# Naive Bayes

Evaluation: Precision, Recall, F-measure

---

## 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 conservative average
- People usually use balanced F1 measure
  - i.e., with $\beta = 1$ (that is, $\alpha = \frac{1}{2}$):  $F = 2PR/(P+R)$

# Classification Methods: Review

- *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 $\gamma{:}d \to c$

44

# Naïve Bayes: Review

- What type of classifier?
- Two simplifying assumptions (one specific to text classification)
- Two types of probabilities

$$c_{NB} = \underset{c \in C}{\operatorname{argmax}} P(c_j) \prod_{x \in X} P(x \mid c)$$

- Learning

45

# 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 $\gamma_c$ to distinguish $c$ from all other classes $c' \in C$
- Given test doc $d$,
  - Evaluate it for membership in each class using each $\gamma_c$
  - $d$ belongs to any class for which $\gamma_c$ returns true

46

# 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 $\gamma_c$ to distinguish $c$ from all other classes $c' \in C$
- Given test doc $d$,
  - Evaluate it for membership in each class using each $\gamma_c$
  - $d$ belongs to the one class with maximum score

47

# 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

Common categories (#train, #test)

48

- 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)

---

## 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

&#3;</BODY></TEXT></REUTERS>

# Confusion matrix c

- For each pair of classes $<c_1, c_2>$ 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 |

# 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:

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

# 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.

52

---

# 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:
- Microaveraged precision:
- Microaveraged score is dominated by score on common classes

53

# 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

54

# Development Test Sets and Cross-validation

Training set    Development Test Set    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 Set   Dev Test

Training Set   Dev Test

Dev Test   Training Set

Test Set

## Statistical Significance

- Suppose we have two classiers, classify1 and classify2.
- Is classify1 better? The "null hypothesis," denoted H0, is that it isn't. But if Accuracy1 >> Accuracy2 (or whatever your evaluation metricis instead of accuracy) we are tempted to believe otherwise.
- How much larger must A1 be than A2 to reject H0?
- Frequentist view: how (im)probable is the observed difference, given H0 = true?

56

# Text Classification

## Practical Issues

# The Real World

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

58

---

# 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

59

# Very little data?

- Use Naive Bayes
  - **On Discriminative vs. Generative classifiers: A comparison of logistic regression and naive Bayes (**Ng and Jordan 2002 NIPS)
- Get more labeled data
  - Find clever ways to get humans to label data for you
- Try semi-supervised machine learning methods

60

# A reasonable amount of data?
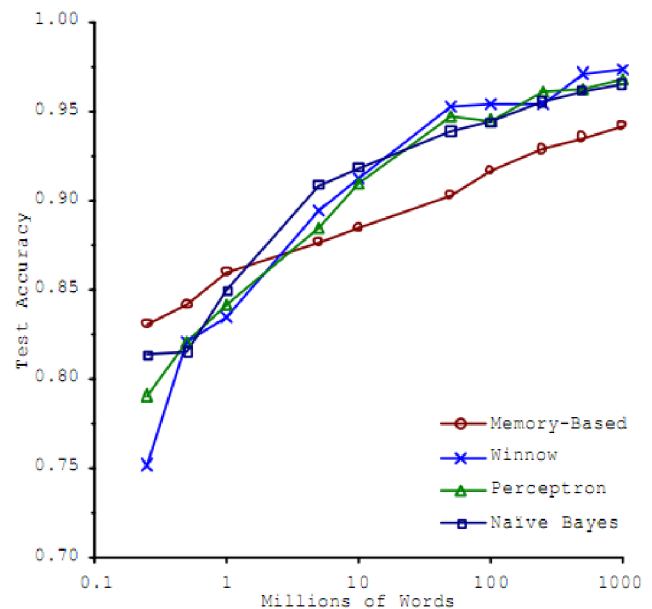
- Try more clever classifiers

61

# A huge amount of data?

- Can achieve high accuracy!
- At a cost (high train or test time for some methods)
- So Naive Bayes can come back into its own again!

62

---

# Accuracy as a function of data size

- With enough data
  - Classifier may not matter



63

Brill and Banko on spelling correction

# 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.

$$c_{NB} = \underset{c_j \in C}{\operatorname{argmax}} \log P(c_j) + \sum_{i \in positions} \log P(x_i \mid c_j)$$

- Model is now just max of sum of weights

---

# 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)