

- ### Tentative Schedule
- 11/1 Machine learning & datamining
 - 11/3 Text categorization & evaluation methods
 - 11/8 Information extraction
 - 11/10 KnowItAll
 - 11/15 ... continued
 - 11/17 Clustering & Focused crawling
 - 11/22 AJAX - Denise Draper
 - 11/24 ---
 - 11/29 Outbreak
 - 12/1 Cryptography / Security
 - 12/6 P2P & Advertising
 - 12/8 Semantic Web
- © Daniel S. Weld 3

- ### Today's Outline
- **Overfitting**
 - **Ensembles**
Learners: The more the merrier
 - **Co-Training**
Supervised learning with few labeled training ex
 - **Clustering**
No training examples
- © Daniel S. Weld 4

- ### 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?
- © Daniel S. Weld 5

- ### Learning = Function Approximation
- E.g., Checkers
V: boards -> evaluation
 - E.g., Handwriting recognition
V: image -> word
 - E.g., Mushrooms
V: mushroom-attributes -> {E, P}
 - OPINE ?
- © Daniel S. Weld 6

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):

Task
 Performance Measure
 Experience

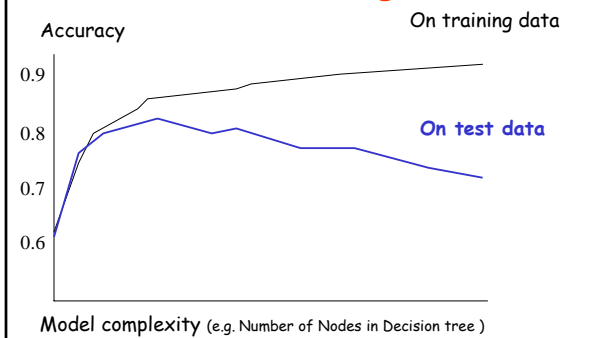
© Daniel S. Weld 7

(Some) Datamining Issues

- What feedback (experience) is available?
- How to represent this experience?
- How avoid overfitting?

© Daniel S. Weld 8

Overfitting



© Daniel S. Weld 9

Overfitting...

- Hypothesis H is *overfit* when $\exists H'$ 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
- **Huge problem in practice**
 Take class in ML or datamining...

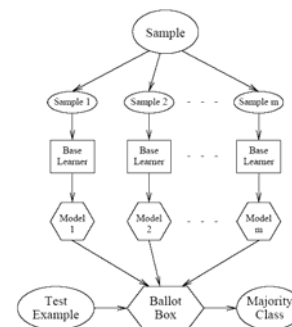
© Daniel S. Weld 10

Ensembles of Classifiers

Bagging
 Cross-validated committees
 Boosting
 Stacking

© Daniel S. Weld 11

Voting



© Daniel S. Weld 12

Ensembles of Classifiers

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

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

© Daniel S. Weld 13

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

© Daniel S. Weld 14

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

© Daniel S. Weld 15

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 16

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

© Daniel S. Weld 17

Co-Training Motivation

- Learning methods need labeled data
 - Lots of $\langle x, f(x) \rangle$ pairs
 - Hard to get... (who wants to label data?)
- But unlabeled data is usually plentiful...
 - Could we use this instead???????

© Daniel S. Weld 18

Co-training

Suppose

- Have *little* labeled data + *lots* of unlabeled
- Each instance has two parts:
 $x = [x_1, x_2]$
 x_1, x_2 conditionally independent given $f(x)$
- Each half can be used to classify instance
 $\exists f_1, f_2$ such that $f_1(x_1) \sim f_2(x_2) \sim f(x)$
- Both f_1, f_2 are learnable
 $f_1 \in H_1, f_2 \in H_2, \exists$ learning algorithms A_1, A_2

© Daniel S. Weld 19

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

Unlabeled Instances

© Daniel S. Weld 20

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

Unlabeled Instances

Lots of Labeled Instances

Hypothesis

© Daniel S. Weld 21

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

© Daniel S. Weld 22

It really works!

- Learning to classify web pages as course pages
 x_1 = bag of words on a page
 x_2 = 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	5.9

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.

© Daniel S. Weld 23

Choosing the Training Experience

- Credit assignment problem:
 - **Direct** training examples: Expensive!
 • E.g. individual checker boards + correct move for each
 • Supervised learning
 - **Indirect** training examples :
 • E.g. complete sequence of moves and final result
 • Reinforcement learning
 - **Unlabeled** training examples
 • Clustering
- Which examples:
 Random, teacher chooses, learner chooses

© Daniel S. Weld 24

Clustering Outline

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

© Daniel S. Weld

25

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

© Daniel S. Weld

26

Document Clustering

- User receives many (200 - 5000) documents from Web search engine
- Group documents in clusters by topic
- Present clusters as interface

© Daniel S. Weld

27

Grouper

www.cs.washington.edu/research/clustering

© Daniel S. Weld

28

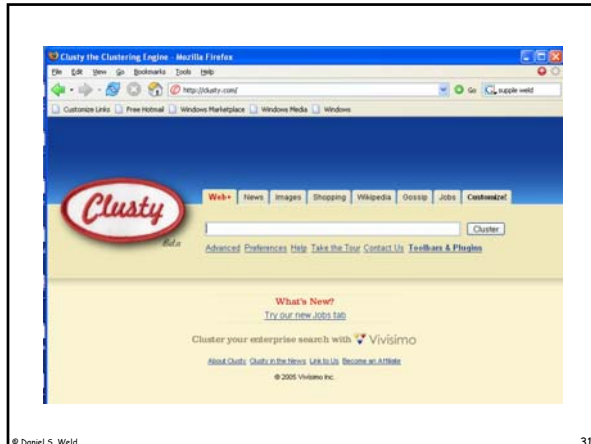
Cluster	Size	Shared Phrases and Sample Document Titles
1	37	Monica Lewinsky (32%), Clinton's scandals (16%), Kenneth Starr Investigation (14%), Hillary Clinton (14%) • Joke Post: Clinton Lewinsky Jokes • The Bill Clinton Information Gateway • Bill Clinton, Monica Lewinsky and Kenneth Starr - the saga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

© Daniel S. Weld

29

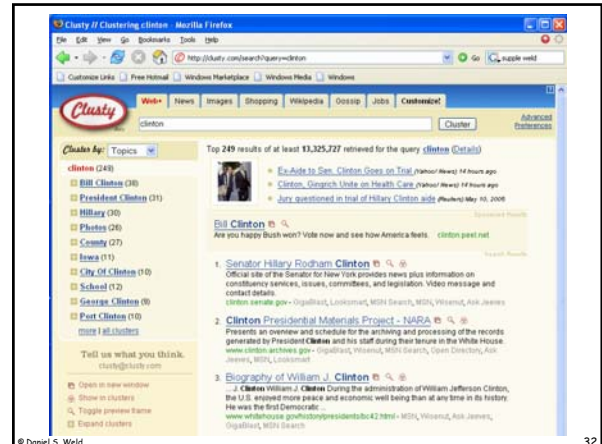
© Daniel S. Weld

30



© Daniel S. Weld

31



© Daniel S. Weld

32

Desiderata

- Coherent cluster
 - Speed
 - Browsible clusters
- Naming

© Daniel S. Weld

33

Main Questions

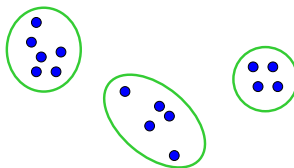
- Is document clustering feasible for Web search engines?
- Will the use of phrases help in achieving high quality clusters?
- Can phrase-based clustering be done quickly?

© Daniel S. Weld

34

1. Clustering

group together similar items
(words or documents)



© Daniel S. Weld

35

Clustering Algorithms

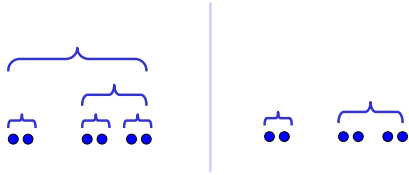
- Hierarchical Agglomerative Clustering
 $O(n^2)$
- Linear-time algorithms
 - K-means (Rocchio, 66)
 - Single-Pass (Hill, 68)
 - Fractionation (Cutting et al, 92)
 - Buckshot (Cutting et al, 92)

© Daniel S. Weld

36

Basic Concepts - 1

- Hierarchical vs. Flat



© Daniel S. Weld

37

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

© Daniel S. Weld

38

Basic Concepts - 3

distance / similarity function
(for documents)

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

© Daniel S. Weld

39

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**

© Daniel S. Weld

40

K-means

- Works when we know k , the number of clusters
- **Idea:**
Randomly pick k points as the "centroids" of the k clusters
- **Loop:**
 - \forall 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

© Daniel S. Weld

Slide from Rao Kambhampati

41

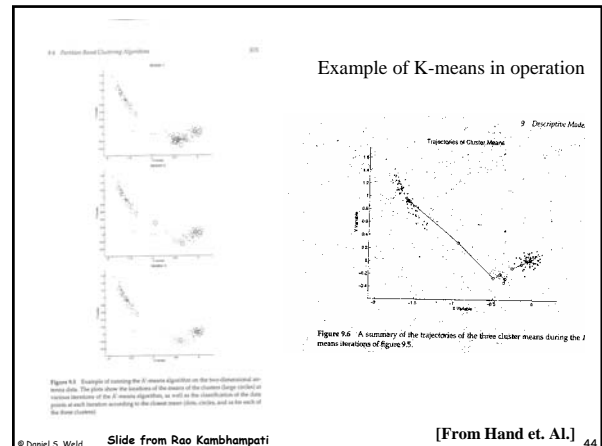
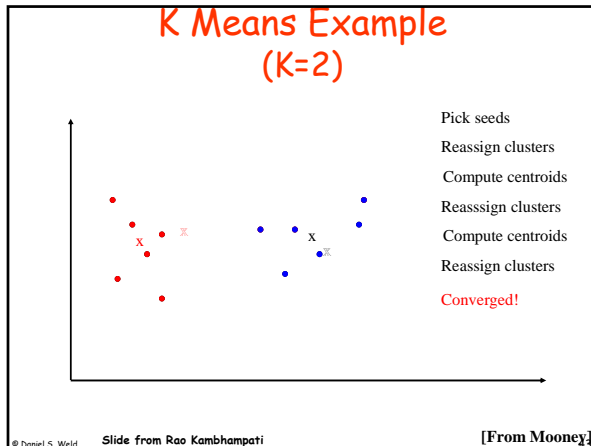
K-means Example

- 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 distance each 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

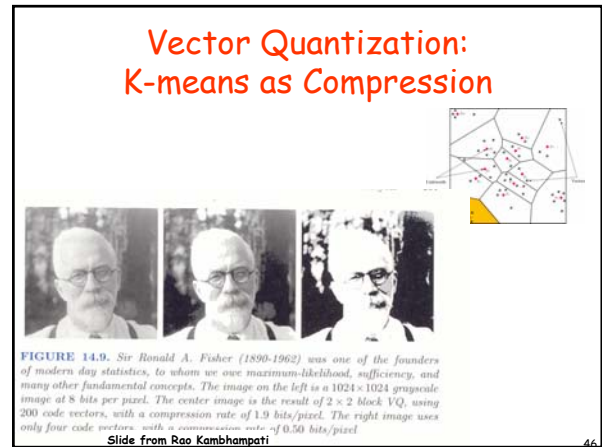
© Daniel S. Weld

Slide from Rao Kambhampati

42



- ## Time Complexity
- Assume computing distance between two instances is $O(m)$ where m is the dimensionality of the vectors.
 - Reassigning clusters: $O(kn)$ distance computations, or $O(knm)$.
 - Computing centroids: Each instance vector gets added once to some centroid: $O(nm)$.
 - Assume these two steps are each done once for I iterations: $O(Iknm)$.
 - 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 45



- ## 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
○	○	○
○	○	○
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 47

- ## 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 48

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

49

HAC Similarity Measures

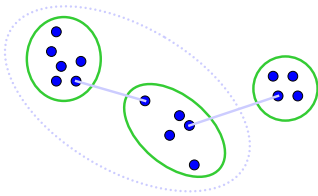
- Single link
- Complete link
- Group average
- Ward's method

© Daniel S. Weld

50

Single Link

- cluster similarity = similarity of two most similar members

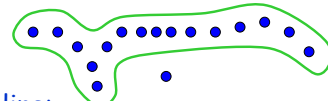


© Daniel S. Weld

51

Single Link

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



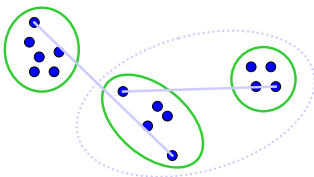
- bottom line:
simple, fast
often low quality

© Daniel S. Weld

52

Complete Link

- cluster similarity = similarity of two least similar members



© Daniel S. Weld

53

Complete Link

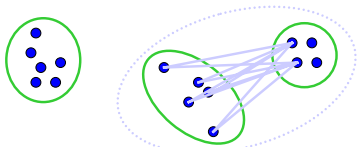
- worst case $O(n^3)$
- fast algo requires $O(n^2)$ space
- no chaining
- bottom line:
typically much faster than $O(n^3)$,
often good quality

© Daniel S. Weld

54

Group Average

- cluster similarity
= average similarity of all pairs



© Daniel S. Weld

55

HAC Often Poor Results - Why?

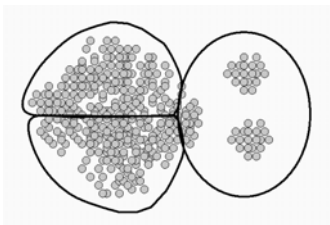
- Often produces single large cluster
- Work best for:
 - spherical clusters; equal size; few outliers
- Text documents:
 - no model
 - not spherical; not equal size; overlap
- Web:
 - many outliers; lots of noise

© Daniel S. Weld

56

Example: Clusters of Varied Sizes

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



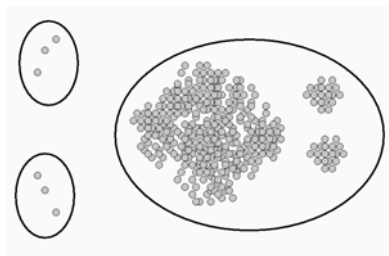
single-link: chaining,
but succeeds on this example

© Daniel S. Weld

57

Example - Outliers

HAC:



© Daniel S. Weld

58

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

© Daniel S. Weld

59

STC Characteristics

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

© Daniel S. Weld

60

STC - Central Idea

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

© Daniel S. Weld

61

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

© Daniel S. Weld

62

Step 1 - Document "Cleaning"

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

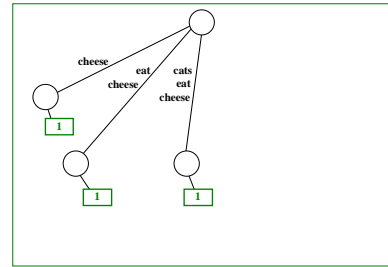
© Daniel S. Weld

63

Suffix Tree

(Weiner, 73; Ukkonen, 95; Gusfield, 97)

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

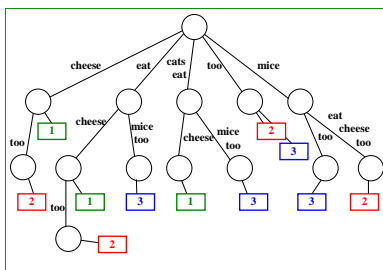


© Daniel S. Weld

64

Example - suffix tree of the strings:

- (1) "cats eat cheese",
- (2) "mice eat cheese too" and
- (3) "cats eat mice too"



© Daniel S. Weld

65

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

© Daniel S. Weld

66

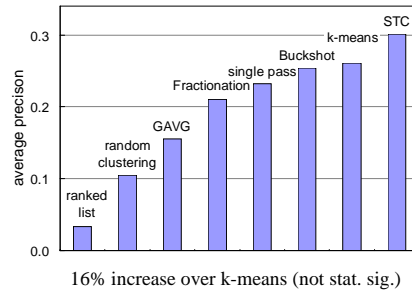
Step 3 - Merging Base Clusters

- 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

© Daniel S. Weld

67

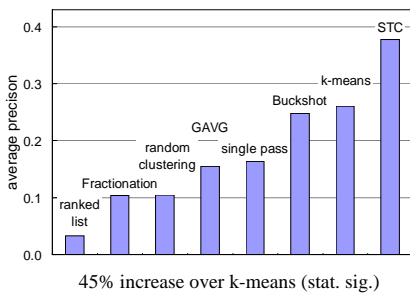
Average Precision - WSR-SNIP



© Daniel S. Weld

68

Average Precision - WSR-DOCS



© Daniel S. Weld

69

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)

© Daniel S. Weld

70

386 documents returned
Dynamic Index:

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

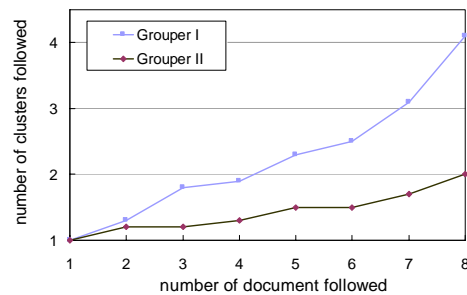
Mark entries of interest above and select next display below

Index Clusters Combined List Zoom In download documents

© Daniel S. Weld

71

Evaluation - Log Analysis



© Daniel S. Weld

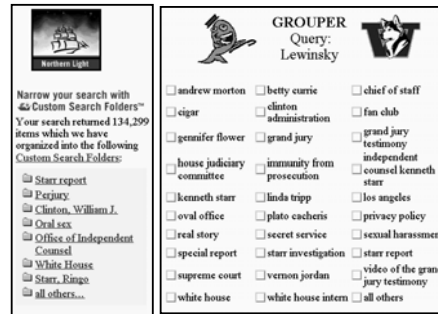
72

Northern Light

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

© Daniel S. Weld

73



The screenshot shows the Northern Light search interface. On the left, there is a search box with the text "Northern Light" and a search button. Below the search box, it says "Narrow your search with Custom Search Folders™" and "Your search returned 134,299 items which we have organized into the following Custom Search Folders:". A list of folders is shown, including "Starr report", "Perjury", "Clinton, William J.", "Oral sex", "Office of Independent Counsel", "White House", "Starr, Einge", and "all others...". On the right, there is a "GROUPER" section with the query "Lewinsky" and a list of related terms, each with a checkbox. The terms include "andrew morton", "betty curie", "chief of staff", "cigar", "clinton administration", "fm club", "geniafer flower", "grand jury", "grand jury testimony", "house judiciary committee", "immunity from prosecution", "independent counsel kenneth starr", "kenneth starr", "linda tripp", "los angeles", "oval office", "plato eacheris", "privacy policy", "real story", "secret service", "sexual harassment", "special report", "starr investigation", "starr report", "supreme court", "vernon jordan", "video of the grand jury testimony", "white house", "white house intern", and "all others".

© Daniel S. Weld

74

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

© Daniel S. Weld

75