# Text Classification and Naive Bayes

Sentiment and Binary Naive Bayes

# Let's do a worked sentiment example!

| | Cat | Documents |
|---|---|---|
| Training | - | just plain boring |
| | - | entirely predictable and lacks energy |
| | - | no surprises and very few laughs |
| | + | very powerful |
| | + | the most fun film of the summer |
| Test | ? | predictable with no fun |

# A worked sentiment example

| | Cat | Documents |
|---|---|---|
| Training | - | just plain boring |
| | - | entirely predictable and lacks energy |
| | - | no surprises and very few laughs |
| | + | very powerful |
| | + | the most fun film of the summer |
| Test | ? | predictable ~~with~~ no fun |

Prior from training:

$P(-) = 3/5$
$P(+) = 2/5$

Drop "with"

Likelihoods from training:

$$P(\text{"predictable"}|-) = \frac{1+1}{14+20}$$

$$P(\text{"predictable"}|+) = \frac{0+1}{9+20}$$

$$P(\text{"no"}|-) = \frac{1+1}{14+20}$$

$$P(\text{"no"}|+) = \frac{0+1}{9+20}$$

$$P(\text{"fun"}|-) = \frac{0+1}{14+20}$$

$$P(\text{"fun"}|+) = \frac{1+1}{9+20}$$

Scoring the test set:

$$P(-)P(S|-) = \frac{3}{5} \times \frac{2 \times 2 \times 1}{34^3} = 6.1 \times 10^{-5}$$

$$P(+)P(S|+) = \frac{2}{5} \times \frac{1 \times 1 \times 2}{29^3} = 3.2 \times 10^{-5}$$

# Optimizing for sentiment analysis

For tasks like sentiment, word **occurrence** is more important 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.

**Binary multinominal naive bayes**, or **binary NB**

- Clip our word counts at 1
- Note: this is different than Bernoulli naive bayes; see the textbook at the end of the chapter.

# Binary Multinomial Naïve Bayes: Learning

- From training corpus, extract *Vocabulary*

- Calculate $P(c_j)$ terms
  - For each $c_j$ in $C$ do
    $docs_j \leftarrow$ all docs with class $= c_j$

$$P(c_j) \leftarrow \frac{|docs_j|}{|\text{total \# documents}|}$$

- Calculate $P(w_k \mid c_j)$ terms
  - $Text_j \leftarrow$ single doc containing all $docs_j$
  - Remove duplicates in each doc:
  - For each word $w_k$ in *Vocabulary*
    - For each word type $w$ in doc
      - Retain only a single instance of $w$
    $n_k \leftarrow$ \# of occurrences of $w_k$ in $Text_j$

$$P(w_k \mid c_j) \leftarrow \frac{n_k + \alpha}{n + \alpha \, |Vocabulary|}$$

# Binary Multinomial Naive Bayes on a test document *d*

- First remove all duplicate words from *d*
- Then compute NB using the same equation:

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

# Binary multinominal naive Baves

**Four original documents:**

- &minus; it was pathetic the worst part was the boxing scenes
- &minus; no plot twists or great scenes
- + and satire and great plot twists
- + great scenes great film

Counts can still be 2! Binarization is within-doc!

# Text Classification and Naive Bayes

## More on Sentiment Classification

# Sentiment Classification: Dealing with Negation

I really **don't** like this movie

Negation changes the meaning of "like" to negative.

Negation can also change negative to positive-ish

- **Don't** dismiss this film
- **Doesn't** let us get bored

# Sentiment Classification: Dealing with 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.

Simple baseline method:

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

# Sentiment Classification: Lexicons

- Sometimes pre-built word lists work magic
- Called **lexicons**
- There are various publically available lexicons

# MPQA Subjectivity Cues Lexicon

Theresa Wilson, Janyce Wiebe, and Paul Hoffmann (2005). Recognizing Contextual Polarity in Phrase-Level Sentiment Analysis. Proc. of HLT-EMNLP-2005.

Riloff and Wiebe (2003). Learning extraction patterns for subjective expressions. EMNLP-2003.

- Home page https://mpqa.cs.pitt.edu/lexicons/subj_lexicon/

- 6885 words from 8221 lemmas, annotated for intensity (strong/weak)
  - 2718 positive
  - 4912 negative
- + : *admirable, beautiful, confident, dazzling, ecstatic, favor, glee, great*
- – : *awful, bad, bias, catastrophe, cheat, deny, envious, foul, harsh, hate*

12

# The General Inquirer

- Home page: http://www.wjh.harvard.edu/~inquirer

- List of Categories:  http://www.wjh.harvard.edu/~inquirer/homecat.htm

- Spreadsheet: http://www.wjh.harvard.edu/~inquirer/inquirerbasic.xls

- Categories:
    - Positiv (1915 words) and Negativ (2291 words)
    - Strong vs Weak, Active vs Passive, Overstated versus Understated
    - Pleasure, Pain, Virtue, Vice, Motivation, Cognitive Orientation, etc

- Free for Research Use

# Bing Liu Opinion Lexicon

Minqing Hu and Bing Liu. Mining and Summarizing Customer Reviews. ACM SIGKDD-2004.

- Bing Liu's Page on Opinion Mining
- http://www.cs.uic.edu/~liub/FBS/opinion-lexicon-English.rar

- 6786 words
    - 2006 positive
    - 4783 negative

14

# Using Lexicons in Sentiment Classification

**Add a feature that gets a count whenever a word from the lexicon occurs**

- **E.g., a feature called "this word occurs in the positive lexicon" or "this word occurs in the negative lexicon"**

Now all positive words (*good, great, beautiful, wonderful*) or negative words count for that feature.

Using 1-2 features isn't as good as using all the words.

- But when training data is sparse or not representative of the test set, dense lexicon features can help

# Naive Bayes in Other tasks: Spam Filtering

- SpamAssassin Features:
  - Mentions millions of (dollar) ((dollar) NN,NNN,NNN.NN)
  - From: starts with many numbers
  - Subject is all capitals
  - HTML has a low ratio of text to image area
  - "One hundred percent guaranteed"

# Naïve Bayes in Language ID

- Determining what language a piece of text is written in.

Features based on character n-grams do very well

- Important to train on lots of varieties of each language (world English, etc)

# Text Classification and Naïve Bayes

- Precision, Recall, and F measure

# Evaluation

- Let's consider just binary text classification tasks

- Imagine you're the CEO of Delicious Pie Company

- You want to know what people are saying about your pies

- So you build a "Delicious Pie" tweet detector
  - Positive class: tweets about Delicious Pie Co
  - Negative class: all other tweets

# The 2-by-2 confusion matrix

*gold standard labels*

|  |  | gold positive | gold negative |  |
|---|---|---|---|---|
| *system output labels* | system positive | true positive | false positive | $\textbf{precision} = \dfrac{tp}{tp+fp}$ |
|  | system negative | false negative | true negative |  |

$$\textbf{recall} = \dfrac{tp}{tp+fn}$$

$$\textbf{accuracy} = \dfrac{tp+tn}{tp+fp+tn+fn}$$

# Evaluation: Accuracy

- Why don't we use **accuracy** as our metric?

- Imagine we saw 1 million tweets
    - 100 of them talked about Delicious Pie Co.
    - 999,900 talked about something else

- We could build a dumb classifier that just labels every tweet "not about pie"
    - It would get 99.99% accuracy!!! Wow!!!!
    - But useless! Doesn't return the comments we are looking for!
    - That's why we use **precision** and **recall** instead

# Evaluation: Precision

- % of items the system detected (i.e., items the system labeled as positive) that are in fact positive (according to the human gold labels)

$$\textbf{Precision} = \frac{\text{true positives}}{\text{true positives} + \text{false positives}}$$

# Evaluation: Recall

- % of items actually present in the input that were correctly identified by the system.

$$\mathbf{Recall} = \frac{\text{true positives}}{\text{true positives} + \text{false negatives}}$$

# Why Precision and recall

- Our dumb pie-classifier
  - Just label nothing as "about pie"

Accuracy=99.99%

but

Recall = 0

- (it doesn't get any of the 100 Pie tweets)

Precision and recall, unlike accuracy, emphasize true positives:

- finding the things that we are supposed to be looking for.

# A combined measure: F

- F measure: a single number that combines P and R:

$$F_\beta = \frac{(\beta^2 + 1)PR}{\beta^2 P + R}$$

- We almost always use balanced $F_1$ (i.e., $\beta = 1$)

$$F_1 = \frac{2PR}{P + R}$$

# Development Test Sets ("Devsets") and Cross-validation
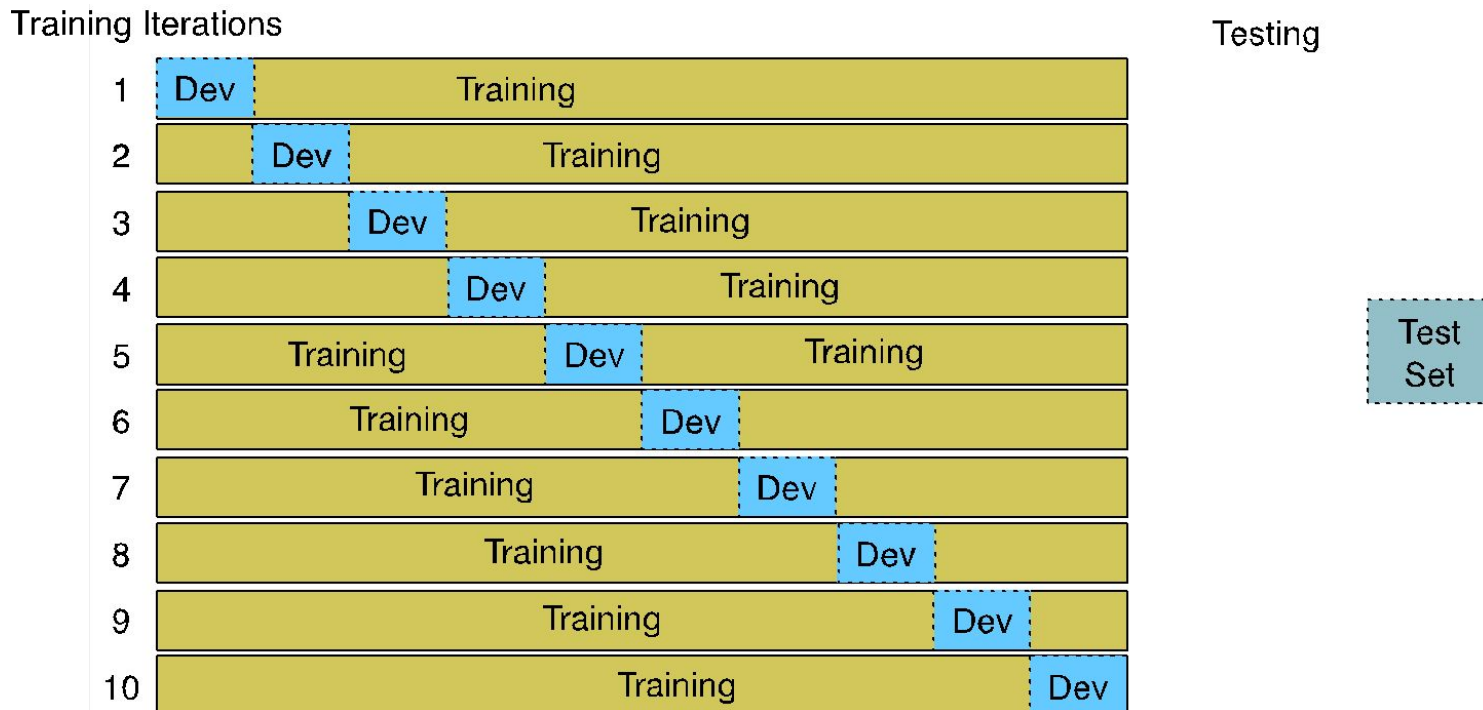
| Training set | Development Test Set | Test Set |
|---|---|---|

- Train on training set, tune on devset, report on testset
  - ==This avoids overfitting== ('tuning to the test set')
  - More conservative estimate of performance
  - But paradox: want as much data as possible for training, and ad much for dev; how to split?

# Cross-validation: multiple splits

- Pool results over splits, Compute pooled dev performance

**Text Classification and Naive Bayes**

Evaluation with more than two classes

# Confusion Matrix for 3-class classification



gold labels

|  | urgent | normal | spam |  |
|---|---|---|---|---|
| **urgent** | 8 | 10 | 1 | $\textbf{precision}_u = \dfrac{8}{8+10+1}$ |
| **normal** | 5 | 60 | 50 | $\textbf{precision}_n = \dfrac{60}{5+60+50}$ |
| **spam** | 3 | 30 | 200 | $\textbf{precision}_s = \dfrac{200}{3+30+200}$ |

*system output*

$$\textbf{recall}_u = \dfrac{8}{8+5+3} \qquad \textbf{recall}_n = \dfrac{60}{10+60+30} \qquad \textbf{recall}_s = \dfrac{200}{1+50+200}$$

# How to combine P/R from 3 classes to get one metric

- ## Macroaveraging:
  - compute the performance for each class, and then average over classes
- ## Microaveraging:
  - collect decisions for all classes into one confusion matrix
  - compute precision and recall from that table.

# Macroaveraging and Microaveraging

**Class 1: Urgent**

|  | true urgent | true not |
|---|---|---|
| system urgent | 8 | 11 |
| system not | 8 | 340 |

**Class 2: Normal**

|  | true normal | true not |
|---|---|---|
| system normal | 60 | 55 |
| system not | 40 | 212 |

**Class 3: Spam**

|  | true spam | true not |
|---|---|---|
| system spam | 200 | 33 |
| system not | 51 | 83 |

**Pooled**

|  | true yes | true no |
|---|---|---|
| system yes | 268 | 99 |
| system no | 99 | 635 |

$$\text{precision} = \frac{8}{8+11} = .42$$

$$\text{precision} = \frac{60}{60+55} = .52$$

$$\text{precision} = \frac{200}{200+33} = .86$$

$$\text{microaverage precision} = \frac{268}{268+99} = .73$$

$$\text{macroaverage precision} = \frac{.42+.52+.86}{3} = .60$$

**Text Classification and Naive Bayes**

Avoiding Harms in Classification

# Harms in sentiment classifiers

- Kiritchenko and Mohammad (2018) found that most sentiment classifiers assign lower sentiment and more negative emotion to sentences with African American names in them.

- This perpetuates negative stereotypes that associate African Americans with negative emotions

# Harms in toxicity classification

- Toxicity detection is the task of detecting hate speech, abuse, harassment, or other kinds of toxic language

- But some toxicity classifiers incorrectly flag as being toxic sentences that are non-toxic but simply mention identities like blind people, women, or gay people.

- This could lead to censorship of discussion about these groups.

# What causes these harms?

- Can be caused by:
  - Problems in the training data; machine learning systems are known to amplify the biases in their training data.
  - Problems in the human labels
  - Problems in the resources used (like lexicons)
  - Problems in model architecture (like what the model is trained to optimized)
- Mitigation of these harms is an open research area
- Meanwhile: **model cards**

# **Model Cards** (Mitchell et al., 2019)

- For each algorithm you release, document:
  - training algorithms and parameters
  - training data sources, motivation, and preprocessing
  - evaluation data sources, motivation, and preprocessing
  - intended use and users
  - model performance across different demographic or other groups and environmental situations

# Language Modeling

## Evaluation and Perplexity

# Evaluation: How good is our model?

- Does our language model prefer good sentences to bad ones?
    - Assign higher probability to "real" or "frequently observed" sentences
        - Than "ungrammatical" or "rarely observed" sentences?
- We train parameters of our model on a **training set**.
- We test the model's performance on data we haven't seen.
    - A **test set** is an unseen dataset that is different from our training set, totally unused.
    - An **evaluation metric** tells us how well our model does on the test set.

# Extrinsic evaluation of N-gram models

- Best evaluation for comparing models A and B
  - Put each model in a task
    - spelling corrector, speech recognizer, MT system
  - Run the task, get an accuracy for A and for B
    - How many misspelled words corrected properly

# Difficulty of extrinsic (in-vivo) evaluation of N-gram models

- Extrinsic evaluation
  - Time-consuming; can take days or weeks
- So
  - Sometimes use **intrinsic** evaluation: **perplexity**
  - Bad approximation
    - unless the test data looks **just** like the training data

# Intuition of Perplexity

- The Shannon Game:
  - How well can we predict the next word?

    I always order pizza with cheese and _____

    The 33rd President of the US was _____

    I saw a _____

  - Unigrams are terrible at this game. (Why?)
- A better model of a text
  - is one which assigns a higher probability to the word that actually occurs

mushrooms 0.1

pepperoni 0.1

anchovies 0.01

....

fried rice 0.0001

....

and 1e-100

# Perplexity

The best language model is one that best predicts an unseen test set

- Gives the highest P(sentence)

Perplexity is the inverse probability of the test set, normalized by the number of words:

$$PP(W) = P(w_1 w_2 ... w_N)^{-\frac{1}{N}}$$

$$= \sqrt[N]{\frac{1}{P(w_1 w_2 ... w_N)}}$$

Chain rule:

$$PP(W) = \sqrt[N]{\prod_{i=1}^{N} \frac{1}{P(w_i | w_1 ... w_{i-1})}}$$

For bigrams:

$$PP(W) = \sqrt[N]{\prod_{i=1}^{N} \frac{1}{P(w_i | w_{i-1})}}$$

**Minimizing perplexity is the same as maximizing probability**

# Perplexity as branching factor

- Let's suppose a sentence consisting of random digits

- What is the perplexity of this sentence according to a model that assign P=1/10 to each digit?

$$
\begin{aligned}
PP(W) &= P(w_1 w_2 \ldots w_N)^{-\frac{1}{N}} \\
&= (\frac{1}{10}^N)^{-\frac{1}{N}} \\
&= \frac{1}{10}^{-1} \\
&= 10
\end{aligned}
$$

# Lower perplexity = better model

- Training 38 million words, test 1.5 million words, WSJ

| N-gram Order | Unigram | Bigram | Trigram |
|---|---|---|---|
| Perplexity | 962 | 170 | 109 |

# Language Modeling

Interpolation, Backoff, and Web-Scale LMs

# Backoff and Interpolation

- Sometimes it helps to use **less** context
  - Condition on less context for contexts you haven't learned much about
- **Backoff:**
  - use trigram if you have good evidence,
  - otherwise bigram, otherwise unigram
- **Interpolation:**
  - mix unigram, bigram, trigram

- Interpolation works better

# Linear Interpolation

- Simple interpolation

$$\hat{P}(w_n|w_{n-2}w_{n-1}) = \lambda_1 P(w_n|w_{n-2}w_{n-1})$$
$$+\lambda_2 P(w_n|w_{n-1})$$
$$+\lambda_3 P(w_n)$$

$$\sum_i \lambda_i = 1$$

- Lambdas conditional on context:

$$\hat{P}(w_n|w_{n-2}w_{n-1}) = \lambda_1(w_{n-2}^{n-1})P(w_n|w_{n-2}w_{n-1})$$
$$+\lambda_2(w_{n-2}^{n-1})P(w_n|w_{n-1})$$
$$+\lambda_3(w_{n-2}^{n-1})P(w_n)$$

# How to set the lambdas?

- Use a **held-out** corpus

| Training Data | Held-Out | Test Data |
|---|---|---|

- Choose λs to maximize the probability of held-out data:
  - Fix the N-gram probabilities (on the training data)
  - <mark>Then search for λs that give largest probability to held-out set:</mark>

$$\log P(w_1 ... w_n \mid M(\lambda_1 ... \lambda_k)) = \sum_i \log P_{M(\lambda_1 ... \lambda_k)}(w_i \mid w_{i-1})$$

# Unknown words: Open versus closed vocabulary tasks

- If we know all the words in advanced
  - Vocabulary V is fixed
  - Closed vocabulary task
- Often we don't know this
  - **Out Of Vocabulary** = OOV words
  - Open vocabulary task
- Instead: create an unknown word token <UNK>
  - Training of <UNK> probabilities
    - Create a fixed lexicon L of size V
    - At text normalization phase, any training word not in L changed to  <UNK>
    - Now we train its probabilities like a normal word
  - At decoding time
    - If text input: Use UNK probabilities for any word not in training

# Huge web-scale n-grams

- How to deal with, e.g., Google N-gram corpus
- Pruning
    - Only store N-grams with count > threshold.
        - Remove singletons of higher-order n-grams
    - Entropy-based pruning
- Efficiency
    - Efficient data structures like tries
    - Bloom filters: approximate language models
    - Store words as indexes, not strings
        - Use Huffman coding to fit large numbers of words into two bytes
    - Quantize probabilities (4-8 bits instead of 8-byte float)