# Machine Learning

CSE 454

---

## Today's Outline

- Brief supervised learning review
- Evaluation
- Overfitting
- Ensembles
    - Learners: The more the merrier
- Co-Training
    - (Semi) Supervised learning with few labeled training ex
- Clustering
    - No training examples

2

---

## Types of Learning

- **Supervised (inductive) learning**
    - Training data includes desired outputs
- **Semi-supervised learning**
    - Training data includes a few desired outputs
- **Unsupervised learning**
    - Training data does not include desired outputs
- **Reinforcement learning**
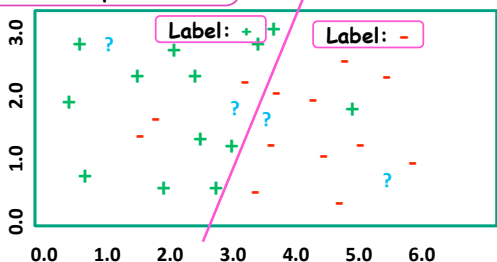    - Rewards from sequence of actions

---

## Supervised Learning

- **Inductive learning** or "Prediction":
    - **Given** examples of a function $(X, F(X))$
    - **Predict** function $F(X)$ for new examples $X$

- Classification
    - $F(X)$ = Discrete
- Regression
    - $F(X)$ = Continuous
- Probability estimation
    - $F(X)$ = Probability$(X)$:

4

---

1

## Classifier

Hypothesis:
Function for labeling examples



Label: **+**    Label: **-**

---

## Bias

- The nice word for prejudice is "bias".
- What kind of hypotheses will you consider?
  - What is allowable **range** of functions you use when approximating?
- What kind of hypotheses do you prefer?

- One idea: Prefer "simplest" hypothesis that is consistent with the data

---

## Naïve Bayes

- Probabilistic classifier:
  - $P(C_i \mid \text{Example})$
- Bias: Assumes all features are conditionally independent given class

$$P(E \mid c_i) = P(e_1 \wedge e_2 \wedge \cdots \wedge e_m \mid c_i) = \prod_{j=1}^{m} P(e_j \mid c_i)$$

- Therefore, we then only need to know $P(\mathbf{e_j} \mid \mathbf{c_i})$ for each feature and category
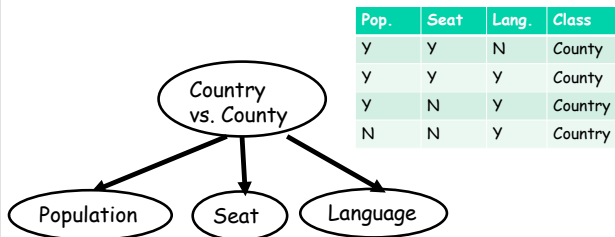
---

## Naïve Bayes for Text

- Modeled as generating a bag of words for a document in a given category
- Assumes that word order is unimportant, only cares whether a word appears in the document
- Smooth probability estimates with Laplace **m**-estimates assuming a uniform distribution over all words ($\mathbf{p} = 1/|\mathbf{V}|$) and $\mathbf{m} = |\mathbf{V}|$
  - Equivalent to a virtual sample of seeing each word in each category exactly once.

## Naïve Bayes

Country vs. County → Population, Seat, Language

| Pop. | Seat | Lang. | Class |
|------|------|-------|---------|
| Y | Y | N | County |
| Y | Y | Y | County |
| Y | N | Y | Country |
| N | N | Y | Country |

Probability(Seat | County)  = ??

Probability(Seat | Country) = ??

---

## Naïve Bayes

Country vs. County → Population, Seat, Language

| Pop. | Seat | Lang. | Class |
|------|------|-------|---------|
| Y | Y | N | County |
| Y | Y | Y | County |
| Y | N | Y | Country |
| N | N | Y | Country |

Probability(Seat | County)  = 2 + 1 / 2 + 1 = 1.0

Probability(Seat | Country) = ??

---

## Naïve Bayes

Country vs. County → Population, Seat, Language

| Pop. | Seat | Lang. | Class |
|------|------|-------|---------|
| Y | Y | N | County |
| Y | Y | Y | County |
| Y | N | Y | Country |
| N | N | Y | Country |

Probability(Seat | County)  = 2 + 1/ 2 + 2  = 0.75

Probability(Seat | Country) = 0 + 1 / 2 + 2 = 0.25

---

## Probabilities: Important Detail!

- $P(\text{spam} \mid E_1 \ldots E_n) = \prod_i P(\text{spam} \mid E_i)$

  **Any more potential problems here?**

- **We are multiplying lots of small numbers Danger of underflow!**
  - $0.5^{57} = 7\ E\ {-}18$

- **Solution? Use logs and add!**
  - $p_1 * p_2 = e^{\log(p1)+\log(p2)}$
  - **Always keep in log form**

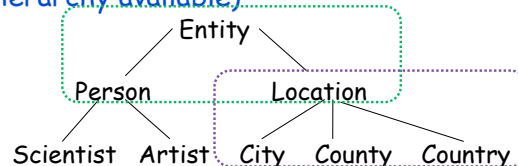## Multi-Class Categorization

- Pick the category with max probability
- Create many 1 vs other classifiers
    Classes = City, County, Country
    Classifier 1 = {City} {County, Country}
    Classifier 2 = {County} {City, Country}
    Classifier 3 = {Country} {City, County}

13

## Multi-Class Categorization

- Use a hierarchical approach (wherever hierarchy available)



14

## Today's Outline

- Brief supervised learning review
- Evaluation
- Overfitting
- Ensembles
    Learners: The more the merrier
- Co-Training
    (Semi) Supervised learning with few labeled training ex
- Clustering
    No training examples

15

## Experimental Evaluation

Question: How do we estimate the performance of classifier on unseen data?

- Can't just at accuracy on training data – this will yield an over optimistic estimate of performance

- Solution: Cross-validation

- Note: this is sometimes called estimating how well the classifier will generalize

16

## Evaluation: Cross Validation

- Partition examples into **k** disjoint sets
- Now create **k** training sets
  Each set is union of all equiv classes **except one**
  So each set has (k-1)/k of the original training data



---

## Cross-Validation (2)

- Leave-one-out
  Use if < 100 examples (rough estimate)
  Hold out one example, train on remaining examples

- 10-fold
  If have 100-1000's of examples

- M of N fold
  Repeat M times
  Divide data into N folds, do N fold cross-validation

---

## Today's Outline

- Brief supervised learning review
- Evaluation
- Overfitting
- Ensembles
  Learners: The more the merrier
- Co-Training
  (Semi) Supervised learning with few labeled training ex
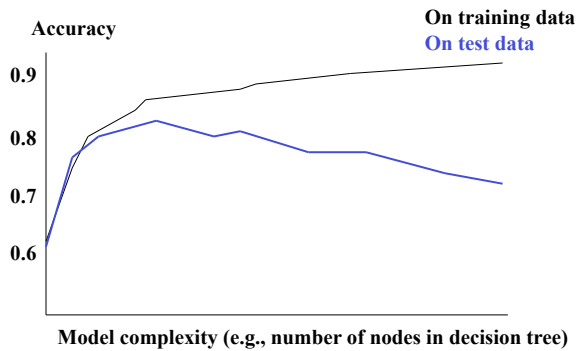- Clustering
  No training examples

19

---

## Overfitting Definition

- Hypothesis H is **overfit** when ∃ H' and
  H has **smaller** error on training examples, but
  H has **bigger** error on test examples
- Causes of overfitting
  Noisy data, or
  Training set is too small
  Large number of features
- Big problem in machine learning
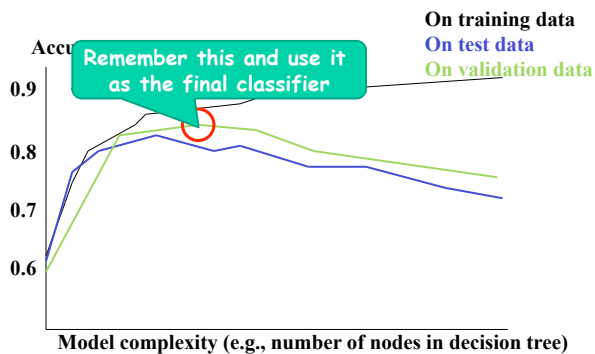- One solution: Validation set

# Overfitting

**Accuracy**

**On training data**
**On test data**

0.9
0.8
0.7
0.6

**Model complexity (e.g., number of nodes in decision tree)**

© Daniel S. Weld          21

# Validation/Tuning Set

- Split data into train and validation set

| | | | | | Tune | Tune | Tune | Test |

- Score each model on the tuning set, use it to pick the 'best' model

# Early Stopping

**On training data**
**On test data**
**On validation data**

**Accu**

**Remember this and use it as the final classifier**

0.9
0.8
0.7
0.6

**Model complexity (e.g., number of nodes in decision tree)**

© Daniel S. Weld          23

# Extra Credit Ideas

- Different types of models
  - Support Vector Machines (SVMs), widely used in web search
  - Tree-augmented naïve Bayes
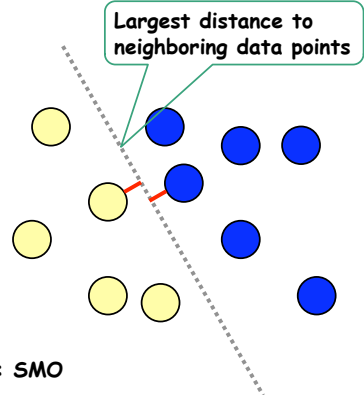- Feature construction

© Daniel S. Weld          24
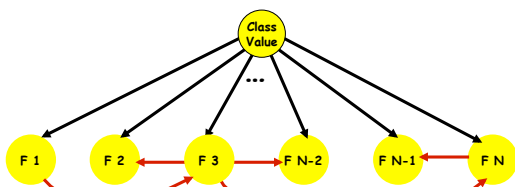
## Support Vector Machines



Which one is best hypothesis?

## Support Vector Machines



Largest distance to neighboring data points

**SVMs in Weka: SMO**

## Tree Augmented Naïve Bayes (TAN)
[Friedman,Geiger & Goldszmidt 1997]



**Models limited set of dependencies
Guaranteed to find best structure
Runs in polynomial time**

## Construct Better Features

- Key to machine learning is having good features

- In industrial data mining, large effort devoted to constructing appropriate features

- Ideas??

28

7

## Possible Feature Ideas

- Look at capitalization (may indicated a proper noun)

- Look for commonly occurring sequences
  - E.g. New York, New York City
  - Limit to 2-3 consecutive words
  - Keep all that meet minimum threshold (e.g. occur at least 5 or 10 times in corpus)

29

## Today's Outline

- Brief supervised learning review
- Evaluation
- Overfitting
- Ensembles
    Learners: The more the merrier
- Co-Training
    (Semi) Supervised learning with few labeled training ex
- Clustering
    No training examples

30

## Ensembles of Classifiers

- Traditional approach: Use one classifier
- Alternative approach: Use lots of classifiers
- Approaches:
  - Cross-validated committees
  - Bagging
  - Boosting
  - Stacking

31

## Voting



32

8

## Ensembles of Classifiers

- Assume
  - Errors are independent (suppose 30% error)
  - Majority vote
- Probability that majority is wrong…
  - = area under binomial distribution



Prob 0.2

0.1

*Ensemble of 21 classifiers*

Number of classifiers in error

- If individual area is 0.3
- **Area under curve for ≥11 wrong is 0.026**
- Order of magnitude improvement!

© Daniel S. Weld                                                                    33

## Constructing Ensembles
## Cross-validated committees

- Partition examples into **k** disjoint equiv classes
- Now create **k** training sets
  - Each set is union of all equiv classes **except one**
  - So each set has (k-1)/k of the original training data

- Now train a classifier on each set



Holdout

© Daniel S. Weld                                                                    34

## Ensemble Construction II
## Bagging

- Generate k sets of training examples
- For each set
  - Draw m examples randomly (with replacement)
  - From the original set of m examples
- Each training set corresponds to
  - 63.2% of original (+ duplicates)
- Now train classifier on each set
- Intuition: Sampling helps algorithm become more robust to noise/outliers in the data

© Daniel S. Weld                                                                    35
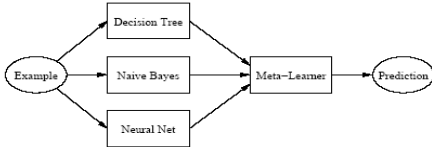
## Ensemble Creation III
## Boosting

- Maintain prob distribution over set of training ex
- Create k sets of training data iteratively:
- On iteration **i**
  - Draw m examples randomly (like bagging)
  - But use probability distribution to bias selection
  - Train classifier number **i** on this training set
  - Test partial ensemble (of **i** classifiers) on all training exs
  - Modify distribution: increase P of each error ex

- Create harder and harder learning problems...
- "Bagging with **optimized** choice of examples"

© Daniel S. Weld                                                                    36

9

# Ensemble Creation IV
## Stacking

- Train several base learners
- Next train meta-learner

> Learns when base learners are right / wrong
> Now meta learner arbitrates



Train using cross validated committees
- Meta-L inputs = base learner predictions
- Training examples = 'test set' from cross validation

---

# Today's Outline

- Overfitting
- Ensembles
> Learners: The more the merrier
- Co-Training
> Supervised learning with few labeled training ex
- Clustering
> No training examples

---

# Today's Outline

- Brief supervised learning review
- Evaluation
- Overfitting
- Ensembles
> Learners: The more the merrier
- Co-Training
> (Semi) Supervised learning with few labeled training ex
- Clustering
> No training examples

---

# Co-Training  Motivation

- Learning methods need labeled data
> Lots of <x, f(x)> pairs
> Hard to get… (who wants to label data?)

- But unlabeled data is usually plentiful…
> Could we use this instead??????

- Semi-supervised learning

10

## Co-training

Suppose

- Have **little** labeled data + **lots** of unlabeled

- Each instance has two parts:
  $x = [x1, x2]$
  $x1, x2$ conditionally independent given $f(x)$

- Each half can be used to classify instance
  $\exists f1, f2$ such that $f1(x1) \sim f2(x2) \sim f(x)$

- Both f1, f2 are learnable
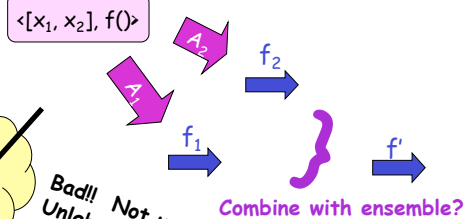  $f1 \in H1,\quad f2 \in H2,\quad \exists$ learning algorithms A1, A2

## Co-training Example

Prof. Domingos

Students: Parag,...

Projects: SRL, Data mining

I teach a class on data mining

CSE 546: Data Mining

Course Description:...

Topics:...

Homework: ...

Jesse

Classes taken:
1. Data mining
2. Machine learning

Research: SRL

## Without Co-training

$f_1(x_1) \sim f_2(x_2) \sim f(x)$

$A_1$ learns $f_1$ from $x_1$
$A_2$ learns $f_2$ from $x_2$

A **Few** Labeled Instances

$\langle[x_1, x_2], f()\rangle$

$A_2$

$f_2$

$A_1$

$f_1$

$[x_1, x_2]$

Bad!! Not using Unlabeled Instances!

Combine with ensemble?

$f'$

Unlabeled Instances

## Co-training

$f_1(x_1) \sim f_2(x_2) \sim f(x)$

$A_1$ learns $f_1$ from $x_1$
$A_2$ learns $f_2$ from $x_2$

A **Few** Labeled Instances

$\langle[x_1, x_2], f()\rangle$

$A_1$

$[x_1, x_2]$

$f_1$

$\langle[x_1, x_2], f_1(x_1)\rangle$

$A_2$

$f_2$

Hypothesis

Unlabeled Instances
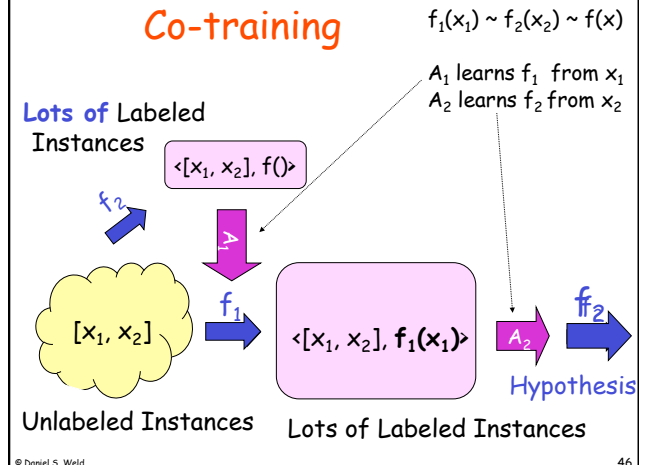
Lots of Labeled Instances

11

## Observations

- Can apply $A_1$ to generate as much training data as one wants
  - If $x_1$ is conditionally independent of $x_2$ / $f(x)$, then the error in the labels produced by $A_1$
    **will look like random noise to $A_2$ !!!**

- Thus **no limit** to quality of the hypothesis $A_2$ can make

45

---

## Co-training

$f_1(x_1) \sim f_2(x_2) \sim f(x)$

$A_1$ learns $f_1$ from $x_1$
$A_2$ learns $f_2$ from $x_2$



**Lots of** Labeled Instances

$\langle [x_1, x_2], f() \rangle$

$x_2$

$A_1$

$[x_1, x_2]$

$f_1$

$\langle [x_1, x_2], \mathbf{f_1(x_1)} \rangle$

$A_2$

$f_2$

Hypothesis

Unlabeled Instances          Lots of Labeled Instances

46

---

## It really works!

- Learning to classify web pages as course pages
  - x1 = bag of words on a page
  - x2 = bag of words from all anchors pointing to a page
- Naïve Bayes classifiers
  - 12 labeled pages
  - 1039 unlabeled

|                     | Page-based classifier | Hyperlink-based classifier | Combined classifier |
|---------------------|-----------------------|----------------------------|---------------------|
| Supervised training | 12.9                  | 12.4                       | 11.1                |
| Co-training         | 6.2                   | 11.6                       | 5.0                 |

Table 2: Error rate in percent for classifying web pages as course home pages. The top row shows errors when training on only the labeled examples. Bottom row shows errors when co-training, using both labeled and unlabeled examples.

47

---

## Today's Outline

- Brief supervised learning review
- Evaluation
- Overfitting
- Ensembles
  - Learners: The more the merrier
- Co-Training
  - (Semi) Supervised learning with few labeled training ex
- Clustering
  - No training examples

48

12

# Clustering Outline

- Motivation
- Document Clustering
- Offline evaluation
- Grouper I
- Grouper II
- Evaluation of deployed systems

# Low Quality of Web Searches

- System perspective:
   small coverage of Web (<16%)
   dead links and out of date pages
   limited resources
- IR perspective
  (relevancy of doc ~ similarity to query):
   very short queries
   huge database
   novice users

# Document Clustering

- User receives many (200 - 5000) documents from Web search engine

- Group documents in clusters
   by topic
- Present clusters as interface

# Grouper

GROUPER
A document clustering interface
for HuskySearch

Results from each engine: 50   Search for   All of these words

*www.cs.washington.edu/research/clustering*

GROUPER
Query: clinton

Documents: 298, Clusters: 15, Average Cluster Size: 16

| Cluster | Size | Shared Phrases and Sample Document Titles |
|---|---|---|
| 1 <br> View Results | 37 | Monica Lewinsky (32%), Clinton's scandals (16%), Kenneth Starr Investigation (14%), Hillary Clinton (14%) <br> ● Joke Post: Clinton Lewinsky Jokes <br> ● The Bill Clinton Information Gateway <br> ● Bill Clinton, Monica Lewinsky and Kenneth Starr – the saga of Bill and Monica. |
| 2 <br> View Results | 20 | Clinton a positive or negative (20%), Clinton/Gore (20%), Presidential Election (20%), election of (20%) <br> ● Republicans for Clinton <br> ● Clinton, Bill – Project Vote Smart <br> ● Clinton Record, The |
| 3 <br> View Results | 8 | Jones's (63%), documents (50%), special (50%); President (37%), Report (37%), legal (37%), Paula (37%) <br> ● Jones v. Clinton Special Report <br> ● Paula Jones Legal Fund <br> ● JONES vs CLINTON |

GROUPER
Query: clinton

Want to be more specific?
Use the phrases found to focus your search!

clinton

Results from each engine: 50   Search for   All of these words

☐ "Monica Lewinsky"          ☐ "Clinton's scandals"

☐ "Kenneth Starr Investigation"   ☐ "Hillary Clinton"

14

# Desiderata

- Coherent cluster
- Speed
- Browsable clusters
    Naming

# Main Questions

- Is document clustering feasible for Web search engines?

- Will the use of phrases help in achieving high quality clusters?

- Can phrase-based clustering be done quickly?

# 1. Clustering

group together similar items
(words or documents)

# Clustering Algorithms

- Hierarchical Agglomerative Clustering
  $O(n^2)$
- Linear-time algorithms
  K-means (Rocchio, 66)
  Single-Pass (Hill, 68)
  Fractionation (Cutting et al, 92)
  Buckshot (Cutting et al, 92)

# Basic Concepts - 1

- Hierarchical vs. Flat

# Basic Concepts - 2

- hard clustering:
  each item in only one cluster
- soft clustering:
  each item has a probability of membership in each cluster
- disjunctive / overlapping clustering:
  an item can be in more than one cluster

16

# Basic Concepts - 3

distance / similarity function
(for documents)

dot product of vectors
number of common terms
co-citations
access statistics
share common phrases

# Basic Concepts - 4

- What is "right" number of clusters?
  - apriori knowledge
  - default value: "5"
  - clusters up to 20% of collection size
  - choose best based on external criteria
  - Minimum Description Length
  - Global Quality Function
- no good answer

# K-means

- Works when we know k, the number of clusters

- Idea:
  Randomly pick k points as the "centroids" of the k clusters
  Loop:
  - ∀ points, add to cluster w/ nearest centroid
  - Recompute the cluster centroids
  - Repeat loop (until no change)

  **Iterative improvement of the objective function:
  Sum of the squared distance from each point
  to the centroid of its cluster**

# K-means Example

- For simplicity, 1-dimension objects and k=2.
  Numerical difference is used as the distance
- Objects: 1, 2,   5, 6,7
- K-means:
  Randomly select 5 and 6 as centroids;
  => Two clusters {1,2,5} and {6,7}; meanC1=8/3, meanC2=6.5
  => {1,2}, {5,6,7}; meanC1=1.5, meanC2=6
  => no change.
  Aggregate dissimilarity
  - (sum of squares of distanceeach point of each cluster from **its** cluster center--(intra-cluster distance)
    $= 0.5^2 + 0.5^2 + 1^2 + 0^2 + 1^2 = 2.5$
    $|1-1.5|^2$

## K Means Example (K=2)



Pick seeds
Reassign clusters
Compute centroids
Reasssign clusters
Compute centroids
Reassign clusters
Converged!

© Daniel S. Weld  **Slide from Rao Kambhampati**  [From Mooney] 69

---

Example of K-means in operation



Figure 9.5 Example of running the K-means algorithm on the two-dimensional antenna data. The plots show the locations of the means of the clusters (large circles) at various iterations of the K-means algorithm, as well as the classification of the data points at each iteration according to the closest mean (dots, circles, and xs for each of the three clusters).

© Daniel S. Weld  **Slide from Rao Kambhampati**  [From Hand et. Al.] 70

---

## Time Complexity

- Assume computing distance between two instances is $O(\mathbf{m})$ where $\mathbf{m}$ is the dimensionality of the vectors.
- Reassigning clusters: $O(\mathbf{kn})$ distance computations, or $O(\mathbf{knm})$.
- Computing centroids: Each instance vector gets added once to some centroid: $O(\mathbf{nm})$.
- Assume these two steps are each done once for $\mathbf{I}$ iterations: $O(\mathbf{Iknm})$.
- Linear in all relevant factors, assuming a fixed number of iterations,
    more efficient than $O(n^2)$ HAC (to come next)

© Daniel S. Weld  **Slide from Rao Kambhampati**  71

---

## Vector Quantization: K-means as Compression



FIGURE 14.9. Sir Ronald A. Fisher (1890-1962) was one of the founders of modern day statistics, to whom we owe maximum-likelihood, sufficiency, and many other fundamental concepts. The image on the left is a 1024×1024 grayscale image at 8 bits per pixel. The center image is the result of 2 × 2 block VQ, using 200 code vectors, with a compression rate of 1.9 bits/pixel. The right image uses only four code vectors, with a compression rate of 0.50 bits/pixel

**Slide from Rao Kambhampati** 72

18

## Problems with K-means

- Need to know k in advance
  - Could try out several k?
    - Cluster tightness increases with increasing K.
      - Look for a kink in the tightness vs. K curve
- Tends to go to local minima that are sensitive to the starting centroids
  - Try out multiple starting points
- Disjoint and exhaustive
  - Doesn't have a notion of "outliers"
    - Outlier problem can be handled by K-medoid or neighborhood-based algorithms
- Assumes clusters are spherical in vector space
  - Sensitive to coordinate changes, weighting etc.

**Why not the minimum value?**

**Example showing sensitivity to seeds**

A   B   C
O   O   O

O   O   O
D   E   F

**In the above, if you start with B and E as centroids you converge to {A,B,C} and {D,E,F}**
**If you start with D and F you converge to {A,B,D,E} {C,F}**

© Daniel S. Weld    **Slide from Rao Kambhampati**                                73

---

## Hierarchical Clustering

- Agglomerative
  bottom-up

Initialize: - **each item a cluster**
Iterate:     - **select two most similar clusters**
               – **merge them**
Halt:         **when have required # of clusters**

© Daniel S. Weld                                74

---

## Bottom Up Example



75

---

## Hierarchical Clustering

- Divisive
  top-bottom

Initialize:  -**all items one cluster**
Iterate:     - **select a cluster (least coherent)**
               - **divide it into two clusters**
Halt:        **when have required # of clusters**

© Daniel S. Weld                                76

19

## Top Down Example

## HAC Similarity Measures

- Single link
- Complete link
- Group average
- Ward's method

## Single Link

- cluster similarity = similarity of two most similar members

## Single Link

- $O(n^2)$
- chaining:



- bottom line:
  simple, fast
  often low quality

## Complete Link

- cluster similarity = similarity of two least similar members



81

## Complete Link

- worst case $O(n^3)$
- fast algo requires $O(n^2)$ space
- no chaining
- bottom line:
  - typically much faster than $O(n^3)$, often good quality

82

## Group Average

- cluster similarity

  = average similarity



83

## HAC Often Poor Results - Why?

- Often produces single large cluster
- Work best for:
  - spherical clusters; equal size; few outliers
- Text documents:
  - no model
  - not spherical; not equal size; overlap
- Web:
  - many outliers; lots of noise

84

21

## Example:  Clusters of Varied Sizes

k-means; complete-link; group-average:



single-link: chaining,
                but succeeds on this example

85

## Example - Outliers

HAC:



86

## Suffix Tree Clustering

### (KDD'97; SIGIR'98)

- Most clustering algorithms aren't **specialized** for text:
  Model document as **set** of words

- STC:
  document = **sequence** of words

87

## STC Characteristics

- Coherent
    phrase-based
    overlapping clusters
- Speed and Scalability
    linear time; incremental
- Browsable clusters
    phrase-based
    simple cluster definition

88

22

## STC - Central Idea

- Identify **base clusters**
    a group of documents that share a phrase
    use a **suffix tree**

- Merge base clusters as needed

89

## STC - Outline

Three logical steps:

1. "Clean" documents
2. Use a **suffix tree** to identify **base clusters** - a group of documents that share a phrase
3. Merge base clusters to form clusters

90

## Step 1 - Document "Cleaning"

- Identify sentence boundaries
- Remove
    HTML tags,
    JavaScript,
    Numbers,
    Punctuation

91

## Suffix Tree
### (Weiner, 73; Ukkonen, 95; Gusfield, 97)

Example - suffix tree of the string: (1) "**cats eat cheese**"



92

23

## Example - suffix tree of the strings:
 (1) "**cats eat cheese**",
 (2) "**mice eat cheese too**" and
 (3) "**cats eat mice too**"

---

## Step 2 - Identify Base Clusters via Suffix Tree

- Build one suffix tree from all sentences of all documents
- Suffix tree node = base cluster
- Score all nodes
- Traverse tree and collect top k (500) base clusters

---

## Step 3 - Merging Base Clusters

- Motivation:  similar documents share multiple phrases
- Merge base clusters based on the overlap of their document sets
- Example (query: "salsa")

"**tabasco sauce**"    docs: **3,4,5,6**
"**hot pepper**"    docs: **1,3,5,6**
"**dance**"    docs: **1,2,7**
"**latin music**"    docs: **1,7,8**

---

## Average Precision - WSR-SNIP



16% increase over k-means (not stat. sig.)

24

## Average Precision - WSR-DOCS



45% increase over k-means (stat. sig.)

## Grouper II

- Dynamic Index:
  Non-merged based clusters
- Multiple interfaces:
  List, Clusters + Dynamic Index (key phrases)
- Hierarchical:
  Interactive "Zoom In" feature
  (similar to Scatter/Gather)

---

**386 documents returned**
**Dynaimc Index:**

| | | |
|---|---|---|
| ❑ **clinton county** (8 docs) | ❑ **clinton crisis** (9 docs) | ❑ **clinton jokes** (15 docs) |
| ❑ **government executive branch clinton administration** (21 docs) | ❑ **hillary clinton** (22 docs) | ❑ **hillary rodham** (13 docs) |
| ❑ **impeach clinton** (9 docs) | ❑ **impeachment** (15 docs) | ❑ **iowa** (10 docs) |
| ❑ **kenneth starr investigation** (11 docs) | ❑ **law** (13 docs) | ❑ **lewinsky scandal** (8 docs) |
| ❑ **monica lewinsky** (11 docs) | ❑ **official** (10 docs) | ❑ **paula jones** (6 docs) |
| ❑ **photos** (6 docs) | ❑ **police department** (7 docs) | ❑ **political** (12 docs) |
| ❑ **port clinton** (9 docs) | ❑ **positive or negative** (7 docs) | ❑ **president** (56 docs) |
| ❑ **president clinton** (34 docs) | ❑ **white house** (7 docs) | ❑ **all others** (60 docs) |

**Mark enteries of interest above and select next display below**

⌃ Index   ⌄ Clusters   ⌄ Combined   ⌄ List   | Zoom In |   ⌐ download documents

| clinton | | New Query |

## Evaluation - Log Analysis

25

## Northern Light

- "Custom Folders"
- 20000 predefined topics in a manually developed hierarchy
- Classify document into topics
- Display "dominant" topics in search results

---

---

## Summary

- Post-retrieval clustering
  - to address low precision of Web searches
- STC
  - phrase-based; overlapping clusters; fast
- Offline evaluation
  - Quality of STC,
  - advantages of using phrases vs. n-grams, FS
- Deployed two systems on the Web
  - Log analysis: Promising initial results

  *www.cs.washington.edu/research/clustering*

---

## Cool Topic

- Internet allows creation of knowledge
  - Aka structured data
- Two dominant techniques
  - ML-based information extraction
    - Google scholar, product search
    - Zoominfo
    - Flipdog
  - Collaborative content authoring
    - Wikipedia
    - Summitpost
    - Amazon reviews (**and** votes on usefulness)
- How integrate the techniques?

26

## Integration Today

- ML first – creates a seed to attract users
- Humans act as proofreaders
    Zoominfo
    Zillow zestimates
    dblife.cs.wisc.edu

- Surely we can do better than this!?

---

## DBLife

---



**Daniel S. Weld**

---

## Total  >  Sum of Parts

- Human corrections
    - → training data
    - → improved ML output
- Active learning to prioritize corrections
- Track author (and ML extractor) reputations
    Learn policy where ML can overwrite human
- Insert javascript code to encourage human fact checking
- Realtime-ML to create "author helper"

## ESP Game

1 MILLION LABELS COLLECTED

**The ESP Game** beta

As seen on CNN and newspapers around the world!

**46 Players Logged In**

**TOP SCORES**

**HOW TO PLAY**

New to the ESP Game?
**Sign up for FREE!**

Already have an account?

Screen Name:

Password:

Sign In

Did you know?
The ESP Game is helping to label all images on the Web!
learn more...

**Play our new game**
NEW **Phetch** NEW

Terms of Service | FAQ | ESP Image Search | Contact Us | Credits

Funded in part by the National Science Foundation (NSF).
© 2005 Carnegie Mellon University, all rights reserved. Patent Pending.

© Daniel S. Weld                                                              109

---

http://www.espgame.org - The ESP Game - Mozilla Firefox

**2:12** Time Left

**The ESP Game**

**OOOO** score

**Taboo Words**
ORANGE
FRUIT
CITRUS
FOOD

**Your Guesses**
SUPERMARKET
YUMMY

**Type your next guess:**

Pass

Your partner has entered a guess

Flag

Applet PlayerClient started

© Daniel S. Weld                                                              110

---

## How Does this Fit In?

© Daniel S. Weld                                                              111

28