\alpha_0, G_0)$
2. Para cada palabra:
 - a) Dibujar un t3pico $\phi_{d,n} \sim G_d$
 - b) Escoger una palabra $w_{d,n} \sim \text{Mult}(\phi_{d,n})$

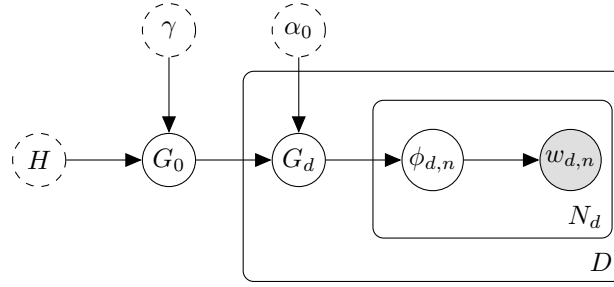


Figura 4.4: Representación gráfica de HDP: círculos denotan variables aleatorias, círculos abiertos denotan parámetros, círculos sombreados denotan variables observadas y los platos indican replicación.

La discretitud a nivel corpus de G_0 asegura que todos los documentos comparten el mismo conjunto de t3picos (*mixture components*). A nivel documento G_d hereda los t3picos de G_0 , pero los pesos de cada t3pico (*mixture proportions*) es específica del documento.

Aplicando *stick breaking construction* se tiene que para el DP dibujado a nivel corpus la siguiente representación:

$$\begin{aligned}
\beta'_k &\sim \text{Beta}(1, \gamma) \\
\beta_k &= \beta'_k \prod_{l=1}^{k-1} (1 - \beta'_l) \\
\phi_k &\sim H \\
G_0(\phi) &= \sum_{k=1}^{\infty} \beta_k \delta_{\phi_k}(\phi)
\end{aligned} \tag{4.8}$$

Así, G_0 es discreto y tiene soporte en los átomos $\phi = \{\phi\}_{k=1}^{\infty}$ con pesos $\beta = \{\beta_k\}_{k=1}^{\infty}$, siendo la distribución de β escrita como $\beta \sim \text{GEM}(\gamma)$. La construcción a nivel documento de G_d es:

$$\begin{aligned}
\pi'_{d,k} &\sim \text{Beta}\left(\alpha_0 \beta_k, \alpha_0 \left(1 - \sum_{l=1}^k \beta_l\right)\right) \\
\pi_{d,k} &= \pi'_{d,k} \prod_{l=1}^{k-1} (1 - \pi'_{d,l}) \\
G_d(\phi) &= \sum_{k=1}^{\infty} \pi_{d,k} \delta_{\phi_k}(\phi)
\end{aligned} \tag{4.9}$$

Donde $\phi = \{\phi_k\}_{k=1}^{\infty}$ son los mismos átomos de G_0 . En la Figura 4.5 se muestra la representación gráfica de esta construcción.

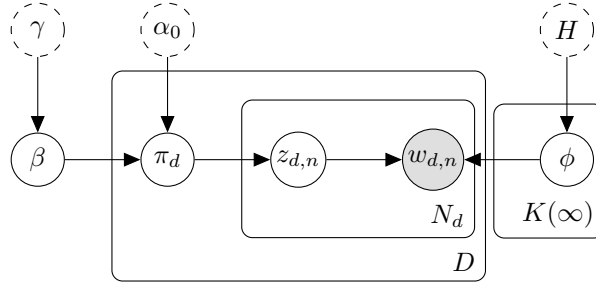


Figura 4.5: Representación gráfica de la construcción stick-breaking de HDP: círculos denotan variables aleatorias, círculos abiertos denotan parámetros, círculos sombreados denotan variables observadas y los platos indican replicación.

Al igual que LDA la distribución posterior es intratable, por lo que en Teh et al. (2005) se marginaliza G_0 y G'_d s afuera, obteniéndose así un nuevo proceso generativo denominado *Chinese restaurant franchise process* (CRF), esta representación permite construir algoritmos eficientes basados en Gibbs Sampling.

4.1.2.3. LDA versus HDP

HDP es un modelo no paramétrico similar en estructura a LDA, la principal desventaja de LDA frente a HDP es que LDA requiere escoger el número de tópicos K por adelantado, por otro lado, HDP el número de tópicos no está acotado y es inferido a partir de los datos. En un enfoque tradicional, se requiere de entrenar múltiples veces LDA para diferentes valores de K y se escoge el que tiene mejor la configuración con mejor desempeño en un conjunto de validación, por lo que LDA termina siendo computacionalmente más costoso que HDP, además este enfoque se vuelve impracticable cuando el conjunto de datos es largo. En el aspecto de cualitativo ambos modelos entregan tópicos igual de consistentes, en métricas de desempeño como *perplexity* HDP suele tener mejor desempeño (Teh et al. (2005)).

4.2. Modelamiento de la evolución de los tópicos en el tiempo

Nuestro objetivo es modelar la evolución en el tiempo de los tópicos, para esto el corpus es dividido en T épocas, en cada época se entrena un modelo de tópicos estático y se obteniéndose T conjuntos de tópicos $\phi = \{\phi_1, \dots, \phi_T\}$, donde $\phi_t = \{\phi_{t,1}, \dots, \phi_{t,K_t}\}$ es el conjunto de tópicos que describen la época t , y K_t el número de tópicos inferido en esa época.

4.2.1. Gráfo de similitud temporal

Para relacionar los tópicos de una época necesitamos una medida de similitud $\rho \in [0, 1]$, con esta mérida de similitud se puede construir un gráfo, donde los nodos son los tópicos de una época y los arcos relacionan tópicos de una época con la siguiente, siendo el peso del arco

la similitud entre los tópicos. Una vez construido el grafo se eliminan las conexiones débiles en base a un umbral $\zeta \in [0, 1]$ a definir, reteniendo solo aquellas conexiones entre tópicos suficientemente similares entre épocas adyacentes, matemáticamente podemos el arco entre los tópicos $\phi_{t,i}$ y $\phi_{t+1,j}$ si $\rho(\phi_{t,i}, \phi_{t+1,j}) \leq \zeta$. Una ilustración conceptual del grafo de similitud es mostrado en la Figura 4.6, este muestra tres épocas consecutivas.

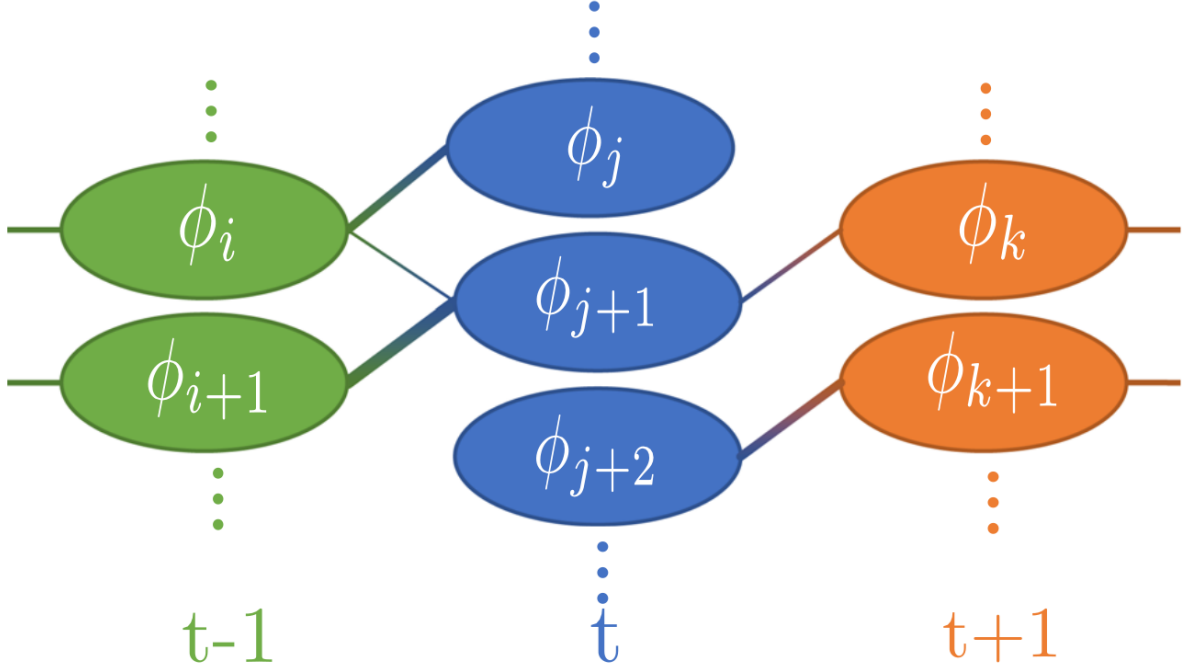


Figura 4.6: Ilustración conceptual del grafo de similitud que modela la dinámica de los tópicos en el tiempo. Un nodo corresponde a un tópico en una época específica; el ancho de los arcos es proporcional a la similitud entre los tópicos, arcos ausentes fueron eliminados por presentar una similitud menor a un umbral.

Esta metodología permite fácilmente detectar desaparición de un tópico, nacimiento de un nuevo tópico, como también dividir o fusionar diferentes tópicos, a continuación se define en detalle cada uno de estos dinamismos:

- **Nacimiento de un tópico:** Si un tópico no tiene ningún arco entrante, por ejemplo, en la Figura 4.6 el tópico ϕ_{j+2} en t .
- **Muerte de un tópico:** Si un tópico no tiene ningún arco saliente, por ejemplo, en la Figura 4.6 el tópico ϕ_j en t .
- **Evolución de un tópico:** Cuando un tópico tiene exactamente un arco de entrada y salida, por ejemplo, en la Figura 4.6 entre las épocas t y $t + 1$ se tiene que el tópico ϕ_{j+2} evoluciona del tópico ϕ_{k+1} .
- **División de un tópico:** Si un tópico tiene más de un arco saliente, por ejemplo, en la Figura 4.6 el tópico ϕ_i de $t - 1$ se divide en $t + 1$ en los tópicos ϕ_j y ϕ_{j+1} .

- **Fusión de un tópico:** Cuando un tópico tiene más de un arco entrante, este tipo de tópicos también pueden ser entendidos como un nuevo tópico, por ejemplo, en la Figura 4.6 los tópicos ϕ_i y ϕ_{i+1} de $t - 1$ forman al tópico ϕ_{j+1} en t .

Un aspecto relevante de esta metodología es definir el umbral de corte ρ , el cual no es fácilmente interpretable, además el umbral depende de la métrica de similitud escogida, dificultando así la comparación entre métricas de similitud. En [Beykikhoshk et al. \(2018\)](#) proponen una alternativa más interpretable para definir el umbral, para esto estiman función de densidad acumulada (cdf) del grafo inicial, donde todos los nodos de una época están conectados con todos los nodos de la época adyacente, sea F_p la cdf sobre las similitudes del grafo inicial, luego sea $\zeta \in [0, 1]$ el punto operante de la cdf, luego eliminamos el arco entre los tópicos $\phi_{t,i}$ y $\phi_{t+1,j}$ si $\rho(\phi_{t,i}, \phi_{t+1,j}) \leq F_p^{-1}(\zeta)$, donde $F_p^{-1}(\zeta)$ es el cuantil ζ de F_p .

4.2.2. Medidas de similitud

Los tópicos son distribuciones de probabilidad sobre un vocabulario fijos de términos. La gran mayoría de medidas de similitud comparan vectores con el mismo dominio y dimensión, esto significa que los tópicos de épocas adyacentes deben compartir el mismo vocabulario, matemáticamente, sea $\phi_{t,i}$ un tópico de la época t y V_t su vocabulario, sea $\phi_{t+1,j}$ un tópico de la época $t + 1$ y V_{t+1} su vocabulario, lo más probable es que existan palabras en V_t que no existan en V_{t+1} y viceversa, para poder comparar tópicos en estas épocas adyacentes se debe construir el vocabulario $V'_{t+1} = V_t \cup V_{t+1}$, luego se aplica *padding* a los vectores $\phi_{t,i}$ y $\phi_{t+1,j}$, es decir, se rellenan con ceros las posiciones de palabras que no están en el vocabulario de su dominio.

La gran desventaja del enfoque anterior es que no captura similitud entre palabras, es decir, dos palabras diferentes que pueden llegar a ser sinónimos ocuparan una posición diferente dentro del vector, siendo no robusta a cuando una palabra esta presente en la época t y no en $t - 1$ por lo que no hay forma de compararla por ejemplo con la palabra de $t - 1$ más similar, por lo que se compara la palabra consigo misma, donde en t tiene un peso distinto de cero y en $t - 1$ un peso nulo. El peor caso sería considerar los vocabularios V_t y V_{t+1} , donde $V_t \cap V_{t+1} = \emptyset$, a pesar de que cada palabra en V_t tiene un sinónimo en V_{t+1} la similitud entre tópicos entre las épocas t y $t + 1$ sería cero.

Para lidiar con el problema anterior en ([Kusner et al., 2015](#)) se propone una medida de distancia llamada Word Mover's Distance (WMD) para comparar dos documentos bajo una representación *bag of words*, donde i y j son los documentos, V_i y V_j los vocabularios, y el peso asociado a cada palabra de un documento es igual a la frecuencia normalizada. Generalizar al caso de tópicos es bastante sencillo, puesto a que estos se construyen bajo una representación *bag of words*, por ejemplo, para comparar el tópico i de la época t con el tópico j de la época $t + 1$, se usan los pesos $\phi_{t,i}$ y $\phi_{t+1,j}$ sobre el vocabulario V_t y V_{t+1} respectivamente. WMD calcula el costo mínimo de transformar un documento en otro, en este caso particular sería el costo mínimo de llevar un tópico a otro, para esto se resuelve un problema de flujo a costo mínimo (MCF), donde los flujos son los pesos $\phi_{t,i}$ y $\phi_{t+1,j}$ y la matriz de costos es una matriz de distancia euclidiana entre los *word embedding* ([Mikolov et al. \(2013\)](#)) de todas las palabras de V_t con V_{t+1} . En la Figura 4.7 se ilustra el espacio en el que viven las palabras de dos tópicos.

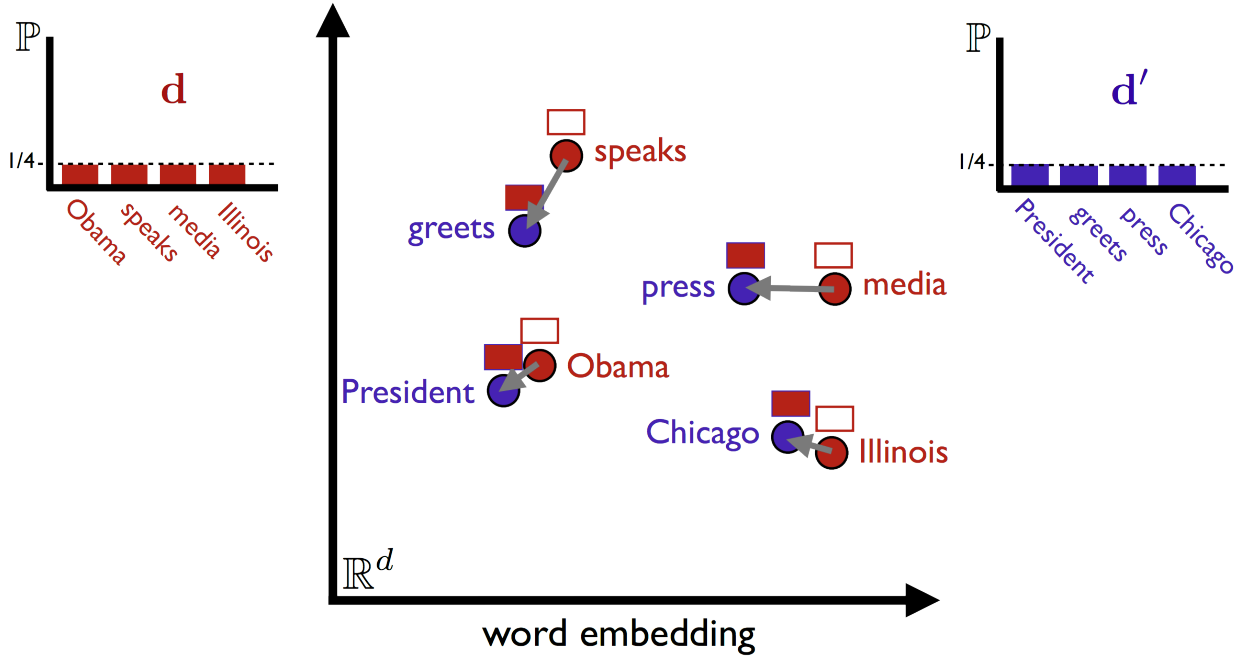


Figura 4.7: Espacio vectorial de los *word embeddings* de las palabras de dos tópicos con un vocabulario de tamaño 4.

Matemáticamente, la WMD entre el tópico i de la época t y el tópico j de la época $t + 1$ viene dado por $WMD(\phi_{i,t}, \phi_{j,t+1})$:

$$\begin{aligned}
 & \underset{T}{\text{minimize}} \sum_{u \in V_t} \sum_{v \in V_{t+1}} c_{u,v} T_{u,v} \\
 & \text{s.t.} \quad \sum_{v \in V_{t+1}} T_{u,v} = \phi_{i,t,u}, u \in V_t \\
 & \quad \sum_{u \in V_t} T_{u,v} = \phi_{j,t+1,v}, v \in V_{t+1} \\
 & \quad T_{u,v} \geq 0, u \in V_t, v \in V_{t+1}
 \end{aligned}$$

Donde $T_{u,v}$ es el flujo que va de la palabra u del tópico i de la época t a la palabra v del tópico j de la época $t + 1$, $\phi_{i,t,u}$, es la probabilidad de la palabra u en el tópico i de la época t , $c_{u,v}$ es el costo de mover una unidad de flujo por el arco (u, v) , el costo entre palabras se mide como la distancia euclidiana entre los *word embedding* de dichas palabras. La primera restricción indica que el flujo que se mueve de una palabra u del tópico i a todas las palabras del tópico j debe sumar su peso ($\phi_{i,t,u}$), la segunda restricción significa que el flujo que se mueve de una palabra v del tópico j a todas las palabras del tópico i debe sumar su peso ($\phi_{j,t+1,v}$). Esta medida de distancia se puede fácilmente transformar en una medida de similitud $\rho(\phi_{i,t}, \phi_{j,t+1}) = \frac{1}{1 + WMD(\phi_{i,t}, \phi_{j,t+1})}$, notar que si la WMD es 0 la similitud es 1 y si es ∞ la similitud es 0.

Capítulo 5

Experimento

5.1. Datos

Para este experimento se cuenta con las fuentes de datos de la Asociación de Aseguradores de Chile (AACH), corresponde a los relatos que las víctimas del robo de sus vehículos dan a las aseguradoras, lo cual corresponde a 49.015 relatos entre el 2011 y 2016.

Para el uso de WMD es necesario contar con *words embeddings*, para esto se utilizaron los *embeddings* de (Pérez, 2019), estos *embeddings* fueron obtenidos utilizando el algoritmo FastText (Bojanowski et al., 2017) sobre el corpus Spanish Billion Word Corpus (SBWC) (Cardellino, 2019). FastText en comparación a otros enfoques para extraer *embeddings* representa los *tokens* a través de n-gramas de caracteres, de esta manera se pueden obtener *embeddings* de *tokens* no vistos durante el entrenamiento a partir de los *embeddings* de los caracteres que lo componen.

5.2. Procesamiento

En minería de texto con el objetivo de extraer el core de palabras del corpus se recurre métodos para reducir el vocabulario, la reducción del vocabulario mejora la significancia estadística de los modelos, puesto que se obtiene un mejor balance entre cantidad de parámetros y cantidad de observaciones, por otro lado puede verse facilitada la interpretación de los tópicos al remover palabras que aportan poca información.

El paso cero en el procesamiento de textos es tokenizar, la tokenización es una operación sobre una cadena de caracteres (*string*) que consiste en dividir el *string* en un conjunto de términos, en este caso la división se hizo por el carácter espacio, como resultado de esto se obtiene una lista de elementos, a cada elemento de esta lista se le denomina *token* que en términos simples puede considerarse como una palabra para el ejemplo mencionado.

Luego, en el primer nivel de procesamiento no interesa hacer distinción entre mayúsculas o minúsculas¹, por ende, los caracteres de cada token son llevados a minúscula, también se

¹ En análisis de sentimiento puede ser interesante ya que las personas suelen expresar mensajes de enfado con letras capitales, por lo que las letras capitales añaden información al análisis.

eliminaron caracteres y tokens que no aportan información, como símbolos de puntuación, correos electrónicos y tokens que contienen números. En la figura 5.1 se observa la distribución acumulada de los tokens del corpus a este nivel de procesamiento, adicionalmente se tiene que el 50 % de los *tokens* del vocabulario ocurren una sola vez, el 80 % tiene una ocurrencia menor o igual a 5 y el 95 % de la distribución acumulada puede ser explicada con 4199 *tokens* (9 %) del vocabulario, se concluye que la distribución es sumamente pesada y es necesario recurrir a métodos adicionales para su reducción.

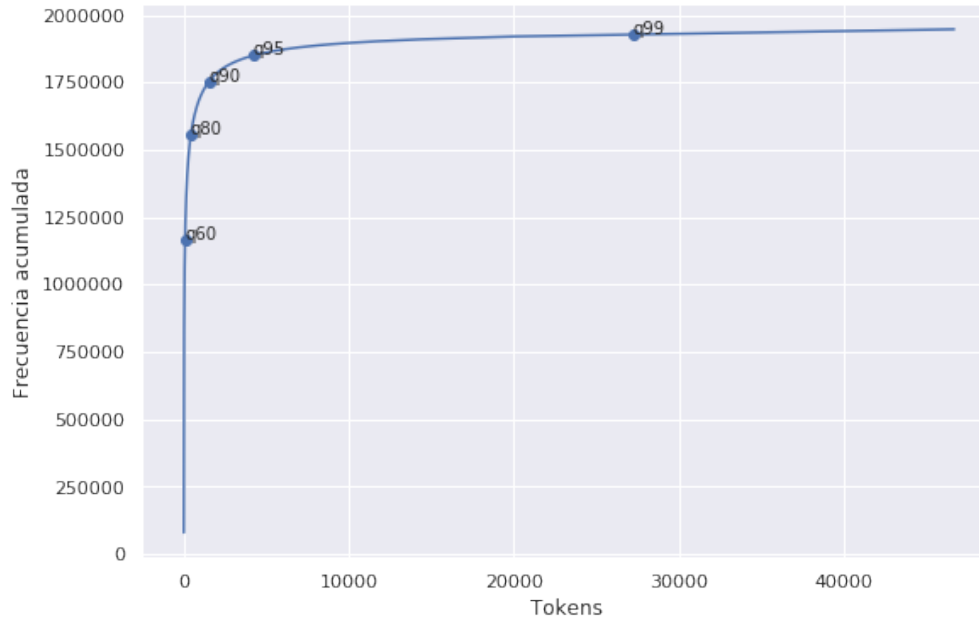


Figura 5.1: Frecuencia acumulada de los tokens únicos aplicando hasta el primer paso de procesamiento. El eje horizontal es el acumulado de tokens únicos en orden decreciente de ocurrencia. Los puntos corresponden a los cuantiles 60 %, 80 %, 90 %, 95 % y 99 %.

En el segundo nivel de procesamiento se eliminaron las stopwords, palabras que aportan poca información, como artículos, preposiciones y conectores, para esto se utilizó la lista de *stopwords* disponible en el paquete NLTK de Python (Bird et al., 2009) la cual cuenta con 313 palabras. Además, esta lista de *stopwords* se alimentó con *stopwords* contextuales, palabras específicas del corpus que aportan poca información, para esto se hizo un etiquetado de las 1000 palabras más frecuentes del corpus incorporando 417 nuevas palabras, algunos ejemplos son palabras que hacen reverencia a vehículo y robo, puesto que todos los documentos corresponden a robos de vehículos.

El tercer nivel de procesamiento consiste en normalizar los tokens para reducir aún más el vocabulario, como métodos de normalización los más utilizados son *stemming* y lematización. *Stemming* es el proceso de llevar una palabra a su raíz (*stem*), en la práctica *stemming* consiste en aplicar un algoritmo basado en ciertas reglas gramaticales para extraer sufijos (Porter et al., 1980), como desventaja es que stemming no tiene en cuenta el contexto de la palabra por lo que la raíz obtenida puede no corresponder a la raíz verdadera de la palabra, además, para el caso de modelamiento de tópicos los tópicos se vuelve más difícil de inter-

pretar, por un lado porque palabras con significado completamente distinto terminan con la misma raíz o porque la raíz encontrada no tiene un significado claro. Por otro lado, lematización es el proceso de agrupar juntas las formas flexionadas de una palabra para que puedan analizarse como un elemento, identificado como lema, su diferencia principal con *stemming* es que opera con conocimiento del contexto de la palabra para discriminar entre palabras que tienen significado diferente dependiendo del *part of speech tagging* (POST) y de una tabla de búsqueda (*lookup table*). Como método de normalización se decidió utilizar lematización en vez de stemming debido a que tiene menos impacto en la interpretación de los tópicos, sin embargo es una operación más intensiva debido a que stemming es un algoritmo basado en reglas simples mientras que en lematización se suele usar redes neuronales recurrentes (RNNs) para el POST y una vez determinadas las etiquetas gramaticales de las palabras en un documento se utiliza una *lookup table* para encontrar el lema correspondiente. La implementación de lematización utilizada es la implementación de lematización en español del paquete spaCy de Python (Honnibal and Montani, 2017).

El cuarto y último nivel de procesamiento corresponde a eliminar tokens con baja frecuencia, puesto que el modelo no será capaz de levantar patrones en tokens que aparecen una única vez o con una ocurrencia poco significativa, luego, como el corpus está particionado en *slices*, se eliminaron aquellos tokens que aparecen en menos de 5 documentos dentro de un *slice*.

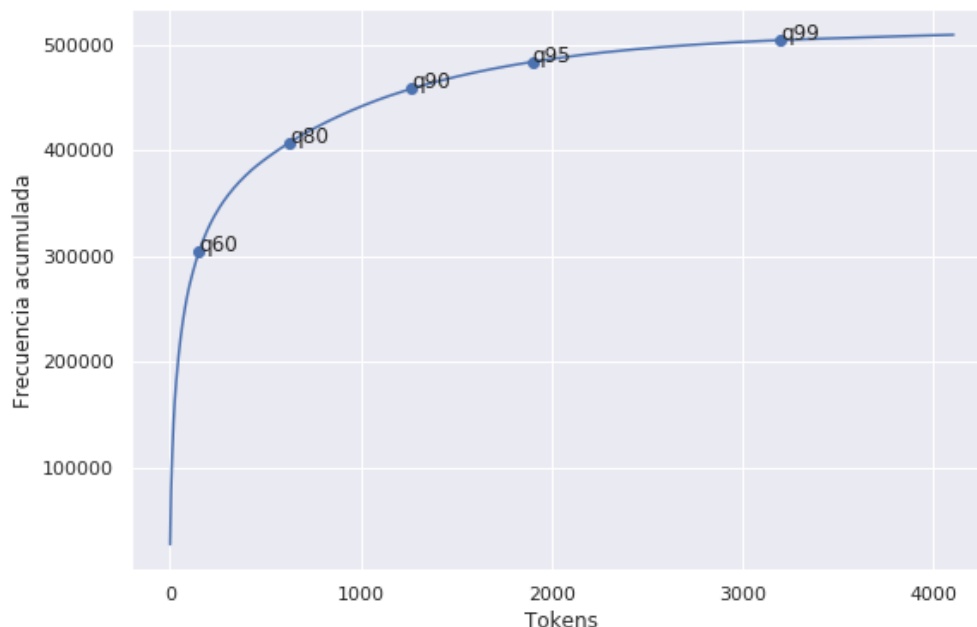


Figura 5.2: Frecuencia acumulada de los tokens únicos aplicando hasta el cuarto nivel de procesamiento. El eje horizontal es el acumulado de tokens únicos en orden decreciente de ocurrencia. Los puntos corresponden a los cuantiles 60 %, 80 %, 90 %, 95 % y 99 %.

En la figura 5.2 se presenta la distribución acumulada del vocabulario hasta el cuarto nivel de procesamiento, en donde se observa que la cola de distribución es bastante menos pesada que bajo el primer nivel de procesamiento, además, como se observa en la tabla 5.1 el tamaño del vocabulario se redujo a menos de un décimo del vocabulario obtenido bajo el

primer nivel de procesamiento y es menos de un décimo del tamaño del corpus, por lo que bajo este nivel de procesamiento es posible desarrollar modelos con mayor fuerza estadística.

| procesamiento | documentos | vocabulario | tokens |
|---------------|------------|-------------|-----------|
| raw | 49.015 | 79.327 | 2.030.980 |
| ch | 49.011 | 46.708 | 1.947.235 |
| ch+s+l+f | 47.993 | 4.106 | 508.987 |

Tabla 5.1: Estadísticas del corpus bajo distintos niveles de procesamiento, **raw**: sin procesamiento, **ch**: eliminación de símbolos de puntuación, correos electrónicos y tokens con números, **ch+s+l+f**: además incluye eliminación de stopwords (s), lematización (l) y eliminación de tokens con baja ocurrencia (f).

En la tabla 5.2 se muestra el detalle del vocabulario para cada una de las *slices* tras procesar el corpus, de aquí se extrae que en promedio un 21.28 % del vocabulario se olvida de una época a otra y un 28.19 % es nuevo, es otras palabras, en promedio alrededor de un 50 % del vocabulario no es común entre tópicos de épocas adyacentes, esto justifica la necesidad de utilizar medidas de similitud que capturen la similitud entre palabras de épocas adyacentes ante la renovación que sufre el vocabulario en el tiempo.

| slice | old_vocabulary | new_vocabulary | %old_tokens | %new_tokens |
|-------|----------------|----------------|-------------|-------------|
| 2 | 1.919 | 1.986 | 23,35 | 26.84 |
| 3 | 1.986 | 2.092 | 22,61 | 27.95 |
| 4 | 2.092 | 2.414 | 18,21 | 33.60 |
| 5 | 2.414 | 2.629 | 19,80 | 28.71 |
| 6 | 2.629 | 2.666 | 22,44 | 23.85 |

Tabla 5.2: Evolución del vocabulario en el tiempo, **old_vocabulary**: corresponde al vocabulario del período $t - 1$, **new_vocabulary**: corresponde al vocabulario del período t , **%old_tokens**: porcentaje de tokens del período $t - 1$ que ya no están en el período t y **%new_tokens**: porcentaje de tokens del período t que no están en el período $t - 1$.

5.3. Análisis cuantitativo de resultados

5.4. Análisis cualitativo de resultados

Capítulo 6

Conclusiones

Bibliografia

- Pete Chapman, Julian Clinton, Randy Kerber, Thomas Khabaza, Thomas Reinartz, Colin Shearer, Rudiger Wirth, et al. Crisp-dm 1.0: Step-by-step data mining guide. *SPSS inc*, 9:13, 2000.
- Susan T Dumais. Latent semantic analysis. *Annual review of information science and technology*, 38(1):188–230, 2004.
- Wei Xu, Xin Liu, and Yihong Gong. Document clustering based on non-negative matrix factorization. In *Proceedings of the 26th annual international ACM SIGIR conference on Research and development in informaion retrieval*, pages 267–273, 2003.
- David M Blei, Andrew Y Ng, and Michael I Jordan. Latent dirichlet allocation. *Journal of machine Learning research*, 3(Jan):993–1022, 2003.
- Yee W Teh, Michael I Jordan, Matthew J Beal, and David M Blei. Sharing clusters among related groups: Hierarchical dirichlet processes. In *Advances in neural information processing systems*, pages 1385–1392, 2005.
- Keith Stevens, Philip Kegelmeyer, David Andrzejewski, and David Buttler. Exploring topic coherence over many models and many topics. In *Proceedings of the 2012 Joint Conference on Empirical Methods in Natural Language Processing and Computational Natural Language Learning*, pages 952–961. Association for Computational Linguistics, 2012.
- David M Blei and John D Lafferty. Dynamic topic models. In *Proceedings of the 23rd international conference on Machine learning*, pages 113–120, 2006.
- Xuerui Wang and Andrew McCallum. Topics over time: a non-markov continuous-time model of topical trends. In *Proceedings of the 12th ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 424–433, 2006.
- Amr Ahmed and Eric P Xing. Timeline: A dynamic hierarchical dirichlet process model for recovering birth/death and evolution of topics in text stream. *arXiv preprint arXiv:1203.3463*, 2012.
- Andrew T Wilson and David G Robinson. Tracking topic birth and death in lda. *Sandia National Laboratories*, 2011.
- Adham Beykikhoshk, Ognjen Arandjelović, Dinh Phung, and Svetha Venkatesh. Discovering topic structures of a temporally evolving document corpus. *Knowledge and Information Systems*, 55(3):599–632, 2018.
- Thomas L Griffiths and Mark Steyvers. Finding scientific topics. *Proceedings of the National academy of Sciences*, 101(suppl 1):5228–5235, 2004.

- Matt Kusner, Yu Sun, Nicholas Kolkin, and Kilian Weinberger. From word embeddings to document distances. In *International conference on machine learning*, pages 957–966, 2015.
- Tomas Mikolov, Ilya Sutskever, Kai Chen, Greg S Corrado, and Jeff Dean. Distributed representations of words and phrases and their compositionality. In *Advances in neural information processing systems*, pages 3111–3119, 2013.
- Jorge Pérez. Fasttext embeddings from SBWC. <https://github.com/dccuchile/spanish-word-embeddings#fasttext-embeddings-from-sbwc>, 2019.
- Piotr Bojanowski, Edouard Grave, Armand Joulin, and Tomas Mikolov. Enriching word vectors with subword information. *Transactions of the Association for Computational Linguistics*, 5:135–146, 2017.
- Cristian Cardellino. Spanish Billion Words Corpus and Embeddings, August 2019. URL <https://crscardellino.github.io/SBWCE/>.
- Steven Bird, Ewan Klein, and Edward Loper. *Natural language processing with Python: analyzing text with the natural language toolkit*. O'Reilly Media, Inc., 2009.
- Martin F Porter et al. An algorithm for suffix stripping. *Program*, 14(3):130–137, 1980.
- Matthew Honnibal and Ines Montani. spaCy 2: Natural language understanding with Bloom embeddings, convolutional neural networks and incremental parsing. To appear, 2017.