

Caso de Estudio: Clasificación de Texto

Clasificación de Texto

Ejemplos:

- Detección de spam: ham or spam?
- **Análisis de Sentimiento:** pos or neg?
 - También neu, none...
 - También valores de intensidad
- Atribución de Autoría: who wrote this?
- Clasificación de artículos: politics, sports or culture?

Clasificación de Texto

Dan Jurafsky



Is this spam?

Subject: Important notice!

From: Stanford University <newsforum@stanford.edu>

Date: October 28, 2011 12:34:16 PM PDT

To: undisclosed-recipients:;

Greats News!

You can now access the latest news by using the link below to login to Stanford University News Forum.

<http://www.123contactform.com/contact-form-StanfordNew1-236335.html>

Click on the above link to login for more information about this new exciting forum. You can also copy the above link to your browser bar and login for more information about the new services.

© Stanford University. All Rights Reserved.

Clasificación de Texto

Dan Jurafsky



Who wrote which Federalist papers?

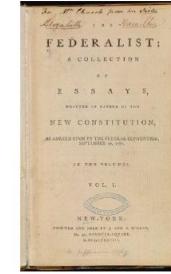
- 1787-8: anonymous essays try to convince New York to ratify U.S Constitution: Jay, Madison, Hamilton.
- Authorship of 12 of the letters in dispute
- 1963: solved by Mosteller and Wallace using Bayesian methods



James Madison



Alexander Hamilton



Clasificación de Texto

Dan Jurafsky



Male or female author?

1. By 1925 present-day Vietnam was divided into three parts under French colonial rule. The southern region embracing Saigon and the Mekong delta was the colony of Cochinchina; the central area with its imperial capital at Hue was the protectorate of Annam...
2. Clara never failed to be astonished by the extraordinary felicity of her own name. She found it hard to trust herself to the mercy of fate, which had managed over the years to convert her greatest shame into one of her greatest assets...

S. Argamon, M. Koppel, J. Fine, A. R. Shimoni, 2003. "Gender, Genre, and Writing Style in Formal Written Texts," *Text*, volume 23, number 3, pp. 321–346

Clasificación de Texto

Dan Jurafsky



Positive or negative movie review?



- unbelievably disappointing



- Full of zany characters and richly applied satire, and some great plot twists



- this is the greatest screwball comedy ever filmed



- It was pathetic. The worst part about it was the boxing scenes.

Clasificación de Texto

Dan Jurafsky



What is the subject of this article?

MEDLINE Article



6

MeSH Subject Category Hierarchy

- Antagonists and Inhibitors
- Blood Supply
- Chemistry
- Drug Therapy
- Embryology
- Epidemiology
- ...

Clasificación de Texto

Dan Jurafsky



Text Classification

- Assigning subject categories, topics, or genres
- Spam detection
- Authorship identification
- Age/gender identification
- Language Identification
- Sentiment analysis
- ...

Clasificación de Texto

Dan Jurafsky



Text Classification: definition

- *Input:*
 - a document d
 - a fixed set of classes $C = \{c_1, c_2, \dots, c_J\}$
- *Output:* a predicted class $c \in C$

Clasificación de Texto

Dan Jurafsky



Classification Methods: Hand-coded rules

- Rules based on combinations of words or other features
 - spam: black-list-address OR (“dollars” AND “have been selected”)
- Accuracy can be high
 - If rules carefully refined by expert
- But building and maintaining these rules is expensive

Clasificación de Texto

Dan Jurafsky



Classification Methods: Supervised Machine Learning

- *Input:*
 - a document d
 - a fixed set of classes $C = \{c_1, c_2, \dots, c_j\}$
 - A training set of m hand-labeled documents $(d_1, c_1), \dots, (d_m, c_m)$
- *Output:*
 - a learned classifier $y: d \rightarrow c$

Clasificación de Texto

Dan Jurafsky



Classification Methods: Supervised Machine Learning

- Any kind of classifier
 - Naïve Bayes
 - Logistic regression
 - Support-vector machines
 - k-Nearest Neighbors
 - ...

Evaluación

Dan Jurafsky



More Than Two Classes: Sets of binary classifiers

- Dealing with **any-of** or **multivalue** classification
 - A document can belong to 0, 1, or >1 classes.
- For each class $c \in C$
 - Build a classifier γ_c to distinguish c from all other classes $c' \in C$
- Given test doc d ,
 - Evaluate it for membership in each class using each γ_c
 - d belongs to **any** class for which γ_c returns true

Evaluación

Dan Jurafsky



More Than Two Classes: Sets of binary classifiers

- One-of or multinomial classification
 - Classes are mutually exclusive: each document in exactly one class
- For each class $c \in C$
 - Build a classifier γ_c to distinguish c from all other classes $c' \in C$
- Given test doc d ,
 - Evaluate it for membership in each class using each γ_c
 - d belongs to the one class with maximum score

Evaluación



Dan Jurafsky

Evaluation: Classic Reuters-21578 Data Set

- Most (over)used data set, 21,578 docs (each 90 types, 200 tokens)
- 9603 training, 3299 test articles (ModApte/Lewis split)
- 118 categories
 - An article can be in more than one category
 - Learn 118 binary category distinctions
- Average document (with at least one category) has 1.24 classes
- Only about 10 out of 118 categories are large
 - Earn (2877, 1087)
 - Acquisitions (1650, 179)
 - Money-fx (538, 179)
 - Grain (433, 149)
 - Crude (389, 189)
 - Trade (369, 119)
 - Interest (347, 131)
 - Ship (197, 89)
 - Wheat (212, 71)
 - Corn (182, 56)

Common categories
(#train, #test)

56

Evaluación

Dan Jurafsky



Reuters Text Categorization data set (Reuters-21578) document

<REUTERS TOPICS="YES" LEWISSPLIT="TRAIN" CGISPLIT="TRAINING-SET" OLDID="12981"
NEWID="798">

<DATE> 2-MAR-1987 16:51:43.42</DATE>

<TOPICS><D>livestock</D><D>hog</D></TOPICS>

<TITLE>AMERICAN PORK CONGRESS KICKS OFF TOMORROW</TITLE>

<DATELINE> CHICAGO, March 2 - </DATELINE><BODY>The American Pork Congress kicks off tomorrow, March 3, in Indianapolis with 160 of the nations pork producers from 44 member states determining industry positions on a number of issues, according to the National Pork Producers Council, NPPC.

Delegates to the three day Congress will be considering 26 resolutions concerning various issues, including the future direction of farm policy and the tax law as it applies to the agriculture sector. The delegates will also debate whether to endorse concepts of a national PRV (pseudorabies virus) control and eradication program, the NPPC said.

A large trade show, in conjunction with the congress, will feature the latest in technology in all areas of the industry, the NPPC added. Reuter

⁵⁷
.</BODY></TEXT></REUTERS>

Evaluación

Dan Jurafsky



Confusion matrix c

- For each pair of classes $\langle c_1, c_2 \rangle$ how many documents from c_1 were incorrectly assigned to c_2 ?
 - $c_{3,2}$: 90 wheat documents incorrectly assigned to poultry

Docs in test set	Assigned UK	Assigned poultry	Assigned wheat	Assigned coffee	Assigned interest	Assigned trade
True UK	95	1	13	0	1	0
True poultry	0	1	0	0	0	0
True wheat	10	90	0	1	0	0
True coffee	0	0	0	34	3	7
True interest	-	1	2	13	26	5
True trade	0	0	2	14	5	10

Evaluación

Dan Jurafsky



Per class evaluation measures

Recall:

Fraction of docs in class i classified correctly:

$$\frac{c_{ii}}{\sum_j c_{ij}}$$

Precision:

Fraction of docs assigned class i that are actually about class i :

$$\frac{c_{ii}}{\sum_j c_{ji}}$$

Accuracy: (1 - error rate)

Fraction of docs classified correctly:

$$\frac{\sum_i c_{ii}}{\sum_j \sum_i c_{ij}}$$

Evaluación

Dan Jurafsky



Micro- vs. Macro-Averaging

- If we have more than one class, how do we combine multiple performance measures into one quantity?
- **Macroaveraging:** Compute performance for each class, then average.
- **Microaveraging:** Collect decisions for all classes, compute contingency table, evaluate.

Evaluación

Dan Jurafsky



Micro- vs. Macro-Averaging: Example

Class 1

	Truth: yes	Truth: no
Classifier: yes	10	10
Classifier: no	10	970

Class 2

	Truth: yes	Truth: no
Classifier: yes	90	10
Classifier: no	10	890

Micro Ave. Table

	Truth: yes	Truth: no
Classifier: yes	100	20
Classifier: no	20	1860

- Macroaveraged precision: $(0.5 + 0.9)/2 = 0.7$
- Microaveraged precision: $100/120 = .83$
- Microaveraged score is dominated by score on common classes

Evaluación

Dan Jurafsky



Micro- vs. Macro-Averaging: Example

Class 1

	Truth: yes	Truth: no
Classifier: yes	10	10
Classifier: no	10	970

Class 2

	Truth: yes	Truth: no
Classifier: yes	90	10
Classifier: no	10	890

Micro Ave. Table

	Truth: yes	Truth: no
Classifier: yes	100	20
Classifier: no	20	1860

- Macroaveraged precision: $(0.5 + 0.9)/2 = 0.7$
- Microaveraged precision: $100/120 = .83$
- Microaveraged score is dominated by score on common classes

61

$$\text{Precision} = \text{tp} / (\text{tp} + \text{fp}) \quad \text{Recall} = \text{tp} / (\text{tp} + \text{fn})$$

Cuestiones Prácticas

Dan Jurafsky



The Real World

- Gee, I'm building a text classifier for real, now!
- What should I do?

Cuestiones Prácticas

Dan Jurafsky



No training data? Manually written rules

If (wheat or grain) and not (whole or bread) then
Categorize as grain

- Need careful crafting
 - Human tuning on development data
 - Time-consuming: 2 days per class

Cuestiones Prácticas

Dan Jurafsky



Very little data?

- Use Naïve Bayes
 - Naïve Bayes is a “high-bias” algorithm (Ng and Jordan 2002 NIPS)
- Get more labeled data
 - Find clever ways to get humans to label data for you
- Try semi-supervised training methods:
 - Bootstrapping, EM over unlabeled documents, ...

Cuestiones Prácticas

Dan Jurafsky



A reasonable amount of data?

- Perfect for all the clever classifiers
 - SVM
 - Regularized Logistic Regression
- You can even use user-interpretable decision trees
 - Users like to hack
 - Management likes quick fixes

Cuestiones Prácticas

Dan Jurafsky



A huge amount of data?

- Can achieve high accuracy!
- At a cost:
 - SVMs (train time) or kNN (test time) can be too slow
 - Regularized logistic regression can be somewhat better
 - So Naïve Bayes can come back into its own again!

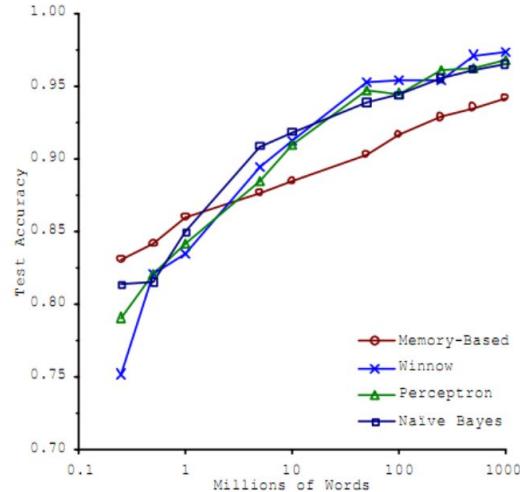
Cuestiones Prácticas

Dan Jurafsky



Accuracy as a function of data size

- With enough data
 - Classifier may not matter



Brill and Banko on spelling correction

Cuestiones Prácticas

Dan Jurafsky



Real-world systems generally combine:

- Automatic classification
- Manual review of uncertain/difficult/"new" cases

Cuestiones Prácticas

Dan Jurafsky



How to tweak performance

- Domain-specific features and weights: *very* important in real performance
- Sometimes need to collapse terms:
 - Part numbers, chemical formulas, ...
 - But stemming generally doesn't help
- Upweighting: Counting a word as if it occurred twice:
 - title words ([Cohen & Singer 1996](#))
 - first sentence of each paragraph ([Murata, 1999](#))
 - In sentences that contain title words ([Ko et al, 2002](#))

Análisis de Sentimiento (Sentiment Analysis)

<https://news.slashdot.org/story/18/05/18/1455211/data-science-is-americas-hottest-job>

Data Science is America's Hottest Job (bloomberg.com)

 Posted by msmash on Friday May 18, 2018 @12:00PM from the feeds-and-speeds dept.   

79

Anonymous readers share a report:

It turns out that even in the wake of Facebook's privacy scandal and other big-data blunders, finding people who can turn social-media clicks and user-posted photos into monetizable binary code is among the biggest challenges facing U.S. industry. People with data science bona fides are among the most sought-after professionals in business, with some data science Ph.Ds commanding as much as \$300,000 or more from consulting firms.

Job postings for data scientists rose 75 percent from January 2015 to January 2018 at Indeed.com, while job searches for data scientist roles rose 65 percent. A growing specialty is "**sentiment analysis**," or finding a way to quantify how many tweets are trashing your company or praising it. A typical data scientist job pays about \$119,000 at the midpoint of salaries and rises to \$168,000 at the 95th percentile, according to staffing agency Robert Half Technology.

Análisis de Sentimiento

Dan Jurafsky



Positive or negative movie review?



- unbelievably disappointing



- Full of zany characters and richly applied satire, and some great plot twists



- this is the greatest screwball comedy ever filmed



- It was pathetic. The worst part about it was the boxing scenes.

Análisis de Sentimiento

Dan Jurafsky



Google Product Search



HP Officejet 6500A Plus e-All-in-One Color Ink-jet - Fax / copier / printer / scanner
\$89 online, \$100 nearby ★★★★☆ 377 reviews
September 2010 - Printer - HP - Inkjet - Office - Copier - Color - Scanner - Fax - 250 sh

Reviews

Summary - Based on 377 reviews

1 star 2 3 4 stars 5 stars

What people are saying

ease of use		"This was very easy to setup to four computers."
value		"Appreciate good quality at a fair price."
setup		"Overall pretty easy setup."
customer service		"I DO like honest tech support people."
size		"Pretty Paper weight."
mode		"Photos were fair on the high quality mode."
colors		"Full color prints came out with great quality."

3

Análisis de Sentimiento

Dan Jurafsky



Bing Shopping

HP Officejet 6500A E710N Multifunction Printer

[Product summary](#) [Find best price](#) **Customer reviews** [Specifications](#) [Related items](#)



\$121.53 - \$242.39 (14 stores)

Compare

Average rating	(Count)
★★★★★	(55)
★★★★★	(54)
★★★★★	(10)
★★★★★	(6)
★★★★★	(23)
★★★★★	(0)

Most mentioned

Topic	(Count)
Performance	(57)
Ease of Use	(43)
Print Speed	(39)
Connectivity	(31)

Show reviews by source

- Best Buy (140)
- CNET (5)
- Amazon.com (3)

[More ▾](#)

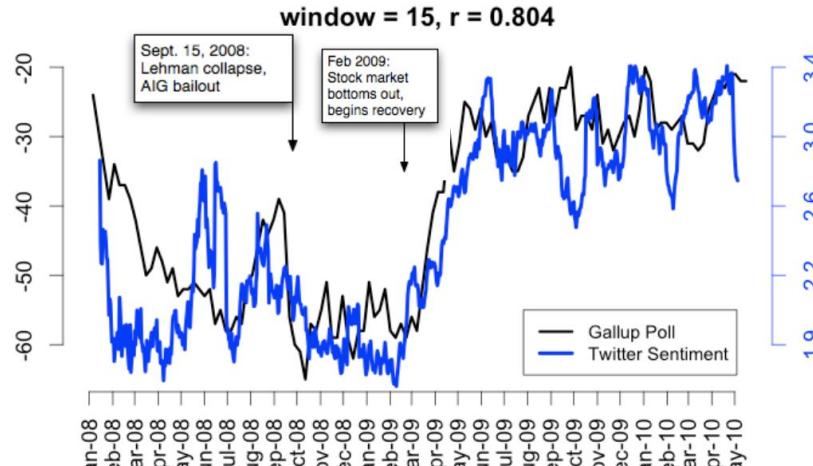
Análisis de Sentimiento

Dan Jurafsky



Twitter sentiment versus Gallup Poll of Consumer Confidence

Brendan O'Connor, Ramnath Balasubramanyan, Bryan R. Routledge, and Noah A. Smith. 2010.
From Tweets to Polls: Linking Text Sentiment to Public Opinion Time Series. In ICWSM-2010



Análisis de Sentimiento

Dan Jurafsky



Target Sentiment on Twitter

Type in a word and we'll highlight the good and the bad

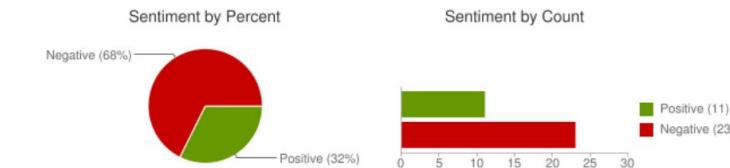
"united airlines"

Search

[Save this search](#)

- [Twitter Sentiment App](#)
- Alec Go, Richa Bhayani, Lei Huang. 2009. Twitter Sentiment Classification using Distant Supervision

Sentiment analysis for "united airlines"



jjacobsen: OMG... Could @United airlines have worse customer service? W8g now 15 minut
Posted 2 hours ago

12345clumsy6789: I hate United Airlines Ceiling!!! Fukn impossible to get my conduit in this d
Posted 2 hours ago

EMLandPRGbelgiu: EML/PRG fly with Q8 united airlines and 24seven to an exotic destination
Posted 2 hours ago

CountAdam: FANTASTIC customer service from United Airlines at XNA today. Is tweet more
Posted 4 hours ago

Análisis de Sentimiento

Dan Jurafsky



Sentiment analysis has many other names

- Opinion extraction
- Opinion mining
- Sentiment mining
- Subjectivity analysis

Análisis de Sentimiento

Dan Jurafsky



Why sentiment analysis?

- *Movie*: is this review positive or negative?
- *Products*: what do people think about the new iPhone?
- *Public sentiment*: how is consumer confidence? Is despair increasing?
- *Politics*: what do people think about this candidate or issue?
- *Prediction*: predict election outcomes or market trends from sentiment

Técnicas



Dan Jurafsky



Sentiment Analysis

- Simplest task:
 - Is the attitude of this text positive or negative?
- More complex:
 - Rank the attitude of this text from 1 to 5
- Advanced:
 - Detect the target, source, or complex attitude types

Técnicas

Dan Jurafsky



Sentiment Classification in Movie Reviews

Bo Pang, Lillian Lee, and Shivakumar Vaithyanathan. 2002. Thumbs up? Sentiment Classification using Machine Learning Techniques. EMNLP-2002, 79–86.

Bo Pang and Lillian Lee. 2004. A Sentimental Education: Sentiment Analysis Using Subjectivity Summarization Based on Minimum Cuts. ACL, 271-278

- Polarity detection:
 - Is an IMDB movie review positive or negative?
- Data: *Polarity Data 2.0:*
 - <http://www.cs.cornell.edu/people/pabo/movie-review-data>

Técnicas

Dan Jurafsky



IMDB data in the Pang and Lee database



when _star wars_ came out some twenty years ago , the image of traveling throughout the stars has become a commonplace image . [...]

when han solo goes light speed , the stars change to bright lines , going towards the viewer in lines that converge at an invisible point .
cool .

october sky offers a much simpler image—that of a single white dot , traveling horizontally across the night sky . [. . .]



“ snake eyes ” is the most aggravating kind of movie : the kind that shows so much potential then becomes unbelievably disappointing .

it’s not just because this is a brian depalma film , and since he’s a great director and one who’s films are always greeted with at least some fanfare .

and it’s not even because this was a film starring nicolas cage and since he gives a brauvara performance , this film is hardly worth his talents .

Técnicas

Dan Jurafsky



Baseline Algorithm (adapted from Pang and Lee)

- Tokenization
- Feature Extraction
- Classification using different classifiers
 - Naïve Bayes
 - MaxEnt
 - SVM

Técnicas

Dan Jurafsky



Sentiment Tokenization Issues

- Deal with HTML and XML markup
- Twitter mark-up (names, hash tags)
- Capitalization (preserve for words in all caps)
- Phone numbers, dates
- Emoticons
- Useful code:
 - [Christopher Potts sentiment tokenizer](#)
 - [Brendan O'Connor twitter tokenizer](#)

Potts emoticons

```
[<>]?
[::=8]
[‐o*']?
()]\(\dDpP/\:\}\{@\|\ \\
| 
()\]\(\dDpP/\:\}\{@\|\ \\
[‐o*']?
[::=8]
[<>]? # optional hat/brow
# eyes
# optional nose
# mouth
#### reverse orientation
# mouth
# optional nose
# eyes
# optional hat/brow
```

Técnicas

Dan Jurafsky



Extracting Features for Sentiment Classification

- How to handle negation
 - I **didn't** like this movie
 - vs
 - I really like this movie
- Which words to use?
 - Only adjectives
 - All words
 - All words turns out to work better, at least on this data

Técnicas

Dan Jurafsky



Negation

Das, Sanjiv and Mike Chen. 2001. Yahoo! for Amazon: Extracting market sentiment from stock message boards. In Proceedings of the Asia Pacific Finance Association Annual Conference (APFA).
Bo Pang, Lillian Lee, and Shivakumar Vaithyanathan. 2002. Thumbs up? Sentiment Classification using Machine Learning Techniques. EMNLP-2002, 79–86.

Add NOT_ to every word between negation and following punctuation:

didn't like this movie , but I



didn't NOT_like NOT_this NOT_movie but I

Técnicas

Dan Jurafsky



Reminder: Naïve Bayes

$$c_{NB} = \operatorname{argmax}_{c_j \in C} P(c_j) \prod_{i \in positions} P(w_i | c_j)$$

$$\hat{P}(w | c) = \frac{\text{count}(w, c) + 1}{\text{count}(c) + |V|}$$

Técnicas

Dan Jurafsky



Binarized (Boolean feature) Multinomial Naïve Bayes

- Intuition:
 - For sentiment (and probably for other text classification domains)
 - Word occurrence may matter more than word frequency
 - The occurrence of the word *fantastic* tells us a lot
 - The fact that it occurs 5 times may not tell us much more.
 - Boolean Multinomial Naïve Bayes
 - Clips all the word counts in each document at 1

Técnicas



Dan Jurafsky

Binarized (Boolean feature) Multinomial Naïve Bayes

- B. Pang, L. Lee, and S. Vaithyanathan. 2002. Thumbs up? Sentiment Classification using Machine Learning Techniques. EMNLP-2002, 79–86.
V. Metsis, I. Androutsopoulos, G. Palouras. 2006. Spam Filtering with Naive Bayes – Which Naive Bayes? CEAS 2006 - Third Conference on Email and Anti-Spam.
K.-M. Schneider. 2004. On word frequency information and negative evidence in Naive Bayes text classification. ICANLP, 474-485.
JD Rennie, L Shih, J Teevan. 2003. Tackling the poor assumptions of naive bayes text classifiers. ICML 2003

- Binary seems to work better than full word counts
 - This is **not** the same as Multivariate Bernoulli Naïve Bayes
 - MBNB doesn't work well for sentiment or other text tasks
 - Other possibility: $\log(\text{freq}(w))$

Técnicas

Dan Jurafsky



Other issues in Classification

- MaxEnt and SVM tend to do better than Naïve Bayes

Problemas

Dan Jurafsky



Problems: What makes reviews hard to classify?

- Subtlety:
 - Perfume review in *Perfumes: the Guide*:
 - “If you are reading this because it is your darling fragrance, please wear it at home exclusively, and tape the windows shut.”
 - Dorothy Parker on Katherine Hepburn
 - “She runs the gamut of emotions from A to B”

Problemas

- Ellipsis e ironía:

@LovNaty Tu vida ha parido a un grandisimo hijo de la gran p... , un maravilloso hombre!!. (NEG)

- Negación, doble negación y “pero”:

No es que ahora no sea feliz, pero antes lo era más (NEG)

Problemas

Dan Jurafsky



Thwarted Expectations and Ordering Effects

- “This film should be brilliant. It sounds like a great plot, the actors are first grade, and the supporting cast is good as well, and Stallone is attempting to deliver a good performance. However, it can’t hold up.”
- Well as usual Keanu Reeves is nothing special, but surprisingly, the very talented Laurence Fishbourne is not so good either, I was surprised.

Vectorización

Dan Jurafsky



The bag of words representation

Y()=C

I love this movie! It's sweet, but with satirical humor. The dialogue is great and the adventure scenes are fun... It manages to be whimsical and romantic while laughing at the conventions of the fairy tale genre. I would recommend it to just about anyone. I've seen it several times, and I'm always happy to see it again whenever I have a friend who hasn't seen it yet.



Vectorización

Dan Jurafsky



The bag of words representation

I **love** this movie! It's **sweet**, but with **satirical** humor. The dialogue is **great** and the adventure scenes are **fun**... It manages to be **whimsical** and **romantic** while **laughing** at the conventions of the fairy tale genre. I would **recommend** it to just about anyone. I've seen it **several** times, and I'm always **happy** to see it **again** whenever I have a friend who hasn't seen it yet.

Y()=C



Vectorización

Dan Jurafsky



The bag of words representation

Y()=C

great	2
love	2
recommend	1
laugh	1
happy	1
...	...



Aumentación de Datos

- Traducción bidireccional desde/hacia varios idiomas.

Original: La verdad es que tiene buena pinta. Investigaré, gracias

Traducciones:

1. La verdad es que parece bueno. Voy a investigar, gracias
2. La verdad es que se ve bien. Voy a investigar, gracias
3. El hecho es que se ven bien. Lo comprobaré, gracias

Word Embeddings

Deep Learning for Natural Language Processing (Stanford)

- Bag-of-words: No están representadas palabras no vistas (OOV words).
- One-hot encoding: todas las palabras tienen similaridad cero entre sí.
- Embeddings:
 - Llevar todas las palabras a una representación vectorial compacta.
 - Aprender los vectores de manera no supervisada con muchos datos.
 - Usar los vectores pre-entrenados como input features en problemas supervisados.