

Text Classification

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



Text Classification: definition

- Input:
 - a document *d*
 - a fixed set of classes $C = \{c_1, c_2, ..., c_J\}$

• *Output*: a predicted class *c* ∈ *C*





Classification Methods: Supervised Machine Learning

- Input:
 - a document d
 - a fixed set of classes $C = \{c_1, c_2, ..., c_J\}$
 - A training set of *m* hand-labeled documents (*d*₁, *c*₁),....,(*d*_m, *c*_m)
- Output:
 - a learned classifier $\gamma: d \rightarrow c$





Classification Methods: Supervised Machine Learning

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





Naïve Bayes Intuition

- Simple ("naïve") classification method based on Bayes rule
- Relies on very simple representation of document
 - Bag of words



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.



The bag of words representation







Multinomial Naïve Bayes Independence Assumptions

$$P(x_1, x_2, \dots, x_n \mid c)$$

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

$$P(x_1,...,x_n | c) = P(x_1 | c) \bullet P(x_2 | c) \bullet P(x_3 | c) \bullet ... \bullet P(x_n | c)$$



Learning the Multinomial Naïve Bayes Model

- First attempt: maximum likelihood estimates
 - simply use the frequencies in the data

$$\hat{P}(c_j) = \frac{doccount(C = c_j)}{N_{doc}}$$

$$\hat{P}(w_i \mid c_j) = \frac{count(w_i, c_j)}{\sum_{w \in V} count(w, c_j)}$$



Multinomial Naïve Bayes: Learning

- From training corpus, extract *Vocabulary*
- Calculate P(c_j) terms
 - For each c_j in C do $docs_j \leftarrow all docs with class = c_j$

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

• Calculate $P(w_k \mid c_j)$ terms

- $Text_j \leftarrow single doc containing all docs_j$
- For each word w_k in *Vocabulary* $n_k \leftarrow \#$ of occurrences of w_k in *Text*_j

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

Dan Jurafsky		Doc	Words	Class
$\hat{P}(c) = \frac{N_c}{N_c}$	Training	1	Chinese Beijing Chinese	С
		2	Chinese Chinese Shanghai	С
Harlanguage Protest		3	Chinese Macao	С
$\hat{P}(w \mid c) = \frac{count(w, c) + 1}{count(w, c) + 1}$		4	Tokyo Japan Chinese	j
count(c)+ V	Test	5	Chinese Chinese Chinese Tokyo Japan	?

Priors:

11

 $P(c) = \frac{3}{4} \frac{1}{4}$

Choosing a class: $P(c|d5) \propto 3/4 * (3/7)^3 * 1/14 * 1/14$ ≈ 0.0003

Conditional Probabilities:

P(Chinese | c) = (5+1) / (8+6) = 6/14 = 3/7P(Tokyo|c) = (0+1) / (8+6) = 1/14P(Japan|c) = (0+1) / (8+6) = 1/14P(Chinese | j) = (1+1) / (3+6) = 2/9P(Tokyo|j) = (1+1) / (3+6) = 2/9P(Japan|i) = (1+1)/(3+6) = 2/9

P(j|d5) ∝
$$1/4 * (2/9)^3 * 2/9 * 2/9 \approx 0.0001$$



Underflow Prevention: log space

- Multiplying lots of probabilities can result in floating-point underflow.
- Since log(xy) = log(x) + log(y)
 - Better to sum logs of probabilities instead of multiplying probabilities.
- Class with highest un-normalized log probability score is still most probable.

$$c_{NB} = \underset{c_{j} \in C}{\operatorname{argmax}} \log P(c_{j}) + \sum_{i \in positions} \log P(x_{i} | c_{j})$$

• Model is now just max of sum of weights



Summary: Naive Bayes is Not So Naive

- 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
 - But we will see other classifiers that give better accuracy



Text Classification: Evaluation



The 2-by-2 contingency table

	correct	not correct
selected	tp	fp
not selected	fn	tn



Precision and recall

• **Precision**: % of selected items that are correct **Recall**: % of correct items that are selected

	correct	not correct
selected	tp	fp
not selected	fn	tn



A combined measure: F

 A combined measure that assesses the P/R tradeoff is F measure (weighted harmonic mean):

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

- The harmonic mean is a very conservative average; see *IIR* § 8.3
- People usually use balanced F1 measure

• i.e., with
$$\beta = 1$$
 (that is, $\alpha = \frac{1}{2}$):

F = 2PR/(P+R)





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 ∈ 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





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 ∈ 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



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



21

Per class evaluation measures

Recall:

Fraction of docs in class *i* classified correctly:

Precision:

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

Accuracy: (1 - error rate) Fraction of docs classified correctly:



 C_{ii}

 $\sum c_{ji}$

 $\sum c_{ii}$



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.





Micro- vs. Macro-Averaging: Example

Class 1

Class 2

Micro Ave. Table

	Truth: yes	Truth: no		Truth: yes	Truth: no	
Classifier: yes	10	10	Classifier: yes	90	10	
Classifier: no	10	970	Classifier: no	10	890	

	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



Development Test Sets and Cross-validation

Training set

Development Test Set



- Metric: P/R/F1 or Accuracy
- Unseen test set
 - avoid overfitting ('tuning to the test set')
 - more conservative estimate of performance
- Cross-validation over multiple splits
 - Handle sampling errors from different datasets
 - Pool results over each split
 - Compute pooled dev set performance

Training Se ⁻	t Dev Test		
Training Set	Dev Test		
Dev Test	Training Set		

Test Set