# ECE8813
# Statistical Natural Language Processing

## Lectures 19-20: Text Categorization

*Chin-Hui Lee*

School of Electrical and Computer Engineering

Georgia Institute of Technology

Atlanta, GA 30332, USA

chl@ece.gatech.edu

# What is Text Classification?

- We are given:
  - a fixed set of categories: $C = \{c_1, c_2, \ldots, c_n\}$
  - a document $d_j \in D$, where $D$ is the domain of documents

- We want to:
  - assign a Boolean value to the pair $<d_j, c_i>$
  - if the value is T, the the $d_j$ is classified under category $c_i$, otherwise it is not

- We essentially want to build categorization functions (classifiers) that assign these values

*Center of Signal and Image Processing*
*Georgia Institute of Technology*

CSIP

# An example: Is this mail spam?

From: lotterias-espana@zwallet.com [mailto:lotterias-espana@zwallet.com]

Sent: Wednesday, June 30, 2004 12:26 PM

Subject: FINAL AWARD WINNING NOTIFICATION

FROM: The Desk of the Managing Director

International Promotion Prize Award Dept.

Ref:LP523275/2003/ES

BATCH:02033/1PD

RE: Final award winning notification.

We are pleased to inform you about the release today the 30th of june 2004 of sweepstake Loteria Primitiva de España held on the 24th may 2004, your name attached to ticket number: 524- 412-56- ES, with serial number 4253/03 drew the lucky number:75-23-58-46-51, which consequently won the lottery in the 3rd category. You have therefore been approved for a lump sum pay out of €500,000.00 euros (five hundred thousand euros) in cash credited to file:lp523275/2003/es. This form is from a total cash prize of €2 million euros share! among the four international lucky winners in this category. furthermore, your lucky winning number falls within our European booklet representative office in Madrid - Spain as indicated in your play coupon. in view of this, your €500,000.00 (five hundred thousand euros) would be released to you by our private security and trust company which had insured your winning in your name with their office in Madrid - Spain, congratulations!

CSIP

# An Example: Language Identification

Die Ausstellung zeigt den Einfluss der Freien Universität auf wissenschafts- und gesellschaftspolitische Entwicklungen im nationalen und internationalen Raum. Im Mittelpunkt stehen die Gründung der FU als Reaktion auf die Relegation, Verhaftung und Drangsalierung demokratisch orientierter Studenten im Jahre 1948, ihre Rolle bei den Studentenunruhen 1968, die Folgen des Mauerfalls 1989 sowie künftige Pläne für den Wissenschaftsstandort Dahlem. Weitere thematische Schwerpunkte sind die Architektur des Universitätsgeländes mit Bauten aus sechs Jahrzehnten, das breit gefächerte Spektrum der angebotenen Wissenschaften, das Leben auf dem Campus sowie Habitus und Ritual der akademischen Welt damals und heute.

Giorno della Memoria - La Casa dello Studente, uno dei luoghi più evocativi legati alle vicende dell'oppressione nazista a Genova e in Liguria e di alto significato morale per la storia della Liberazione, sarà aperto al pubblico per iniziativa dell'Università di Genova e dell'ERSU e permetterà la visita alle "celle" della sede del Comando delle S.S. (1943-1945) 31 gennaio 2005

- Here the decision may be more than binary: given a set of languages (English, French, Italian, German, Spanish, Portuguese, etc.) what language does a given text belong to?

*Center of Signal and Image Processing*
*Georgia Institute of Technology*

CSIP

# Text Classification Examples

- Assign categories to web pages
  - e.g. sports:football, news:world:asia, finance, etc.
- Find the genre of a given web page
  - e.g. research page, news article, review page, etc.
- Categories may be binary
  - "spam", non-spam"
  - "interesting-to-me", "not-interesting-to-me"
  - "appropriate-for-kids", "not-appropriate-for-kids"
  - etc.

*Center of Signal and Image Processing*
*Georgia Institute of Technology*

CSIP

# Applications

- ## Document organisation
  - e.g. a newspaper that wants to "classified adds" put into categories such as "Car sales", "Property Rental", "Personals", etc.

- ## Text filtering
  - classify a stream of incoming documents depending on their relevance to the information consumer
  - typically a binary case (relevant – not relevant)
  - common to have a profile for the information consumer
    - the profile can be updated depending on the consumer's implicit or explicit relevance assessments on the provided information (adaptive filtering)

CSIP

# Applications (Cont.)

- Word sense disambiguation
  - e.g. "bank": financial institution, or river bank?
  - we can view word occurrence contexts as documents, and word senses as categories
  - we have a number of "documents" put in the correct "categories", and try to find the correct word sense for a new incoming word occurrence context
- Hierarchical categorisation of web pages
  - automatically classify pages under the hierarchical catalogue of e.g. Yahoo
  - searchers may find it easier to navigate in a hierarchy
  - the hypertextual nature of web pages is useful (one can take into advantage the links between pages)
  - the hierarchical structure of the categories is also useful
    - e.g. decompose the classification problem to a number of branching decisions at each internal node
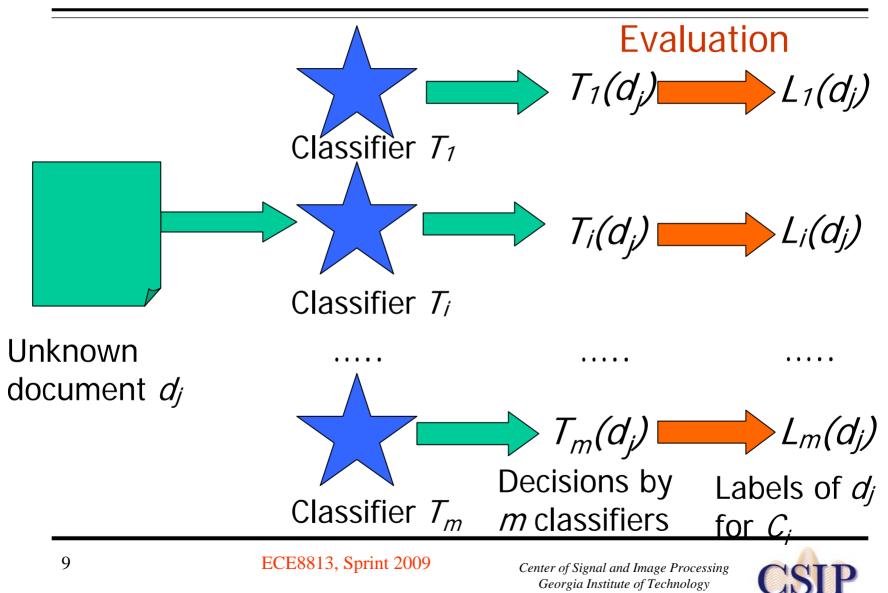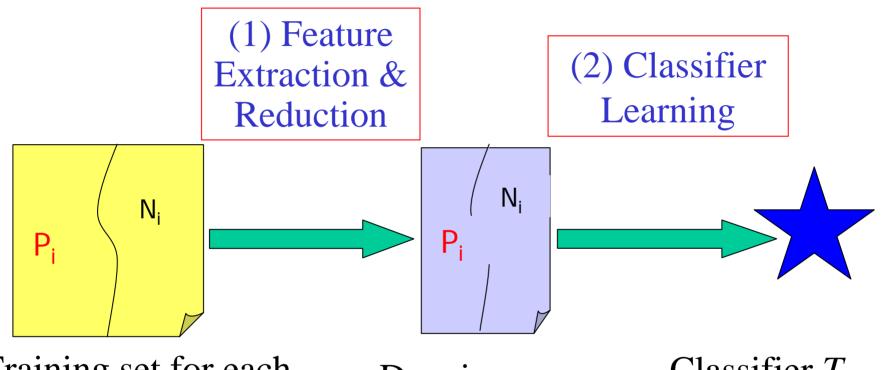
CSIP

# The Main Approach to Classification

- The machine learning approach
  - build a class for classifier for a category $c_i$ by observing the properties of the set of documents manually classified under ci (learning)
  - from these properties, get the properties that an unseen document should have in order to be classified under $c_i$
  - this is a case of supervised learning

- The knowledge engineering approach
  - need a large set of rules *if <> then <category>*
  - rules manually constructed
  - major drawback: *knowledge acquisition bottleneck*, i.e. how do you deal with new categories, different domain, etc.

*Center of Signal and Image Processing*
*Georgia Institute of Technology*

CSIP

# Text Categorization – Topic Identification



Unknown document $d_j$

Classifier $T_1$

Classifier $T_i$

.....

Classifier $T_m$

Evaluation

$T_1(d_j)$ → $L_1(d_j)$

$T_i(d_j)$ → $L_i(d_j)$

.....

$T_m(d_j)$ → $L_m(d_j)$

Decisions by $m$ classifiers

Labels of $d_j$ for $C_i$

ECE8813, Sprint 2009

*Center of Signal and Image Processing*
*Georgia Institute of Technology*

CSIP

# Text Categorization: Training Classifiers



(1) Feature Extraction & Reduction

(2) Classifier Learning

$N_i$

$P_i$

$N_i$

$P_i$

Training set for each category $C_i$, i= 1,…,m. (Positive +Negative)

Doc. in new feature space

Classifier $T_i$ for category $C_i$

*Center of Signal and Image Processing*
*Georgia Institute of Technology*

CSIP

# Multi-Class vs. Binary Decision Rule

- Multi-class (MC) classification

$$C(X) = \arg\max_j g_j(X;W), \quad 1 \le j \le m$$

if $\quad g_j(X;W) > g_{i \ne j}(X;W)$

- Special case: Binary classifier with LDF
  (*C+: positive class, C-: negative class*)

$$\begin{cases} f(W,X) \ge 0 & \text{label C+} \\ \text{Others} & \text{label C-} \end{cases}$$

➤ **Decision rule is a *discrete, non-differential function of the classifier parameters (need MFoM to optimize)***

*Center of Signal and Image Processing*
*Georgia Institute of Technology*

CSIP

# A Text Categorization Scenario

- Suppose you want to buy a cappuccino maker as a gift on the web

  – try Google for "cappuccino maker"

  – try "Yahoo! Shopping" for "cappuccino maker"

*Center of Signal and Image Processing*
*Georgia Institute of Technology*

CSIP

# Google Search Results

ECE8813, Sprint 2009

*Center of Signal and Image Processing*
*Georgia Institute of Technology*

# Yahoo Search Results



ECE8813, Sprint 2009    *Center of Signal and Image Processing*
*Georgia Institute of Technology*

# Observations

- Broad indexing & speedy search alone are not enough

- Organizational view of data is critical for effective retrieval

- Categorized data are easy for user to browse

- Category taxonomies become most central in well-known web sites (Yahoo!, Lycos, ...)

ECE8813, Sprint 2009   *Center of Signal and Image Processing*
*Georgia Institute of Technology*

CSIP

# Categorization/Classification

Given:

A description of an instance, $x \in X$, where $X$ is the *instance language* or *instance space*

Issue: how to represent text documents?

Example: A fixed set of categories:

$C = \{c_1, c_2, \ldots, c_n\}$

Determine:

The category of $x$: $c(x) \in C$, where $c(x)$ is a *categorization function* whose domain is $X$ and whose range is $C$

We want to know how to build categorization functions ("classifiers"), and often involve computing a score, or a goodness-of-fit function for each $x$ and each $c(x) \in C$

*Center of Signal and Image Processing*
*Georgia Institute of Technology*

CSIP

# Document Classification (Topic ID)

**Test Data:**

"planning language proof intelligence"

**Classes:**

(AI)          (Programming)          (HCI)

| ML | Planning | Semantics | Garb.Coll. | Multimedia | GUI |

**Training Data:**

| learning | planning | programming | garbage | ... | ... |
| intelligence | temporal | semantics | collection | | |
| algorithm | reasoning | language | memory | | |
| reinforcement | plan | proof... | optimization | | |
| network... | language... | | region... | | |

(Note: in real life there is often a hierarchy, not present in the above problem statement; and you get papers on ML approaches to Garb. Coll.)

ECE8813, Sprint 2009

*Center of Signal and Image Processing*
*Georgia Institute of Technology*

CSIP

# Text Categorization Examples

- Assign labels to each document or web-page:
- Labels are most often topics such as Yahoo-categories
  - *e.g., "finance," "sports," "news>world>asia>business"*
- Labels may be genres
  - *e.g., "editorials" "movie-reviews" "news"*
- Labels may be opinion
  - *e.g., "like", "hate", "neutral"*
- Labels may be domain-specific binary
  - *e.g., "interesting-to-me" : "not-interesting-to-me"*
  - *e.g., "spam" : "not-spam"*
  - *e.g., "contains adult language" :"doesn't"*

*Center of Signal and Image Processing*
*Georgia Institute of Technology*

CSIP

# Text Categorization Applications

- Web pages organized into category hierarchies
- Journal articles indexed by subject categories (e.g., the Library of Congress, MEDLINE, etc.)
- Responses to Census Bureau occupations
- Patents archived using *International Patent Classification*
- Patient records coded using international insurance categories
- E-mail message filtering
- News events tracked and filtered by topics

ECE8813, Sprint 2009

*Center of Signal and Image Processing*
*Georgia Institute of Technology*

CSIP

# Cost of Manual Text Categorization

- Yahoo!
  - 200 (?) people for manual labeling of Web pages
  - using a hierarchy of 500,000 categories
- MEDLINE (National Library of Medicine)
  - $2 million/year for manual indexing of journal articles
  - using MEdical Subject Headings (18,000 categories)
- Mayo Clinic
  - $1.4 million annually for coding patient-record events
  - using the International Classification of Diseases (ICD) for billing insurance companies
- US Census Bureau decennial census (1990: 22 million responses)
  - 232 industry categories and 504 occupation categories
  - $15 million if fully done by hand

CSIP

# Fast Entry is a Must to Compete

- Suppose you were starting a web search company, what would it take to compete with established engines?

    – You need to be able to establish a competing hierarchy *fast*

    – You will need a relatively *cheap* solution.  (Unless you have investors that want to pay millions of dollars just to get off the ground)

*Center of Signal and Image Processing*
*Georgia Institute of Technology*

CSIP

# *Semi*-Automatic Labeling

- Humans can encode knowledge of what constitutes membership in a category
- This encoding can then be automatically applied by a machine to categorize new examples
- For example...Text in a Web Page

"Saeco revolutionized *espresso* brewing a decade ago by introducing Saeco SuperAutomatic *machines*, which go from bean to *coffee* at the touch of a button. The all-new Saeco Vienna Super-Automatic home coffee and *cappucino machine* combines top quality with low price!"

CSIP

# Rule-based Approach to TC

- Rules
  - Rule 1:
    (*espresso* **or** *coffee* **or** *cappucino* ) **and** *machine\** $\Rightarrow$ *Coffee Maker*
  - Rule 2:
    *automat\** **and** *answering* **and** *machine\** $\Rightarrow$ *Phone*
  - Rule ...

- Experience has shown that defining rules by hands is
  - too time consuming
  - too difficult
  - inconsistency issues (as the rule set gets large)

# Expert System for TC (Late 1980s)

**Expert system for text categorization (late 1980s)**

ECE8813, Sprint 2009

*Center of Signal and Image Processing*
*Georgia Institute of Technology*

# From Knowledge Engineering to Statistical Learner

### DTree induction for text categorization (since 1994)

"eager learning"

training documents (categorized)

Induction algorithm (C5)

new document

Decision Trees (equivalent to Boolean rules)

| candidate category | system decision |
|---|---|
| Coffee Maker | Yes |
| Cookware | No |
| Phone | No |
| ... | ... |

output of the system

ECE8813, Sprint 2009

*Center of Signal and Image Processing*
*Georgia Institute of Technology*

CSIP

# A Comparison: Another Familiar Story

- For US Census Bureau Decennial Census 1990
  - 232 industry categories and 504 occupation categories
  - $15 million if fully done by hand

- Define classification rules manually:
  - Expert System AIOCS
  - Development time: 192 person-months (2 people, 8 years)
  - Accuracy = 47%

- Learn classification function
  - Nearest Neighbor classification (Creecy '92: 1-NN)
  - Development time: 4 person-months (Thinking Machine)
  - Accuracy = 60%

*Center of Signal and Image Processing*
*Georgia Institute of Technology*

CSIP

# An Example: Predicting Topics of News Stories

- Given: Collection of example news stories already labeled with a category (topic)

- Task: Predict category for news stories not yet labeled

- For our example, we'll only get to see the headline of the news story

- We'll represent categories using colors (All examples with the same color belong to the same category)

*Center of Signal and Image Processing*
*Georgia Institute of Technology*

CSIP

# Our Labeled Examples

| | | | | |
|---|---|---|---|---|
| Amatil Proposes Two-for-Five Bonus Share Issue | Citibank Norway Unit Loses Six Mln Crowns in 1986 | Japan Ministry Says Open Farm Trade Would Hit U.S. | Vieille Montagne Says 1986 Conditions Unfavourable | Jardine Matheson Said It Sets Two-for-Five Bonus Issue Replacing "B" Shares |
| Anheuser-Busch Joins Bid for San Miguel | Italy's La Fondiaria to Report Higher 1986 Profits | Isuzu Plans No Interim Dividend | Senator Defends U.S. Mandatory Farm Control Bill | Bowater Industries Profit Exceed Expectations |

ECE8813, Sprint 2009

*Center of Signal and Image Processing*
*Georgia Institute of Technology*

CSIP

# Topic Prediction

*Center of Signal and Image Processing*
*Georgia Institute of Technology*

# Topic Prediction with Evidence

| | | | | |
|---|---|---|---|---|
| | | Senate Panel Studies Loan Rate, Set Aside Plans | | |

| | | | | |
|---|---|---|---|---|
| Amatil Proposes Two-for-Five Bonus Share Issue | Citibank Norway Unit Loses Six Mln Crowns in 1986 | Japan Ministry Says Open Farm Trade Would Hit U.S. | Vieille Montagne Says 1986 Conditions Unfavourable | Jardine Matheson Said It Sets Two-for-Five Bonus Issue Replacing "B" Shares |
| Anheuser-Busch Joins Bid for San Miguel | Italy's La Fondiaria to Report Higher 1986 Profits | Isuzu Plans No Interim Dividend | Senator Defends U.S. Mandatory Farm Control Bill | Bowater Industries Profit Exceed Expectations |

ECE8813, Sprint 2009

*Center of Signal and Image Processing*
*Georgia Institute of Technology*

CSIP

# Handling Documents with Multiple Classes

| | | | | |
|---|---|---|---|---|
| Amatil Proposes Two-for-Five Bonus Share Issue | Citibank Norway Unit Loses Six Mln Crowns in 1986 | Japan Ministry Says Open Farm Trade Would Hit U.S. | Vieille Montagne Says 1986 Conditions Unfavourable | Jardine Matheson Said It Sets Two-for-Five Bonus Issue Replacing "B" Shares |
| Anheuser-Busch Joins Bid for San Miguel | Italy's La Fondiaria to Report Higher 1986 Profits | Isuzu Plans No Interim Dividend | Senator Defends U.S. Mandatory Farm Control Bill | Bowater Industries Profit Exceed Expectations |

ECE8813, Sprint 2009

*Center of Signal and Image Processing*
*Georgia Institute of Technology*

CSIP

# Document Representation

- Usually, an example is represented as a series of feature-value pairs.  The features can be arbitrarily abstract (as long as they are easily computable) or very simple.

- For example, the features could be the set of all words and the values, their number of occurrences in a particular document.

| Japan Firm Plans to Sell U.S. Farmland to Japanese | → Representation → | Farmland:1 Firm:1 Japan:1 Japanese:1 Plans:1 Sell:1 To:2 U.S.:1 |
|---|---|---|

ECE8813, Sprint 2009

*Center of Signal and Image Processing*
*Georgia Institute of Technology*

CSIP

# Performance Evaluation

- Suppose we have a set $D$ of labeled documents that we use as our training set for 1-NN.  We need an idea of how well this system will perform in the future.  So, we go through $D$ and make predictions for each document

  - What will our accuracy be?

  - Is this a fair assessment of its performance? (i.e. is it likely that the performance will be within a small tolerance of what we've estimated)

*Center of Signal and Image Processing*
*Georgia Institute of Technology*

CSIP

# Classification Performance Measures

- Given *n* test documents and *m* classes in consideration, a classifier makes $n \times m$ binary decisions. A two-by-two contingency table can be computed for each class

|            | truly YES | truly NO |
|------------|-----------|----------|
| system YES | a         | b        |
| system NO  | c         | d        |
|            |           |          |

CSIP

# Classification Performance Measures

- Recall = a/(a+c) where a + c > 0 (o.w. undefined).
  - Did we find all of those that belonged in the class?

- Precision = a/(a+b) where a+b>0 (o.w. undefined).
  - Of the times we predicted it was "in class", how often are we correct?

- Accuracy = (a + d) / n
  - When one classes is overwhelmingly in the majority, this may not paint an accurate picture.

- Others: miss, false alarm (fallout), error, F-measure, area under PR ROC curve, break-even point, ...

*Center of Signal and Image Processing*
*Georgia Institute of Technology*

CSIP

# Global Performance Measures

- Global Performance Measures

| Category set $C = \{C_1,\dots,C_m\}$ | | Manual Labels | |
|---|---|---|---|
| | | C+ | C- |
| Classifier Judgments | C+ | $TP = \sum_{i=1}^{m} \mathrm{TP}_i$ | $FP = \sum_{i=1}^{m} \mathrm{FP}_i$ |
| | C- | $FN = \sum_{i=1}^{m} \mathrm{FN}_i$ | $TN = \sum_{i=1}^{m} \mathrm{TN}_i$ |

ECE8813, Sprint 2009

*Center of Signal and Image Processing*
*Georgia Institute of Technology*

# Local Performance Measures in TC

- Local Performance Measures for Category $C_i$

| Category $C_i$ | | Manual Labels | |
|---|---|---|---|
| | | C+ | C- |
| Classifier Judgments | C+ | $TP_i$ | $FP_i$ |
| | C- | $FN_i$ | $TN_i$ |

$$Pr_i = \frac{TP_i}{TP_i + FP_i}$$

$$Re_i = \frac{TP_i}{TP_i + FN_i}$$

$$F_{1i} = \frac{2\,Re_i\,Pr_i}{Re_i + Pr_i}$$

- *Precision*, *Recall* and *F1*

ECE8813, Sprint 2009

*Center of Signal and Image Processing*
*Georgia Institute of Technology*

CSIP

# Summary Performance Measures

- Micro-averaging

$$\text{Pr}^u = \frac{TP}{TP + FP}, \ \text{Re}^u = \frac{TP}{TP + FN}, \quad F_1^{\mu} = \frac{2TP}{FP + FN + 2TP}.$$

- Macro-averaging

$$\text{Pr}^M = \frac{\sum_{i=1}^{m} \text{Pr}_i}{m}, \quad \text{Re}^M = \frac{\sum_{i=1}^{m} \text{Re}_i}{m},$$

$$F_1^M = \frac{2\,\text{Re}^M \, \text{Pr}^M}{\text{Re}^M + \text{Pr}^M} = \frac{2\sum_{i=1}^{m} \text{Re}_i \sum_{i=1}^{m} \text{Pr}_i}{m(\sum_{i=1}^{m} \text{Pr}_i + \sum_{i=1}^{m} \text{Re}_i)}.$$

ECE8813, Sprint 2009

*Center of Signal and Image Processing*
*Georgia Institute of Technology*

CSIP

# Summary of Performance Measure

$$P = \frac{\sum_{i=1}^{|C|} TP_i}{\sum_{i=1}^{|C|} (TP_i + FP_i)}$$

$$R = \frac{\sum_{i=1}^{|C|} TP_i}{\sum_{i=1}^{|C|} (TP_i + FN_i)}$$

microaveraging

$$P = \frac{\sum_{i=1}^{|C|} P_i}{|C|}$$

$$R = \frac{\sum_{i=1}^{|C|} R_i}{|C|}$$

macroaveraging

- These two methods can give different results
- It is essential to make clear which method one uses when reporting P and R values for classification

*Center of Signal and Image Processing*
*Georgia Institute of Technology*

CSIP

# Hold-out Sets (Validation Data)

- Estimating our performance on data we used in training is likely to give us a very skewed estimate of the final system's performance.  As a result, if we have a set of labeled data, $D$, we typically split it into a training set, $D_{train}$, and a *hold-out* set, $D_{test}$

- $D_{train}$ is the only data given to the classifier for training. $D_{test}$ can then be used to estimate performance independently.  Once performance estimates are used to choose the best classifier, the final classifier is usually trained over all of $D$ before deployment (more data generally means better performance – so our estimate was pessimistic)

CSIP

# Empirically Tuning Parameters

- When parameters need to be empirically tuned as a part of training (e.g. choosing *k*), the performance of each possible choice needs to be estimated. For the same reasons as above, the classifier cannot simply check the performance on $D_{train}$ to estimate future performance. Therefore $D_{train}$ is usually subdivided into a portion used to train and another portion used for picking optimal parameters (usually referred to as the *validation* set)

- After setting the parameters, the classifier trains over all of $D_{train}$ before returning to the function that will evaluate its performance over $D_{test}$

*Center of Signal and Image Processing*
*Georgia Institute of Technology*

CSIP

# Approaches to Automated Text Categorization

- Regression based on Least Squares Fit (1991)
- Nearest Neighbor Classification (1992)
- Bayesian Probabilistic Models (1992)
- Symbolic Rule Induction (1994)
- Neural Networks (1995)
- Rocchio approach (traditional IR, 1996)
- Support Vector Machines (1997)
- Boosting or Bagging (1997)
- Hierarchical Language Modeling (1998)
- First-Order-Logic Rule Induction (1999)
- Maximum Entropy (1999)
- Hidden Markov Models (1999)
- Error-Correcting Output Coding (1999)
- Maximal Figure-of-Merit Learning (2003)
  - ...

*Center of Signal and Image Processing*
*Georgia Institute of Technology*

CSIP

# Classification Types

- ## Single label vs. multi-label
  - Exactly 1 category assigned to each document vs. 0 to |C|
- ## Binary vs. multi-way classification
  - Binary: a special case of single label, $d_j \in D$ is assigned either to $c_i$ or to its complement (e.g.spam – non spam)
- ## Document-pivoted (DPC) vs. category pivoted (CPC)
  - Given a $d_j \in D$ , we want to find all the $c_i \in C$ under which it should be classified (document-pivoted)
    - DPC is suitable when documents become available at different moments in time, e.g. filtering e-mail
  - Given a $c_i \in C$ , we want to find all the $d_j \in D$ under that should be classified under it (category-pivoted)
    - CPC is suitable when new categories are likely to be be added to $C$

*Center of Signal and Image Processing*
*Georgia Institute of Technology*

CSIP

# Related Work on Classifier Design

- Decision Tree: available tools, C4.5, CART, ID3

  Linear discriminative function:     $f(X,W) = \sum_{i=1}^{D} w_i x_i - w_0$

- *K*-Nearest Neighbor (*k*NN)
- Naïve Bayes: simple distributions for each class
- Support Vector Machine (SVM)
- Linear Discriminative Function (LDF)
- Artificial Neural Networks (ANN)
- Tree Classifiers (CART)
- Semantic Perceptron Net (SPN)
- Others: HMM, kernels, Discriminative Training

*Center of Signal and Image Processing*
*Georgia Institute of Technology*

CSIP

# Linear vs. Nonlinear Classifiers

- Linear classifiers if
  - all data points can be correctly classified by a linear decision boundary
  - simpler, less parameters

- Non-linear otherwise
  - more accurate
  - more complicated, more parameters

*Center of Signal and Image Processing*
*Georgia Institute of Technology*

CSIP

# Linear Case – An Example

Linear Decision boundary

● Class1
● Class2

ECE8813, Sprint 2009

*Center of Signal and Image Processing*
*Georgia Institute of Technology*

CSIP

# Nonlinear Case – An Example



Non Linear Classifier

● Class1
○ Class2

ECE8813, Sprint 2009

*Center of Signal and Image Processing*
*Georgia Institute of Technology*

CSIP

# Linear: Support Vector Machines (SVM)

- Find the hyperplane that maximizes the margin between negative and positive training examples

- Lines represent decision surfaces

- Decision surface $\sigma_1$ is the best possible one
  - middle element of the widest possible set of parallel decision surfaces
  - min. distance to any training example is maximum
  - Small boxes indicate the support vectors, the set of training examples that are used in the decision

From (Sebastiani, 2002)

*Center of Signal and Image Processing*
*Georgia Institute of Technology*

CSIP

# Supporting Vector Machines

- Strengths
    - very effective classification
    - can scale up to data of high dimensionality
    - dimensionality reduction is normally not needed
- Weaknesses
    - can be computationally expensive, but efficient algorithms have been proposed

*Center of Signal and Image Processing*
*Georgia Institute of Technology*

**CSIP**

# Key Components of Nearest Neighbor

- "Similar" item: We need a functional definition of "similarity" if we want to apply this automatically.

- How many neighbors do we consider?

- Does each neighbor get the same weight?

- All categories in neighborhood? Most frequent only? How do we make the final decision?

CSIP

# Nearest Neighbor Classification

- *Instance-Based Learning, Lazy Learning*
  - well-known approach to pattern recognition
  - initially by Fix and Hodges (1951)
  - theoretical error bound analysis by Duda & Hart (1957)
  - applied to text categorization in early 90's
  - strong baseline in benchmark evaluations
  - among top-performing methods in TC evaluations
  - scalable to large TC applications

*Center of Signal and Image Processing*
*Georgia Institute of Technology*

CSIP

# 1-Nearest Neighbor

- Looking back at our example

  – Did anyone try to find the most similar labeled item and then just guess the same color?

  – This is 1-Nearest Neighbor

| Senate Panel Studies Loan Rate, Set Aside Plans |
|---|

| | | | | |
|---|---|---|---|---|
| Amatil Proposes Two-for-Five Bonus Share Issue | Citibank Norway Unit Loses Six Mln Crowns in 1986 | Japan Ministry Says Open Farm Trade Would Hit U.S. | Vieille Montagne Says 1986 Conditions Unfavourable | Jardine Matheson Said It Sets Two-for-Five Bonus Issue Replacing "B" Shares |
| Anheuser-Busch Joins Bid for San Miguel | Italy's La Fondiaria to Report Higher 1986 Profits | Isuzu Plans No Interim Dividend | Senator Defends U.S. Mandatory Farm Control Bill | Bowater Industries Profit Exceed Expectations |

*Center of Signal and Image Processing*
*Georgia Institute of Technology*

CSIP

# 1-Nearest Neighbor (Graphically)

1-NN: assign "x" (new point) to the class of it nearest neighbor



assign "x" to "white"

decision surface divided by points
("Voronoi diagram")

*Center of Signal and Image Processing*
*Georgia Institute of Technology*

# K-Nearest Neighbor: *Majority* Voting Scheme

K-Nearest Neighbor using a *majority* voting scheme



k=1: majority vote for "white"

k=5; majority vote for "black"

k=10: even votes for both

ECE8813, Sprint 2009

*Center of Signal and Image Processing*
*Georgia Institute of Technology*

# K-NN : Weighted-Sum Voting Scheme

k-NN using a weighted-sum voting scheme



kNN (k = 5)

Assign "white" to x because the weighted sum of "whites" is larger then the sum of "blacks".

Each neighbor is given a weight according to its nearness.

*Center of Signal and Image Processing*
*Georgia Institute of Technology*

CSIP

# Category Scoring for Weighted-Sum

- The score for a category is the sum of the similarity scores between the point to be classified and all of its k-neighbors that belong to the given category.

- To restate: 
$$score\,(c\mid x) = \sum_{d \in kNN\ of\ x} sim\,(x,d)\,I(d,c)$$

where *x* is the new point; *c* is a class (*e.g. black* or *white);*
*d* is a classified point among the k-nearest neighbors of *x*;
*sim(x,d)* is the similarity between *x* and *d;*
*I(d,c)* = 1 iff point *d* belongs to class *c*;
*I(d,c)* = 0 otherwise.

*Center of Signal and Image Processing*
*Georgia Institute of Technology*

CSIP

# The kth Nearest Neighbor Decision Rule (Fix and Hodges, 1951)

- Define a metric to measure "closeness" between any two points

- Fix $k$ (empirically chosen)

- Given a new point $x$ and a training set of classified points
  - Find the $k$ nearest neighbors (kNN) to $x$ in the training set
  - Classify $x$ as class $y$ if more of the nearest neighbors are in class $y$ than in any other classes (*majority vote*)

CSIP

# kNN for Text Categorization
## (Yang, SIGIR-1994)

- Represent documents as points (vectors)

- Define a similarity measure for pair-wise documents

- Tune parameter $k$ for optimizing classification effectiveness

- Choose a voting scheme (e.g., weighted sum) for scoring categories

- Threshold on the scores for classification decisions

*Center of Signal and Image Processing*
*Georgia Institute of Technology*

CSIP

# **Thresholding for Classification Decisions**

- Alternative thresholding strategies:
  - Rcut: For each document to be categorized, rank candidate categories by score, and assign YES to the top-$m$ categories (where $m$ is some fixed number)

  - Pcut: Applies only when we have a whole batch of documents to be categorized.  Make the category assignments proportional to the category distribution in the training set (i.e. if 1/4$^{th}$ of the training documents were in the category "Coffee Maker" then we will assign 1/4$^{th}$ of the documents in this batch to the "Coffee Maker" category)

  - Scut: For each category, choose a threshold score (empirically).  Any document with a category score that surpasses its respective threshold will be predicted to be a member of that category

*Center of Signal and Image Processing*
*Georgia Institute of Technology*

CSIP

# Key Components (Revisited)

- Functional definition of "similarity"
  - e.g. cos, Euclidean, kernel functions, ...

- How many neighbors do we consider?
  - Value of $k$ determined *empirically* (see methodology section)

- Does each neighbor get the same weight?
  - Weighted-sum or not

- All categories in neighborhood?   Most frequent only? How do we make the final decision?
  - Rcut, Pcut, or Scut

*Center of Signal and Image Processing*
*Georgia Institute of Technology*

CSIP

# Pros of kNN

- Simple and effective (among top-5 in benchmark evaluations)
    - Non-linear classifier (vs linear)
    - Local estimation (vs global)
    - Non-parametric (very few assumptions about data)
    - Reasonable similarity measures (borrowed from IR)

- Computation (time & space) linear to the size of training data

- Low cost for frequent re-training, i.e., when categories and training documents need to be updated (common in Web environment and e-commerce applications)

*Center of Signal and Image Processing*
*Georgia Institute of Technology*

CSIP

# Cons of kNN:

- Online response is typically slower than *eager learning* algorithms
  - Trade-off between off-line training cost and online search cost

- Scores are not normalized (probabilities)
  - Comparing directly to and combining with scores of other classifiers is an open problem

- Output not good in explaining why a category is relevant
  - Compared to DTree, for example (take this with a grain of salt)

*Center of Signal and Image Processing*
*Georgia Institute of Technology*

CSIP

# Bayesian Methods

- Learning and classification methods based on probability theory

- Bayes theorem plays a critical role in probabilistic learning and classification

- Build a *generative model* that approximates how data is produced

- Uses *prior* probability of each category given no information about an item

- Categorization produces a *posterior* probability distribution over the possible categories given a description of an item

CSIP

# Bayes' Rule once more

$$P(C, X) = P(C \mid X)P(X) = P(X \mid C)P(C)$$

$$P(C \mid X) = \frac{P(X \mid C)P(C)}{P(X)}$$

ECE8813, Sprint 2009

*Center of Signal and Image Processing*
*Georgia Institute of Technology*

CSIP

# *Maximum a posteriori* Decision Rule

$$h_{MAP} \equiv \operatorname*{argmax}_{h \in H} P(h \mid D)$$

$$= \operatorname*{argmax}_{h \in H} \frac{P(D \mid h)P(h)}{P(D)}$$

$$= \operatorname*{argmax}_{h \in H} P(D \mid h)P(h)$$

CSIP

# *Maximum likelihood* Hypothesis

If all hypotheses are a priori equally likely, we only need to consider the $P(D|h)$ term:

$$h_{ML} \equiv \underset{h \in H}{\mathrm{argmax}}\, P(D \mid h)$$

CSIP

# Naive Bayes Classifiers

Task: Classify a new instance $D$ based on a tuple of attribute values $D = \langle x_1, x_2, \ldots, x_n \rangle$ into one of the classes $c_j \in C$

$$c_{MAP} = \operatorname*{argmax}_{c_j \in C} P(c_j \mid x_1, x_2, \ldots, x_n)$$

$$= \operatorname*{argmax}_{c_j \in C} \frac{P(x_1, x_2, \ldots, x_n \mid c_j) P(c_j)}{P(x_1, x_2, \ldots, x_n)}$$

$$= \operatorname*{argmax}_{c_j \in C} P(x_1, x_2, \ldots, x_n \mid c_j) P(c_j)$$

*Center of Signal and Image Processing*
*Georgia Institute of Technology*

CSIP

# Naïve Bayes Classifier: Assumption

- $P(c_j)$: Can be estimated from the frequency of classes in the training examples

- $P(x_1, x_2, \ldots, x_n / c_j)$: O($|X|^n \cdot |C|$) parameters, Could only be estimated if a very, very large number of training examples was available

- Naïve Bayes Conditional Independence Assumption:
  - Assume that the probability of observing the conjunction of attributes is equal to the product of the individual probabilities $P(x_i|c_j)$.

*Center of Signal and Image Processing*
*Georgia Institute of Technology*

CSIP

# The Naïve Bayes Classifier



$$P(X_1, \ldots, X_5 \mid C) = P(X_1 \mid C) \bullet P(X_2 \mid C) \bullet \cdots \bullet P(X_5 \mid C)$$

*Center of Signal and Image Processing*
*Georgia Institute of Technology*

CSIP

# Learning the Model



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

$$\hat{P}(c_j) = \frac{N(C = c_j)}{N}$$

$$\hat{P}(x_i \mid c_j) = \frac{N(X_i = x_i, C = c_j)}{N(C = c_j)}$$

*Center of Signal and Image Processing*
*Georgia Institute of Technology*

CSIP

# Smoothing to Avoid Overfitting

$$\hat{P}(x_i \mid c_j) = \frac{N(X_i = x_i, C = c_j) + 1}{N(C = c_j) + k}$$

# of values of $X_i$

- Somewhat more subtle version

overall fraction in data where $X_i = x_{i,k}$

$$\hat{P}(x_{i,k} \mid c_j) = \frac{N(X_i = x_{i,k}, C = c_j) + m p_{i,k}}{N(C = c_j) + m}$$

extent of "smoothing"

ECE8813, Sprint 2009    *Center of Signal and Image Processing*
*Georgia Institute of Technology*

CSIP

# Class Conditional Multinomial NB



- Effectively, the probability of each class is done as a class-specific unigram language model

ECE8813, Sprint 2009

*Center of Signal and Image Processing*
*Georgia Institute of Technology*

CSIP

# Basic NB Classifiers to Classify Text

- Attributes are text positions, values are words

$$c_{NB} = \underset{c_j \in C}{\operatorname{argmax}} P(c_j) \prod_i P(x_i \mid c_j)$$

$$= \underset{c_j \in C}{\operatorname{argmax}} P(c_j) P(x_1 = \text{"our"} \mid c_j) \cdots P(x_n = \text{"text"} \mid c_j)$$

- Still too many possibilities
- Assume that classification is *independent* of the positions of the words
  - Use same parameters for each position
  - Result is bag of words model (over tokens not types)

*Center of Signal and Image Processing*
*Georgia Institute of Technology*

# Naïve Bayes: Learning

- From training corpus, extract *Vocabulary*
- Calculate required $P(c_j)$ and $P(x_k \mid c_j)$ terms
  - For each $c_j$ in $C$ do
    - $docs_j \leftarrow$ subset of documents for which the target class is $c_j$
    - $$P(c_j) \leftarrow \frac{\mid docs_j \mid}{\mid \text{total \# documents} \mid}$$

  $Text_j \leftarrow$ single document containing all $docs_j$
  for each word $x_k$ in *Vocabulary*

  $n_k \leftarrow$ number of occurrences of $x_k$ in $Text_j$

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

ECE8813, Sprint 2009

*Center of Signal and Image Processing*
*Georgia Institute of Technology*

# Naïve Bayes: Classifying

- positions ← all word positions in current document which contain tokens found in *Vocabulary*

- Return $c_{NB}$, where

$$c_{NB} = \underset{c_j \in C}{\mathrm{argmax}}\ P(c_j) \prod_{i \in positions} P(x_i \mid c_j)$$

*Center of Signal and Image Processing*
*Georgia Institute of Technology*

# Vector Space Representation



ECE8813, Sprint 2009

*Center of Signal and Image Processing*
*Georgia Institute of Technology*

CSIP

# Word-Document Co-Occurrence

- Given - *N* documents, vocabulary size *M*

- Generate a word-documents co-occurrence matrix **W**

$$\mathbf{W} = \begin{array}{c} \\ w_1 \\ w_2 \\ \vdots \\ w_M \end{array}$$

$d_1 \quad d_2 \quad \ldots \ldots \quad d_N$

$n_{ij} \implies n_{i\cdot} \text{ (row sum)}$

$n_{\cdot j} \text{ (column sum)}$

*Center of Signal and Image Processing*
*Georgia Institute of Technology*

CSIP

# LSA Count in the Column Vector

- A trick from Information Retrieval
  - Each **document** (paragraph or sentence) in the training document corpus is a length-$M$ vector

aardvark  abacus  abandoned  abbot  abduct  above  ...  zygote  zymurgy

(0,   3,   3,   1,   0,   7,   . . .   1,   0)

a single document

*Center of Signal and Image Processing*
*Georgia Institute of Technology*

CSIP

# LSA Mathematical Framework

- LSA Matrix (also known as Routing Matrix) $C$

$$c_{ij} = (1 - \varepsilon_i) n_{ij} / n_{\cdot j} \text{ (scaling and normalization)}$$

- number of times word $w_i$ occurs in $A_j$ : $\quad n_{ij}$
- total number of words present in $A_j$ : $\quad n_{\cdot j} \text{ (column sum)}$
- total number of $w_i$ occurs in $A$ : $\quad n_{i\cdot} \text{ (row sum)}$
- "indexing" power of $w_i$ in corpus $A$ : $\quad \eta_i = 1 - \varepsilon_i$
- normalized entropy:

$$\varepsilon_i = -\frac{1}{\log N} \sum_{j=1}^{N} \frac{n_{ij}}{n_{i\cdot}} \log \frac{n_{ij}}{n_{i\cdot}} \quad 0 \leq \varepsilon_i \leq 1$$

$$\begin{cases} \varepsilon_i = 0 & \text{if } n_{ij} = n_{i\cdot} \quad \text{maximum indexing power} \\ \varepsilon_i = 1 & \text{if } n_{ij} = \frac{n_{i\cdot}}{N} \quad \text{no power (equally probable)} \end{cases}$$

ECE8813, Sprint 2009

*Center of Signal and Image Processing*
*Georgia Institute of Technology*

CSIP

# Semantic Similarity Measure

- To find similarity between two documents, project them in LS space

- Then calculate the cosine measure between their projection

- With this measure, various problems can be addressed e.g., natural language understanding, cognitive modeling etc.

*Center of Signal and Image Processing*
*Georgia Institute of Technology*

CSIP

# Confidence Scoring

- Inner Product: $s(x, y) = x \bullet y^t$

- Cosine:

$$s(x, y) = \frac{x \bullet y^t}{|x \| y|} \quad \text{or} \quad \cos^{-1}[s(x, y)]$$

- Confidence Scoring: Sigmoid function fitting

$$Conf\ (s; \alpha, \beta) = [1 + e^{-(\alpha s + \beta)}]^{-1}$$

- Other Scores
  - Euclidean, Manhattan, etc.

- Generalized Scores
  - between any two vectors: $s(x, y) = f(x, y; \Gamma)$

*Center of Signal and Image Processing*
*Georgia Institute of Technology*

CSIP

# Similarity in LSA

- The vector of a passage is the vector sum of the vectors standing for the words it contains

- Similarity of any two words or two passages is computed as the cosine between them in the semantic space:
  - Identical meaning: value of cosine = 1
  - Unrelated meaning: value of cosine = 0
  - Opposite meaning: value of cosine = -1

- Number of dimensions used is an important issue
  - Small dimensions (small singular values) represent local unique components
  - Large dimensions capture similarities and differences

*Center of Signal and Image Processing*
*Georgia Institute of Technology*

CSIP

# A Simple Binary Tree Classifier



C+
(Positive)

C-
(Negative)

● Root node

● Non-leaf node

● Leaf node

($X$: feature vector, $W$: parameters of the classifier)

$$f(X,W) = \sum_{i=1}^{D} w_i x_i - w_0$$

Decision rule:

$$\begin{cases} f(X,W) \geq 0, & \text{label C+} \\ \text{Otherwise} & \text{label C-} \end{cases}$$

ECE8813, Sprint 2009

*Center of Signal and Image Processing*
*Georgia Institute of Technology*

CSIP

# Multi-Class vs. Binary Decision Rule

- Multi-class (MC) classification

$$C(X) = \arg\max_j g_j(X;W), \quad 1 \le j \le m$$

$$g_j(X;W) > g_{i \ne j}(X;W) \quad X \in C_j$$

- Special case: Binary classifier with LDF *(C+: positive class, C-: negative class)*

$$\begin{cases} f(W,X) \ge 0 & \text{label C+} \\ \text{Others} & \text{label C-} \end{cases}$$

➢ **Decision rule is a *discrete, non-differential function of the classifier parameters (need MFoM to optimize)***

*Center of Signal and Image Processing*
*Georgia Institute of Technology*

CSIP

# Task and Experimental Setup

- *ModApte* split version of *Reuters-21578* corpus
  - lexicon: 10118 words, remove 319 stop-words and words occurred less than 4 times
  - corpus clean-up: remove documents which are not labeled by topics, miss topics, or are labeled by topics only occurred in training or test corpus
  - final experiments setup: 7,770 training documents, 3,019 test documents, 90 topics
  - some topics have little data for training or testing and with conflict labels in some cases

ECE8813, Sprint 2009

*Center of Signal and Image Processing*
*Georgia Institute of Technology*

CSIP

# Performance vs. LSI Feature Dimension



- MFoM Classifier performs better than the best SVM

ECE8813, Sprint 2009

*Center of Signal and Image Processing*
*Georgia Institute of Technology*

CSIP

# Experimental Results
## - Properties of MFoM Learning



**Figure 3. GPD convergence for category 'acq' (feature dimension: 400, X-axis: number of the iteration, Y-axis: $F_1$ measure for the positive class over training samples)**

ECE8813, Sprint 2009

*Center of Signal and Image Processing*
*Georgia Institute of Technology*

# Separation before and after MFoM (Gao, Wu and Lee, *SIGIR-2003*)

ECE8813, Sprint 2009

*Center of Signal and Image Processing*
*Georgia Institute of Technology*

# Performance Comparison (SIGIR2003)

|  | $k$-NN | SVM | Binary $F_1$-MFoM |
|---|---|---|---|
| micR | 0.834 | 0.812 | 0.857 |
| micP | 0.881 | 0.914 | 0.914 |
| micF$_1$ | 0.857 | 0.860 | 0.884 |
| macF$_1$ | 0.524 | 0.525 | 0.556 |

# Binary vs. MC TC (ICML04)

| Category | # of Training instances | Binary MFoM | MC MFoM |
|---|---|---|---|
| Income | 9 | 0.429 | 0.600 |
| Oat | 8 | 0.167 | 0.500 |
| Platinum | 5 | 0.286 | 0.833 |
| Potato | 3 | 0.333 | 0.750 |
| Sun-meal | 1 | 0.000 | 0.667 |

$F_1$ -based comparison:
Multi-Class MFoM works much better for small training sizes

ECE8813, Sprint 2009    *Center of Signal and Image Processing*
*Georgia Institute of Technology*

CSIP

# From Text to Multimedia Documents

- Property of raw multimedia patterns
  - Mostly fuzzy low-level signal representations
  - Hard to locate segmentation and object boundaries
- Definition of common sets of fundamental units
  - No obvious fundamental alphabets and words
  - Precision and coverage of multimedia tokenization
- Extraction of multimedia document feature vectors
  - Dimensionality, discrimination ability and trainability
- What are the missing links?
  - Shannon's information theory perspective (1951)
  - Finding acoustic, audio, visual "alphabets" and "words"

*Center of Signal and Image Processing*
*Georgia Institute of Technology*

CSIP

# Event Representation & Topic Classification

- Video: speech, audio, image, text, and others

**Speech** → [ ASR ] →

**Text** →

**Image** → [ AIA ] →

[ Morphological Filtering ] → [ Query-Vector Extraction ] → [ Text Categorization ] →

ASR: Automatic Speech Recognition
AIA: Automatic Image Annotation

CSIP

# Common Technology Thread: DSP, Feature Extraction & Classifier Learning

**Speech /Image /Audio** → **Media Tokenization** → **Text** → **LSA-Based FE/SVD** → **Feature** → **Audiovisual Classification** → **Results**

- Media Tokenization ← **A/V Alphabet Model**
- LSA-Based FE/SVD ← **A/V Word List**
- Audiovisual Classification ← **TC Classifier**

**Text Document Training Set** → **TC Classifier Learning** → **TC Classifier**

**First Step: Define alphabets and training alphabet models**

CSIP

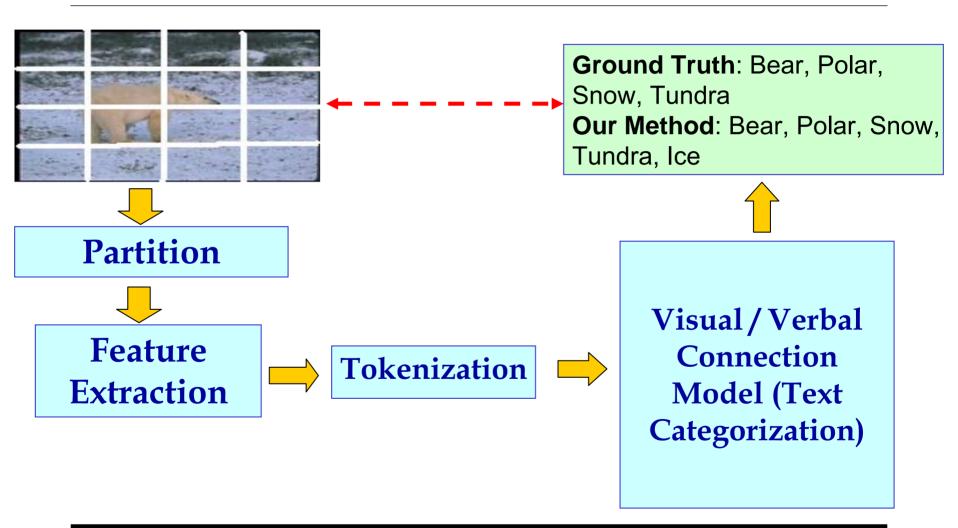# Automatic Image Annotation (AIA)

- A process associating concepts or keywords to images describing their visual content

- AIA can be used to make queries based on image concepts (Google-style keyword search)

Image/languag e connection

…,boat, sea, sky, beach,…

Verbal annotation

*Center of Signal and Image Processing*
*Georgia Institute of Technology*

CSIP

# Automatic Image Annotation



**Ground Truth**: Bear, Polar, Snow, Tundra
**Our Method**: Bear, Polar, Snow, Tundra, Ice

**Partition**

**Feature Extraction** → **Tokenization** → **Visual / Verbal Connection Model (Text Categorization)**

ECE8813, Sprint 2009

*Center of Signal and Image Processing*
*Georgia Institute of Technology*

CSIP

# Text Representation of Images

- Given a visual lexicon, $A=\{A_1, A_2, …, A_M\}$, with $M$ visual terms, an image document can be represented by $V=\{V_1, V_2, …, V_M\}$, each component being statistics of visual term occurred in the particular image document

- SVD can be applied to reduce the dimension, $M$

- Semantic concept modeling for image annotation

  – Semantic concept set, $C = \{C_j, \ 1 \le j \le N\}$, $N$: total concepts. Each concept has a discriminant function, $g_j(X; \Lambda_j)$, to be trained. Multiple relevant keywords are assigned to an image $X$, according to the rule,

*Center of Signal and Image Processing*
*Georgia Institute of Technology*

CSIP

# Music and Speech Connection

- Krishna and Sreenivas (2004) drew parallels between music and speech

    - Speech recognition ≈ music transcription

    - Instrument recognition ≈ speaker recognition

    - "Cocktail" separation ≈ instrument separation

    - Genre classification ≈ language classification

- Perceptual results do exist that give support to the link between music and language, but the debate is still continuing

CSIP

# Some references

- ## If you only read one article/reference:
  - Sebastiani, F. Machine learning in automated text categorization. ACM Computing Surveys, 34(1):1-47, 2002

- ## Worth having a look at:
  - Yang, Y. and Pedersen, J.O. A comparative study on feature selection in text categorization. In Proceedings of the 14th International Conference on Machine Learning, pages 412-420, 1997.
  - Dumais, S. and Chen, H. Hierarchical classification of web content. In Proceedings of the 23rd ACM SIGIR Conference, pages 256-263, 2000.
  - Lewis, D. D. An evaluation of phrasal and clustered representations on a text categorization task. In Proceedings of the 15th ACM SIGIR Conference, pages 37-50, 1992.

*Center of Signal and Image Processing*
*Georgia Institute of Technology*

CSIP

# Some References (Cont.)

- Mitchell, T. Machine learning. McGraw-Hill, 1997.
- Yang, Y. and Liu, X. A re-examination of text categorization methods. In Proceedings of ACM SIGIR, 1999.
- Lewis, D.D. Evaluating and optimizing autonomous text classification systems. In Proceedings of ACM SIGIR, 1995.
- Joachims, T. Text categorization with support vector machines: learning with many relevant features. In Proceedings of 10th European Conference on Machine Learning, pages 137-142, 1998.
- Hearst, M.A. Trends and discoveries: support vector macines. In IEEE Intelligent Systems, July/August 1998, pages 18-28.
- Yang, Y., Slattery, S., Ghani, R. A study of approaches to hypertext categorization. Journal of Intelligent Information Systems, 18(2/3):219-241, 2002.
- Gao, S., Wu, W., Chua, T.-S., Lee, C.-H. "A maximal figure-of-merit learning approach to trext categorization,". *Proc. of SIGIR*, 2003.

*Center of Signal and Image Processing*
*Georgia Institute of Technology*

CSIP

# Summary

- Today's Class
  - Text categorization

- Next Classes
  - Information retrieval
  - Labs 4-5 on PoS tagging and document clustering
  - Spring break: March 16-20, catch-up time
  - After break: IR, PCFG, probabilistic parsing

- Reading Assignments
  - Manning and Schutze, Chapters 14-16

*Center of Signal and Image Processing*
*Georgia Institute of Technology*

CSIP