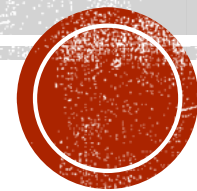
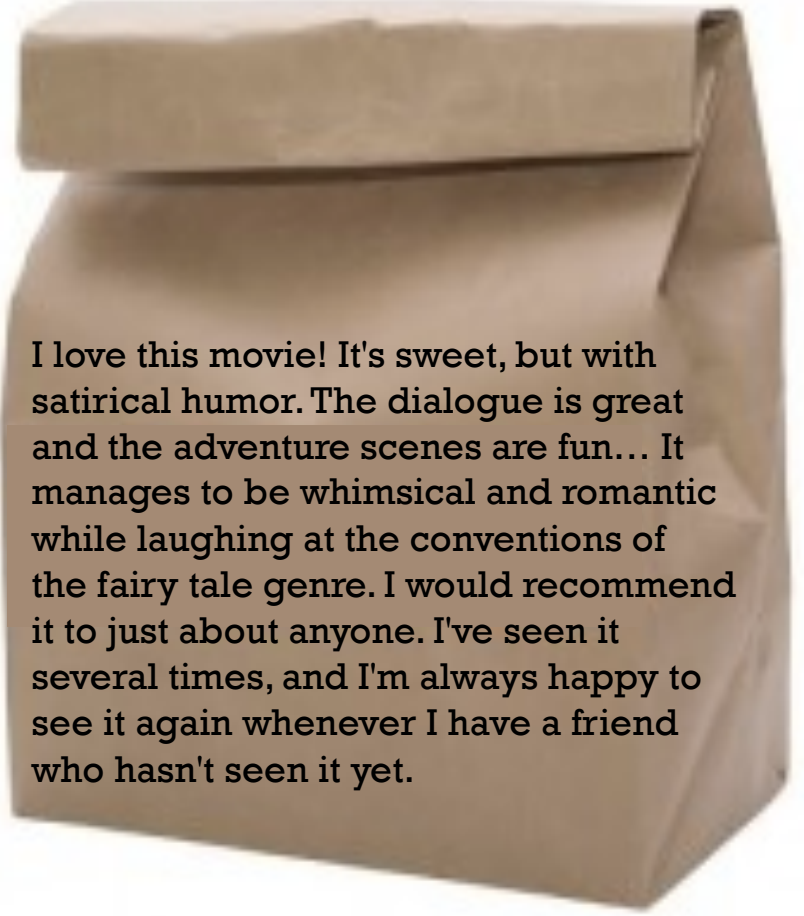


IMPROVING TEXT CLASSIFICATION



BAG OF WORDS

- **Bag of Words:** Assume position doesn't matter
- **Conditional Independence:** Assume the feature probabilities $P(x_i | c_j)$ are independent given the class 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.



ADVANTAGES

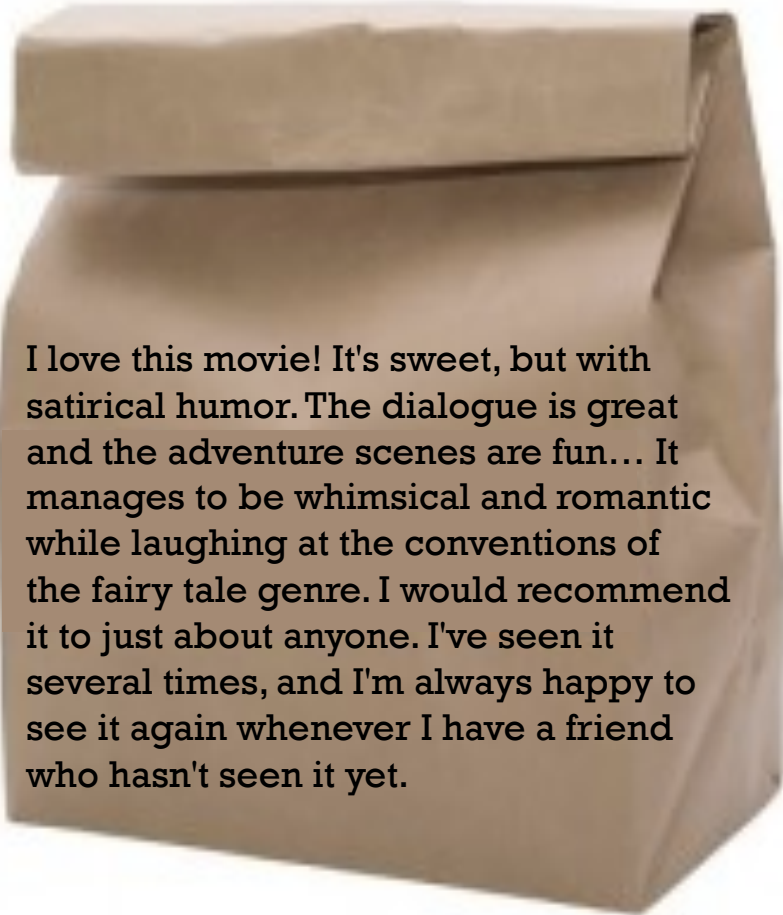
- Very Fast, low storage requirements
- Robust to Irrelevant Features
 - Irrelevant Features cancel each other without affecting results
- Very good in domains with many equally important features
 - Decision Trees suffer from *fragmentation* in such cases – especially if little data
- Optimal if the independence assumptions hold:
 - If assumed independence is correct, then it is the Bayes Optimal Classifier for problem
- A good dependable baseline for text classification



POSSIBLE IMPROVEMENTS

- Stop words
- Stemming
- TF-IDF
- Bigrams
- More





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.

$$\gamma \begin{pmatrix} \text{it} \\ \text{I} \\ \text{the} \\ \text{to} \\ \text{and} \\ \text{seen} \\ \text{yet} \\ \text{love} \\ \text{movie} \\ \text{sweet} \\ \dots \end{pmatrix} = c$$



STOP WORDS

- Commonly used words
- May not contribute to how text is classified
- Possible wasted time and space
- Making stop word list
 - Collect most frequent terms
 - Keep those with relevant semantic content
- Pre-process all documents to remove stop words

a	but
about	by
above	can't
after	cannot
again	could
against	couldn't
all	did
am	didn't
an	do
and	does
any	doesn't
are	doing
aren't	don't
as	down
at	during
be	each
because	few
been	for
before	from
being	further
below	had
between	hadn't
both	has



car
cars

movie
movies

laugh
laughs
laughing
laughed
laughter

replace
replaces
replaced
replacement

swim
swims
swam
swimming
swimmer

hope
hopes
hoped
hoping
hopeful
hopefully
hopefulness
hopeless
hopelessly
hopelessness

pony
ponies
ponied

quick
quicker
quickest
quickly

hero
heroes
heroic
heroically



STEMMING

- Many forms of a word may be used
- Stemming: Chop off the end of words
- Crude heuristic process
- Porter's algorithm
 - 5 step process
 - Changes ending based on some conditions

- Example rules:

- sses → ss
- ies → i
- ss → ss
- s →

possesses → possess
movies → movi
possess → possess
cars → car

- (m>1) ement →

replacement → replac
cement → cement



TERM FREQUENCY

- d_1 : This is the best movie of the year.
- d_2 : This movie features a great story with a great cast.
- “this”
 - Raw count:
 - $f(\text{“this”}, d_1) = 1$
 - $f(\text{“this”}, d_2) = 1$
 - Frequency:
 - $tf(\text{“this”}, d_1) = 1/8$
 - $tf(\text{“this”}, d_2) = 1/10$

Term	Term count	Term	Term count
this	1	this	1
is	1	a	2
the	2	features	1
best	1	great	2
movie	1	movie	1
of	1	with	1
year	1	story	1
		cast	1



INVERSE DOCUMENT FREQUENCY

- A measure of how much information a word provides
 - If a word is in most (or all) documents, it doesn't help tell us how to classify the document
- Ratio of how many documents in corpus and how many documents with a given word

$$\text{idf}(t, D) = \log \frac{N}{|\{d \in D : t \in d\}|}$$

- “this”
 - $\text{idf}(\text{“this”}, D) = \log 2/2 = 0$
- “best”
 - $\text{idf}(\text{“best”}, D) = \log 2/1 = 0.301$

Term	Term count	Term	Term count
this	1	this	1
is	1	a	2
the	2	features	1
best	1	great	2
movie	1	movie	1
of	1	with	1
year	1	story	1
		cast	1

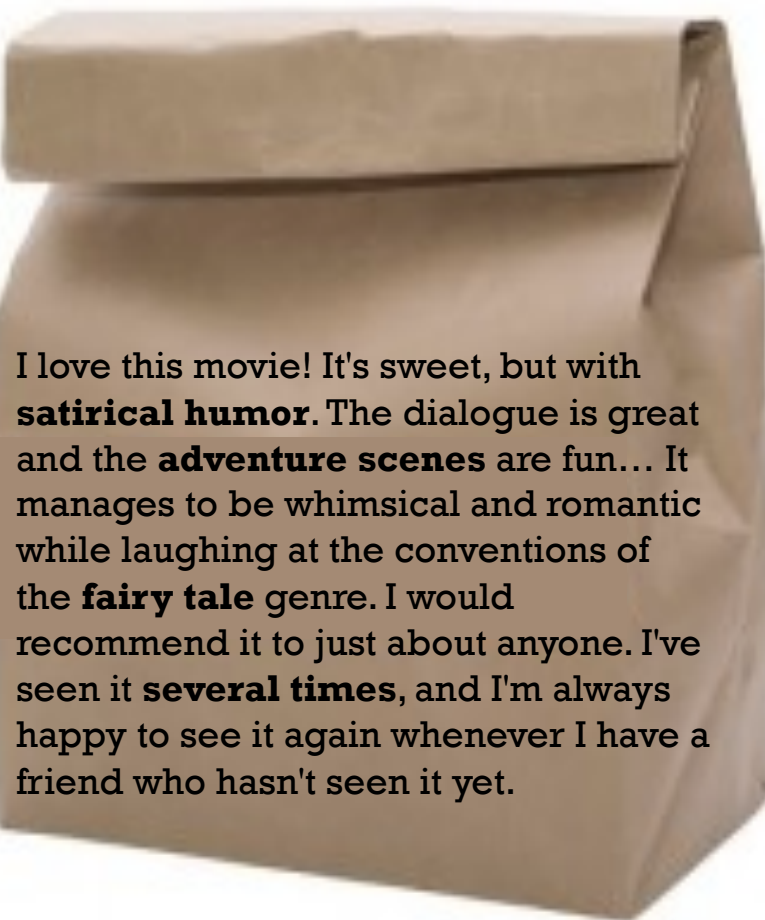


TF-IDF

- Term frequency – inverse document frequency
- $\text{tfidf}(t,d) = \text{tf}(t,d) * \text{idf}(t,D)$
- “this”
 - $\text{tfidf}(\text{“this”},d_1) = \text{tf}(\text{“this”},d_1) * \text{idf}(\text{“this”},D)$
 $= 0.125 * 0 = 0$
 - $\text{tfidf}(\text{“this”},d_2) = \text{tf}(\text{“this”},d_2) * \text{idf}(\text{“this”},D)$
 $= 0.1 * 0 = 0$
- “best”
 - $\text{tfidf}(\text{“best”},d_1) = \text{tf}(\text{“best”},d_1) * \text{idf}(\text{“best”},D)$
 $= 0.125 * 0.301 = 0.038$
 - $\text{tfidf}(\text{“best”},d_2) = \text{tf}(\text{“best”},d_2) * \text{idf}(\text{“best”},D)$
 $= 0 * 0.301 = 0$

Term	Term count	Term	Term count
this	1	this	1
is	1	a	2
the	2	features	1
best	1	great	2
movie	1	movie	1
of	1	with	1
year	1	story	1
		cast	1





BIGRAMS

this movie	6	
I have	3	
of the	7	
fairy tale	2	
and the	8	
satirical humor	1	= c
adventure scenes	1	
several times	1	
seen it	3	
a friend	4	
...		



MORE

- Capitalization
- Punctuation
- Contractions
- Misspellings
- Jargon

