



Análisis de Sentimientos

Ing. Juan M. Rodríguez
jmrodriguez@lsia.fi.uba.ar

Análisis de Sentimientos - introducción



Criticas

Decepcionante.

=> **Negativo**

Aburrida.

=> **Negativo**

Opiniones

Personajes memorables y bien desarrollados.

=> **Positivo**

Comentarios

Una gran puesta en escena que no

=> **Positivo**

defraudará.

=> **Negativo**

Increíblemente predecible.

Análisis de Sentimientos - introducción



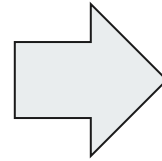
velocidad de impresión

precio

facilidad de uso

soporte al cliente

etc.



Aspectos

o

Atributos

Análisis de Sentimientos - aspectos

Hoteles

Restaurantes

Customer Rating

76%



Cleanliness:



Staff:



Location:



Value:



Recommended:

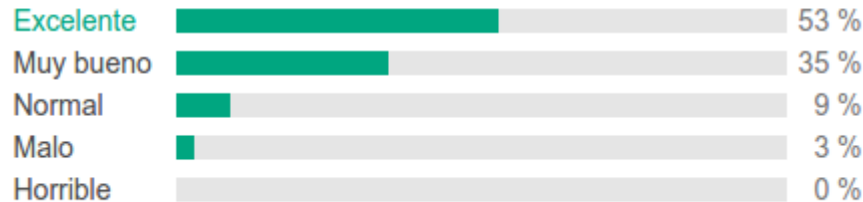


Rated By 247
Customer/S

[Read All Ideal Hotel
Paris Reviews](#)

4,5

53 opiniones



Análisis de Sentimientos - tareas complejas



Estimar la confianza del consumidor:

"La confianza del consumidor es un **indicador económico** que mide el grado de optimismo que los consumidores sienten sobre el estado general de la economía y sobre su situación financiera personal. Qué tan seguras se sienten las personas sobre la estabilidad de sus ingresos determina sus **actividades de consumo** y por lo tanto sirve como uno de los indicadores claves en la forma general de la economía." (*Wikipedia*)

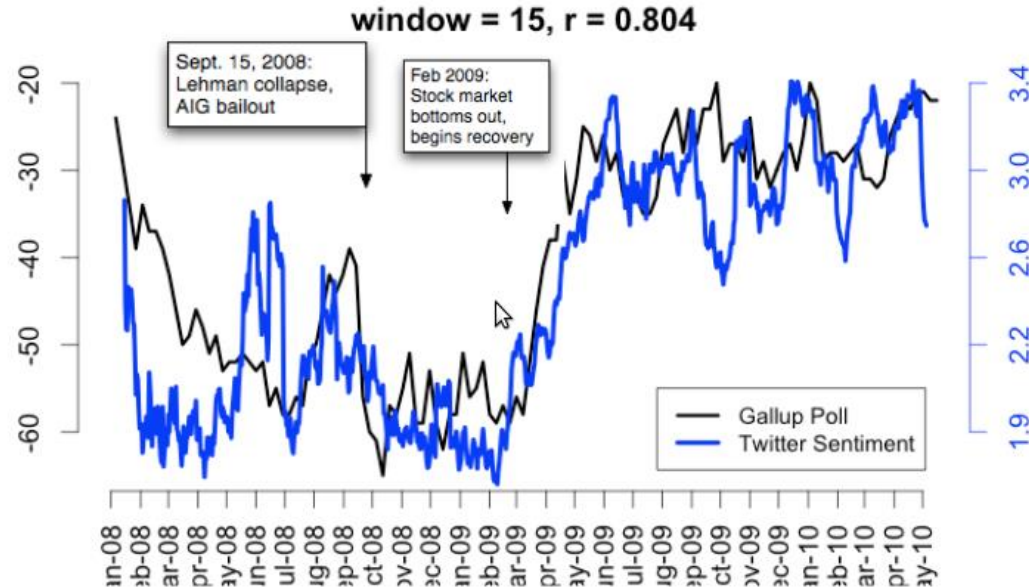
Análisis de Sentimientos - tareas complejas

Estimar la confianza del consumidor:

Gallup

Twitter

relación: 0.804



Análisis de Sentimientos - tareas complejas



predecir el mercado de valores

se usaron dos herramientas

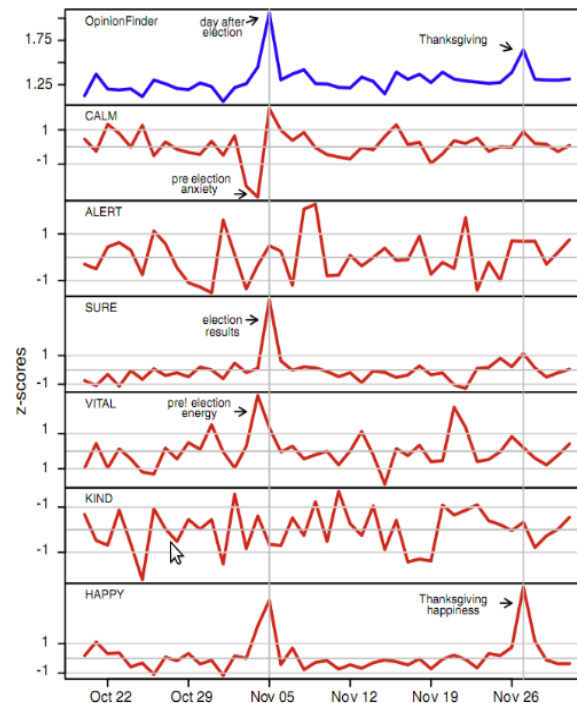
- **OpinionFinder** que mide opinión negativa versus opinión positiva
- **Google-Profile of Mood States** (perfil Google de estados de ánimo) que mide el humor en término de 6 dimensiones: calmado, alerta, seguro, vital, amable y feliz.

Análisis de Sentimientos - tareas complejas

predecir el mercado de valores:

Dos herramientas

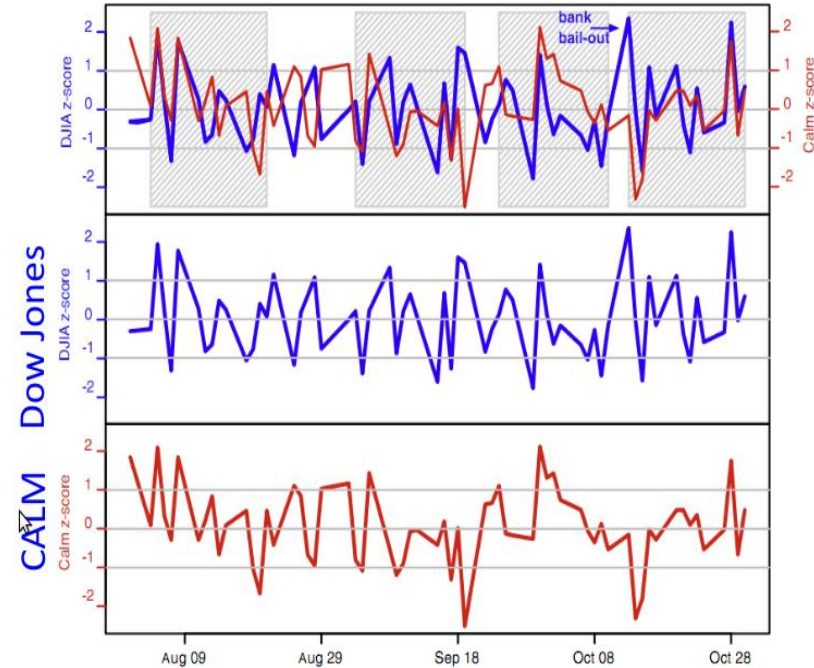
- **OpinionFinder** que mide opinión negativa versus opinión positiva
- y
- **Google-Profile of Mood States** (perfil Google de estados de ánimo) que mide el humor en término de 6 dimensiones: calmado, alerta, seguro, vital, amable y feliz.



Análisis de Sentimientos - tareas complejas

Predecir el mercado de valores:

Se descubrió, entre otras cosas, que "la calma" podía predecir el Índice Dow Jones con 3 días de anticipación



Análisis de Sentimientos - herramientas *online*

<http://www.sentiment140.com>

Sentiment140

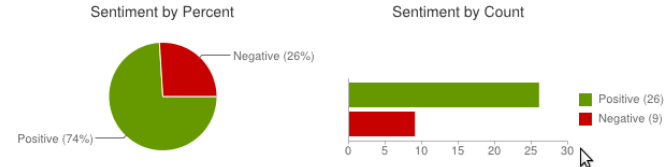
Tweet 696

Me gusta

+1 156

Spanish ▾

Sentiment analysis for Samsung



Tweets about: Samsung

[AllisonGualdron](#): Más phablets y tablets de gran pantalla y alta resolución, lo prepara Sa

Posted: 20 seconds ago

[jsnp1982](#): @SkyAlertSoporte en **Samsung Ace** no deja actualizarlo, lo borre y al buscarl

Posted: 3 minutes ago

[rodoaguilar](#): Tengo más de un mes que rooté mi galaxy tab 2 de **Samsung** y que le insta

Análisis de Sentimientos - sinónimos



Extracción de opinión (opinion extraction)

Minería de opiniones (opinion mining.)

Minería de sentimientos (sentiment mining,)

Análisis de subjetividad (subjectivity analysis)

Análisis de Sentimientos - usos



¿Para qué podemos utilizar el análisis de sentimientos?

- **Películas:** ¿Es esta crítica positiva o negativa?
- **Productos:** ¿Qué piensa la gente del nuevo iPhone?
- **Sentimientos públicos:** ¿Cómo es la confianza del consumidor?, ¿crece de forma impar? , etc.
- **Política:** ¿Qué piensa la gente acerca de este candidato o de esta situación?
- **Predicción:** predecir el resultado de una elección o una tendencia de mercado a partir de los sentimientos.

Tipología de Scherer de los estados afectivos

- **Emoción:** Respuesta relativamente corta del organismo a estímulos externos. Ejemplos de emoción son la **ira**, la **tristeza**, la **alegría**, el **miedo**, la **vergüenza**, el **orgullo**, la **alegría** y la **desesperación**.
- **Estado de ánimo:** sentimiento subjetivo de baja intensidad y larga duración. Ejemplos de estados de ánimo : **alegre**, **triste**, **irritable**, **apático**.
- **Postura interpersonal:** posición afectiva respecto a otra persona en una interacción específica. Ejemplos de posturas interpersonales son: **distante**, **frío**, **cálido**, **de apoyo** y **de desprecio**.
- **Actitudes:** preferencia o predisposición de una persona respecto a otras personas u objetos. Ejemplos de actitudes son: **simpatía**, **amor**, **odio**, **deseo** y **valoración**.
- **Rasgos de personalidad:** tendencias en el comportamiento típico de una persona. Ejemplo de rasgos de personalidad: **nervioso**, **ansioso**, **imprudente**, **taciturno**, **hostil**, **envidioso** y **celoso**.

Análisis de Sentimientos - definición



- **Actitudes:** preferencia o predisposición de una persona respecto a otras personas u objetos. Ejemplos de actitudes son: **simpatía, amor, odio, deseo y valoración.**

Cuando se realiza análisis de sentimientos, en verdad lo que se está haciendo es detectando actitudes, es decir se está detectando la "preferencia o predisposición de una persona respecto a otras personas u objetos"

Análisis de Sentimientos - definición de tareas



1. El portador de dicha actitud
2. El destinatario de dicha actitud.
3. El tipo de actitud:
 - a. simpatía, **amor**, **odio**, **deseo** y **valoración** (lista de candidatas)
 - b. o una polaridad ponderada: **positiva**, **negativa** o **neutral** (y a veces un valor asociado)
4. El texto que contiene la actitud (documento u oración)

Análisis de Sentimientos - definición de tareas



Resumen

- Una **tarea sencilla** de análisis de sentimientos consiste en:
 - determinar si la actitud de un texto es positiva o negativa.
- Una **tarea un poco más compleja** de análisis de sentimientos consiste en:
 - puntuar la actitud de un texto de 1 a 5.
- Una **tarea realmente avanzada** de análisis de sentimientos consiste en:
 - detectar el portador, el destinatario y el tipo de actitud de un texto.

Algoritmo de Pang y Lee

Detección de la polaridad

<http://www.cs.cornell.edu/people/pabo/movie-review-data/>

Algoritmo de Pang y Lee

Hace mucho tiempo no veía una peli tan buena como esta. Original y hermosa desde todo punto de vista, las canciones, los chistes y los efectos visuales. Increíble como mejoraron en los últimos años con las expresiones de los personajes animados! Una obra maestra desde donde se la mire! No entendía cómo Monster University no había sido nominada al oscar, pero luego de ver Frozen, entiendo que Monster's no tenía nada que hacer. Es de las pelis que, para aquellos que crecimos con los clásicos de disney como el rey leon, tarzan o el jorobado de Notre dame, nos reavivan el -:- que tenemos adentro. Si o si verla en 3D.



Algoritmo de Pang y Lee

La película es malísima. Está entre las peores que he visto en mi vida, sino es la peor. Es lineal, predecible, aburrida. Solo apta para menores de 13 años. Deberían dedicarse a otra cosa todos los que intervinieron para que exista. Realmente es muy recomendable que NO LA VEAN. Después no digan que no les advertí.



Algoritmo de Pang y Lee



1. Tokenización del texto
2. Extracción de características (palabras o frases claves)
3. Clasificación utilizando distintos algoritmos de clasificación:
 - a. Naïve Bayes
 - b. MaxEnt
 - c. SVM

Tokenización, problemas comunes



- Lidar con los **tags XML o HTML**
- Tener que reconocer las **marcas de Twitter** (si queremos sacar información de ahí) como los nombres de usuario y los *hash tags*
- **El uso de mayúsculas.** Generalmente nos va a interesar conservar las mayúsculas de las palabras en las distintas fases del algoritmo.
- Números de teléfono y fechas.
- **Emoticones:** es muy útil detectar los emoticones cuando se está haciendo análisis de sentimientos.

Extracción de características



En esta etapa tenemos dos problemas:

- ¿Cómo lidiar con la negación?
 - **No** me gustó esta película.
 - Me gustó esta película.
- ¿Qué conviene usar?
 - todas las palabras
 - solo los adjetivos

Extracción de características



¿Qué conviene usar?

- todas las palabras
- solo los adjetivos

Se demostró que al menos con la información de **IMDB**, es **conveniente utilizar todas las palabras**. Se obtienen así mejores resultados y en términos generales diría que siempre conviene utilizar todas las palabras ya que a veces los sustantivos y los verbos nos dan información valiosa sobre el juicio de valor de un crítica.

Extracción de características

Cómo lidiar con la negación

No me gustó esta película, pero yo...



No NO_me NO_gustó NO_esta NO_película pero yo

Clasificación - Naïve Bayes



$$C_{\text{map}} = \underset{c_j \in C}{\operatorname{argmax}} P(c_j) \prod_{i \in \text{Posiciones}} P(x_i | c_j)$$

Los x_i son las características del documento, es decir las palabras. Ya que un documento es representado como una "bolsa" de palabras.

Clasificación - Naïve Bayes



Y cada una de las $P(x_i | c_j)$ es calculada, usando *Laplace smoothing*:

$$P'(w | c) = \frac{\text{cantidad}(w | c) + 1}{\text{cantidad}(c) + |V+1|}$$

Clasificación - Naïve Bayes multinomial binarizada (o booleana)

Antes de comenzar a calcular las probabilidades de las clases y a contar las palabras vamos a recorrer uno por uno todos los documentos en el conjunto de entrenamiento y prueba y a **eliminar las palabras duplicadas**.

En términos generales **esta variante del algoritmo da mejores resultados** que la versión tradicional que cuenta todas las ocurrencias de las palabras.

Clasificación - Entrenamiento

Primera iteración:

P	E	E	E	E	E	E	E	E	E
---	---	---	---	---	---	---	---	---	---

Segunda iteración:

E	P	E	E	E	E	E	E	E	E
---	---	---	---	---	---	---	---	---	---

Tercera iteración:

E	E	P	E	E	E	E	E	E	E
---	---	---	---	---	---	---	---	---	---

Cuarta iteración:

E	E	E	P	E	E	E	E	E	E
---	---	---	---	---	---	---	---	---	---

Quinta iteración:

E	E	E	E	P	E	E	E	E	E
---	---	---	---	---	---	---	---	---	---

Sexta iteración:

E	E	E	E	E	P	E	E	E	E
---	---	---	---	---	---	---	---	---	---

...

Buscar el que maximice la precisión, el recall y la medida F1

Algoritmo de Pang y Lee - problemas



Sutilezas:

Si usted está leyendo esto porque es su fragancia favorita, por favor úsela exclusivamente en su casa y cierre bien las ventanas.

Expectativas frustradas:

La película debería ser excelente ya que cuenta con grandes actores y una banda sonora fantástica, sin embargo es terriblemente aburrida.

Lexicón de sentimientos

The General Inquirer



- 1915 palabras en la categoría: “positivas”
- 2291 palabras en la categoría “negativas”
- Clasificaciones complejas como por ejemplo: Fuerte vs. Débil o Activa vs. Pasiva
- Está en inglés
- Es gratis para su uso en investigación.

<http://www.wjh.harvard.edu/~inquirer/homecat.htm>

LIWC (Linguistic Inquiry and Word Count)



- 2300 palabras
- Más de 70 clases.
- Soporta idioma español.
- Tiene clasificaciones complejas.
- No es gratuito, tiene dos versiones con costos de US\$29.95 y US\$89.95.

<http://liwc.net/liwcspanol/descriptiontable1.php>

MPQA Subjectivity Cues Lexicon



- Existe desde 2006
- 2718 palabras en la categoría “positivas”
- 4912 en la categoría “negativas”
- Está en idioma inglés
- Indica la intensidad de la palabra (fuerte/ debil)
- Se distribuye bajo licencia GNU GPL

http://mpqa.cs.pitt.edu/lexicons/subj_lexicon/

Bing Liu Opinion Lexicon



- Existe desde 2004
- 2006 palabras en la categoría “positivas”
- 4783 en la categoría “negativas”
- Está en idioma ingles
- Solo tiene las categorías: positiva y negativa

- Página del creador: <http://www.cs.uic.edu/~liub/>
- Lexicón: <http://www.cs.uic.edu/~liub/FBS/opinion-lexicon-English.rar>

SentiWordNet



- Clasificación: positiva, negativa u objetiva (pudiendo una palabra tener al mismo tiempo valores negativos y positivos)
- Está basada en WordNet (3.0 la última versión)
- Está en idioma inglés
- Se distribuye bajo una licencia: "ShareAlike" de Creative Commons

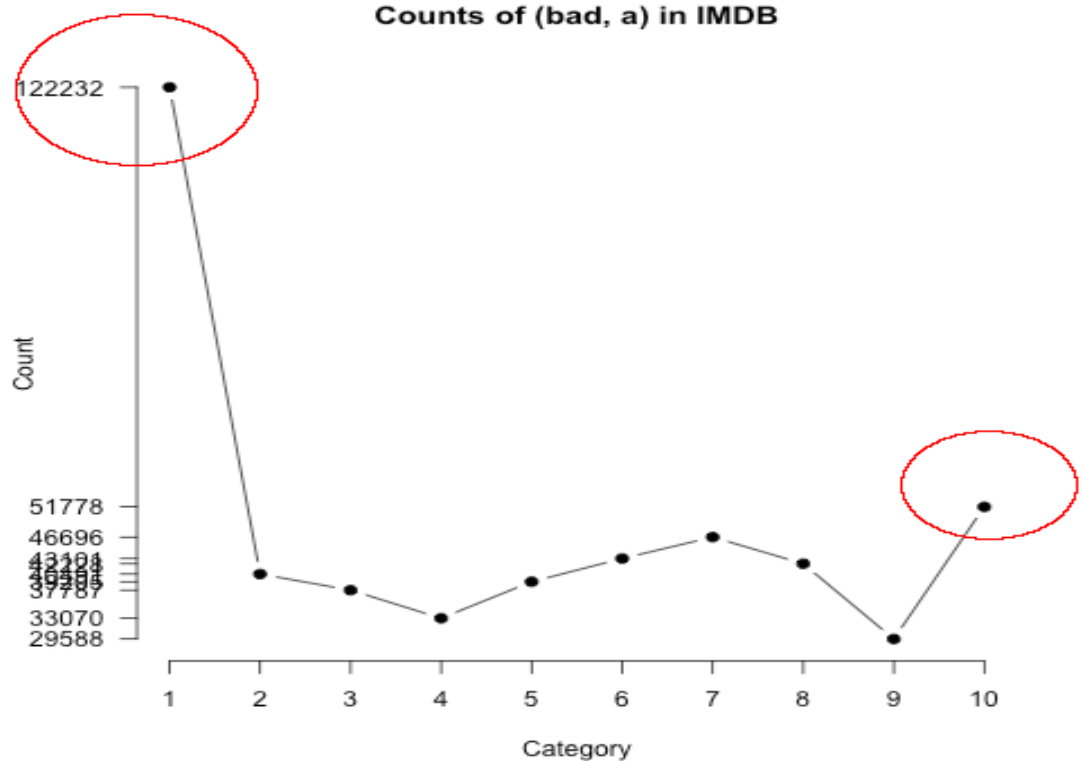
<http://sentiwordnet.isti.cnr.it/>

Desacuerdos entre distintos Lexicons

	MPQA	Opinion Lexicon	Inquirer	SentiWordNet	LIWC
MPQA	—	33/5402 (0.6%)	49/2867 (2%)	1127/4214 (27%)	12/363 (3%)
Opinion Lexicon		—	32/2411 (1%)	1004/3994 (25%)	9/403 (2%)
Inquirer			—	520/2306 (23%)	1/204 (0.5%)
SentiWordNet				—	174/694 (25%)
LIWC					—

Desacuerdos entre distintos Lexicons

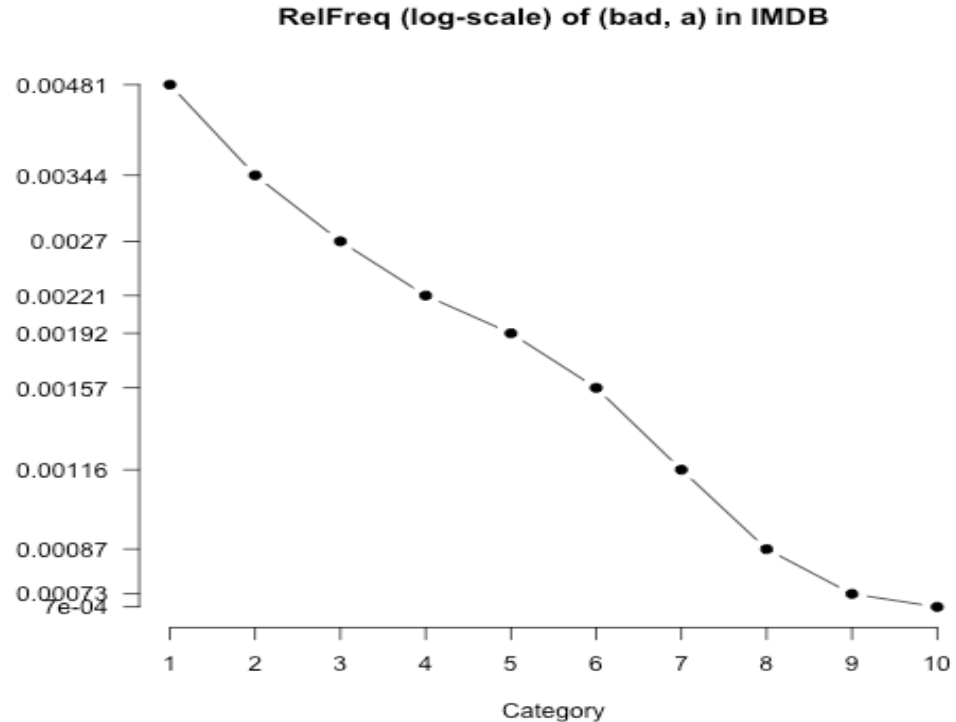
¿Por qué observamos estas diferencias?



Desacuerdos entre distintos Lexicons

Frecuencia relativa

<http://sentiment.christopherpotts.net/lexicons.html>



Creación de un lexicon propio



Qué necesitamos:

- Un puñado de ejemplos previamente clasificados
- Algunas reglas escritas a mano que identifiquen ciertos patrones en una frase.

Algoritmo de Hatzivassiloglou y McKeown para la ampliación de un lexicón

Algoritmo de Hatzivassiloglou y McKeown



Adjetivos unidos por "y" tienen la misma polaridad:

- Justo y legitimo
- corrupto y brutal

Adjetivos unidos por "pero" tienen distinta polaridad:

- justo pero brutal
- corrupto pero legitimo
- hermosa pero malvada

Algoritmo de Hatzivassiloglou y McKeown



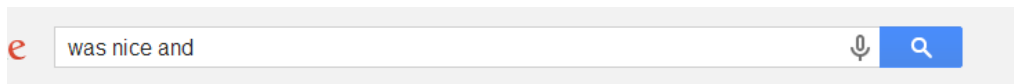
Teniendo esta idea en mente, idearon un **algoritmo en 4 pasos**:

- 1) Construyeron un Lexicón a mano con 1336 adjetivos:
657 positivos y
679 negativos
- 2) Buscaron en Google cada uno de los adjetivos con la formula:
“**was <adjetivo> and**” y recolectaron la palabra que seguía a continuación.

Luego lo repitieron con "**but**" en vez de "and".

Algoritmo de Hatzivassiloglou y McKeown

Ejemplo: "was nice and":



Web Imágenes Videos Noticias Shopping Más ▾ Herramientas de búsqueda

Página 3 de alrededor de 4,450,000,000 resultados (0,32 segundos)

Traducción it **was nice** to see you español | Diccionario ...

diccionario.reverso.net/ingles.../it%20was%20nice%20to%20see%20you ▾

traducción it **was nice** to see you en espanol, diccionario Ingles - Espanol, definición, consulte también 'Nice', 'Nice', 'niche', 'NIC'

Good Omens: The **Nice and Accurate** Prophecies of Agnes ...

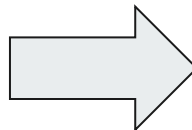
www.amazon.com › ... › Humor › Satire ▾ Traducir esta página

Good Omens: The **Nice and Accurate** Prophecies of Agnes Nutter, Witch [Neil Gaiman, Terry Pratchett] on Amazon.com. *FREE* shipping on qualifying offers.

The room **was nice and clean** - TripAdvisor

www.tripadvisor.co.uk/LocationPhotoDirectLink-g5... ▾ Traducir esta página

Sunprime Coral Suites And Spa, Playa de las Americas Picture: The room **was nice and clean** - Check out TripAdvisor members' 1159 candid photos and ...

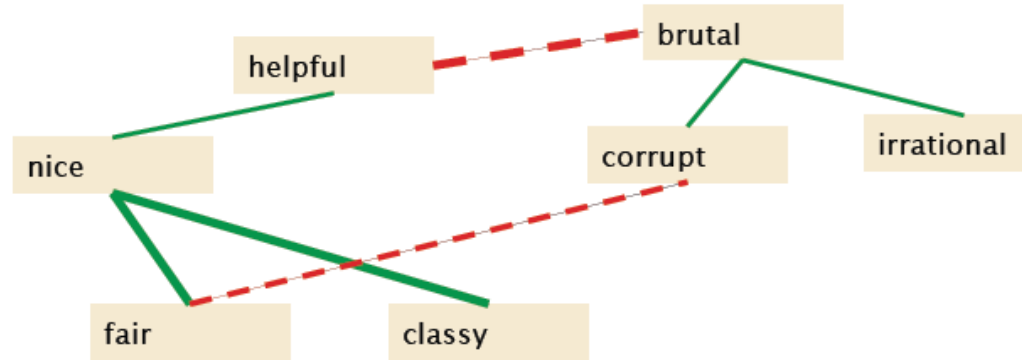


Accurate

Clean

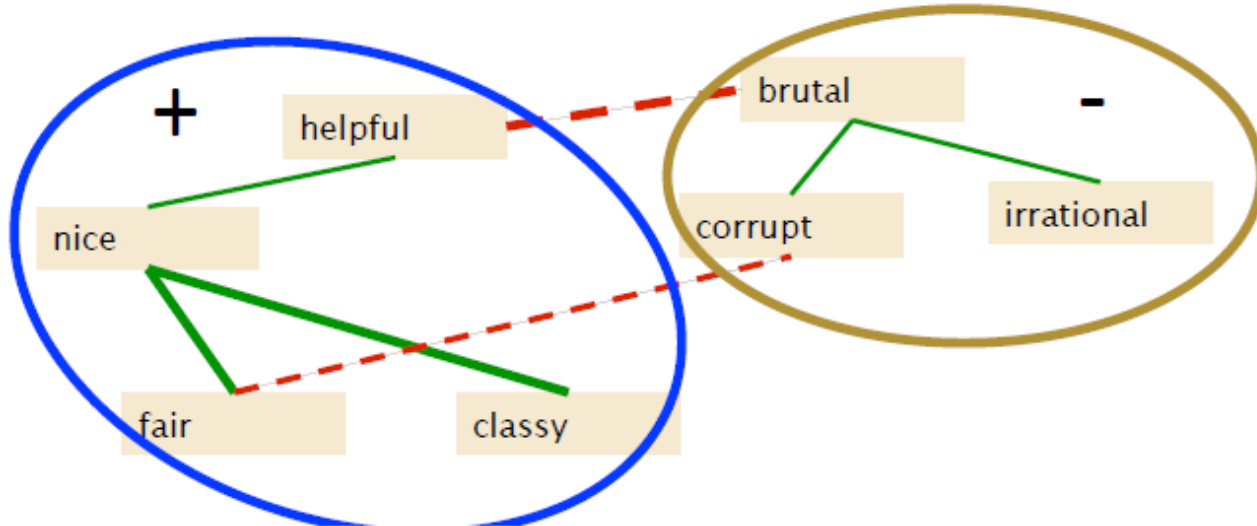
Algoritmo de Hatzivassiloglou y McKeown

3) Construyeron un mapa que vinculaba las palabras similares entre sí, mostrando también las que tenían sentido opuesto:



Algoritmo de Hatzivassiloglou y McKeown

4) Finalmente buscaron una forma de separar el mapa creado intentando que queden dos conjuntos bien diferenciados:



Algoritmo de Hatzivassiloglou y McKeown



Si bien lograron ampliar considerablemente el Lexicón original, el nuevo Lexicón contenía algunos **errores**, es decir palabras mal catalogadas.

Es por ello que este algoritmo necesariamente necesita de un paso extra que consista en la **revisión de los datos obtenidos**.

Algoritmo de Turney para obtener la polaridad de frases

Algoritmo de Turney



Este algoritmo se puede desglosar en 3 pasos principales:

1. Extraer frases de opiniones/críticas (*reviews*) y armar un Lexicón de frases
2. Aprender la polaridad de cada frase
3. Puntuar las críticas según el promedio de las polaridades de sus frases

Algoritmo de Turney - extracción de frases



Primer Palabra	Segunda Palabra	Tercer Palabra (no se extrajo)
Adjetivo	Sustantivo (plural o singular)	Cualquier palabra
Adverbio	Adjetivo	No Sustantivo
Adjetivo	Adjetivo	No Sustantivo
Sustantivo (plural o singular)	Adjetivo	No Sustantivo
Adverbio	Verbos	Cualquier palabra

Algoritmo de Turney - polaridad de las frases



Para verificar la polaridad de las frases, se verificó cuan cerca aparecían estas de palabras con polaridad ya conocida como por ejemplo: "**excelente**" y "**pobre**".

La idea detrás de esto es que la **co-ocurrencia** de una frase junto con la palabra "**excelente**" o bien con la palabra "**pobre**" no es casualidad, sino que la frase misma tiene una carga de valor, una polaridad.

Algoritmo de Turney - polaridad de las frases



Pointwise mutual information:

El **PMI** de un par de valores **x** e **y** pertenecientes a dos variables aleatorias discretas: **X** e **Y** respectivamente, cuantifica la discrepancia entre la probabilidad de su coincidencia dada su distribución conjunta y su distribución individual y asumiendo su independencia.

Matemáticamente:

$$\text{PMI}(X, Y) = \text{Log}_2 \frac{P(x, y)}{P(x)P(y)}$$

Algoritmo de Turney - polaridad de las frases



En otras palabras: cuanto más posible es que el evento **X** aparezca vinculado al evento **Y** a que aparezcan ambos de forma independiente entre sí.

Pointwise mutual information entre dos palabras:

$$\text{PMI}(\textit{palabra1}, \textit{palabra2}) = \text{Log}_2 \frac{P(\textit{palabra1}, \textit{palabra2})}{P(\textit{palabra1})P(\textit{palabra2})}$$

Algoritmo de Turney - polaridad de las frases

Turney utilizó el buscador **Altavista.com** para obtener estos valores, pero la forma que utilizemos para contar estos resultados dependerá de nuestro conjunto de datos de entrenamiento.

$$P(\text{palabra}) = \frac{\text{cantidad de resultados para "palabra"}}{\text{cantidad total}}$$

$$P(\text{palabra}) = \frac{\# \text{palabra}}{N}$$

Algoritmo de Turney - polaridad de las frases



$$P(\text{palabra1}, \text{palabra2}) =$$

$$= \frac{\text{cantidad de resultados para "palabra1 NEAR palabra2"}}{(\text{cantidad total})^2}$$

$$P(\text{palabra1}, \text{palabra2}) = \frac{\#(\text{palabra1 NEAR palabra2})}{N^2}$$

NEAR significa "cerca de", indica que **palabra1** apareció a no más de **N** palabras de distancia de **palabra2** en un texto dado.

Algoritmo de Turney - polaridad de las frases

$$\text{PMI}(\text{palabra1}, \text{palabra2}) = \text{Log}_2 \frac{\frac{\#(\text{palabra1 NEAR palabra2})}{N^2}}{\frac{\#(\text{palabra1}) \#(\text{palabra2})}{N * N}}$$

Los denominadores se cancelan

Algoritmo de Turney - polaridad de las frases



$$\text{PMI}(\text{palabra1}, \text{palabra2}) = \log_2 \left(\frac{\#(\text{palabra1 NEAR palabra2})}{\#(\text{palabra1}) \#(\text{palabra2})} \right)$$

Calculando la polaridad de una frase:

$$\text{Polaridad}(\text{frase}) = \text{PMI}(\text{frase}, \text{"excelente"}) - \text{PMI}(\text{frase}, \text{"pobre"})$$

En resumen:

$$\text{Polaridad}(\text{frase}) = \log_2 \frac{\#(\text{frase NEAR excelente}) \#(\text{pobre})}{\#(\text{frase NEAR pobre}) \#(\text{excelente})}$$

Algoritmo de Turney - polaridad de las frases



Ejemplos:

Fraser	Polaridad
online service	2.8
online experience	2.3
low fees	0.33
inconveniently located	-1.5
Average	0.32

Algoritmo de Turney - puntuar las críticas



Armar una lista de frases y asociarles el valor obtenido con los cálculos anteriores para luego descomponer las críticas en sus frases y realizar el promedio.

Algoritmo de Turney - conclusiones



Sobre 410 opiniones de *Epinions*

- Exactitud promedio: 74%
 - críticas cinematográficas: 66%
 - críticas sobre bancos y automóviles: 80% y el 84%
 - críticas sobre viajes: intermedio.

Aspectos

Aspectos



¿Cómo detectar más de un sentimiento en la misma frase?

¿Qué sucede cuando tenemos frases como la siguiente?

¡La comida era excelente pero el servicio pésimo!

Aspectos - Método de Minqing Hu y Bing Liu



- Frecuencia
- Reglas

Frecuencia: Buscaron todas las **frases frecuentes** en las críticas de un lugar dado.

Por ejemplo para un mismo restaurante encontraron que se repetía muchas veces la frase: "**tacos de pescado**". A este tipo de frases las llamaron **aspectos**, **atributos** o bien **objetos de sentimiento** ya que representan el objeto al cual se está criticando.

Aspectos - Método de Minqing Hu y Bing Liu



- Frecuencia
- Reglas

Reglas: Filtraron todas esas frases frecuentes con algunas reglas como: "ocurre después de una palabra que indica sentimientos".

Ejemplo: "geniales **tacos de pescado!**" => indica que "tacos de pescado" es muy probablemente un aspecto.

Aspectos - Método de Minqing Hu y Bing Liu



Encontraron, para críticas sobre los siguientes sitios, los siguientes aspectos:

Sitio	Aspectos
Casino	casino, buffet, piscina, resort, camas
Peluquería infantil	corte, trabajo, experiencia, niños
Restaurante Griego	comida, vino, servicio, aperitivos, cordero
Grandes tiendas	selecciones, departamentos, ventas, negocios, ropa

Aspectos - Método de Minqing Hu y Bing Liu



Consideraciones finales

- El aspecto puede no ser mencionado en una sentencia
- Para hoteles/restaurantes los aspectos son generalmente conocidos y fáciles de identificar
- Es posible utilizar una clasificación supervisada:
 - Para pequeños corpus, se puede utilizar una clasificación manual de aspectos: comida, decoración, servicio, precio, **nada**
- Y luego entrenar un clasificador:
 - dada una sentencia, tiene alguno de estos aspectos: "comida, decoración, servicio, precio, nada"

Aspectos - Método de Minqing Hu y Bing Liu



Consideraciones finales

- Si la cantidad de críticas no está balanceada (entre positivas y negativas o entre los rangos elegidos)
 - no se puede utilizar el estimador: precisión
 - hay que utilizar F-score visto anteriormente
- Si el desbalanceo es muy pronunciado se puede degradar severamente el rendimiento del clasificador
 - Dos soluciones comunes:
 - tomar un muestreo parejo
 - penalizar más severamente al clasificador por un error al categorizar la clase más rara.

Preguntas...