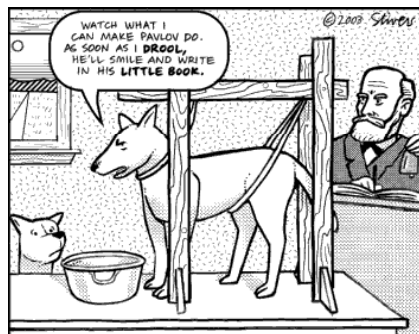


Machine Learning



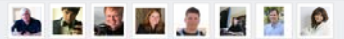
CSE 454

Search Engine News

blekko beta

examples: cure for headaches | global warming liberal

check out my slashtags:



slashtag chatter find slashtags what is a slashtag

slash the web

watch the demo video

blekko | machine learning

1 to 20 of 30M web results for machine learning

1 Machine learning - Wikipedia
tag | edit links cache | chatter | spam
Machine learning is a scientific discipline that is concerned with the design and development of algorithms that allow computers to evolve behaviors based on empirical data.
en.wikipedia.org/wiki/Machine_learning

2 UCI Machine Learning Repository
tag | edit links cache | chatter | spam
A repository of databases, domain theories and data generators that are used by the machine learning community for the empirical analysis of machine learning algorithms at ics.uci.edu/~mlearn/ Repository.html

3 Machine Learning textbook
tag | edit links cache | chatter | spam
A textbook by Tom Mitchell. McGraw Hill, 1997. Machine Learning, Tom Mitchell, McGraw Hill, 1997. Machine Learning is the study of computer algorithms that improve automatically cs.cmu.edu/~tom/mlbook.html

4 Journal of Machine Learning Research Homepage
tag | edit links cache | chatter | spam
Publishes scholarly articles in all areas of machine learning electronically. Paper volume release 6 times annually. Website contains content in different file formats, information and jmlr.csail.mit.edu

5 The International Machine Learning Society - About

blekko | machine learning

1 to 20 of 30M web results for machine learning

1 Machine learning - Wikipedia
tag | edit links cache | chatter | spam
Machine learning is a scientific discipline that is concerned with the design and development of algorithms that allow computers to evolve behaviors based on empirical data.
en.wikipedia.org/wiki/Machine_learning

2 UCI Machine Learning Repository
tag | edit links cache | chatter | spam
A repository of databases, domain theories and data generators that are used by the machine learning community for the empirical analysis of machine learning algorithms at ics.uci.edu/~mlearn/ Repository.html

3 Machine Learning textbook
tag | edit links cache | chatter | spam
A textbook by Tom Mitchell. McGraw Hill, 1997. Machine Learning, Tom Mitchell, McGraw Hill, 1997. Machine Learning is the study of computer algorithms that improve automatically cs.cmu.edu/~tom/mlbook.html

4 Journal of Machine Learning Research Homepage
tag | edit links cache | chatter | spam
Publishes scholarly articles in all areas of machine learning electronically. Paper volume release 6 times annually. Website contains content in different file formats, information and jmlr.csail.mit.edu

5 The International Machine Learning Society - About

blekko | machine learning

1 to 20 of 30M web results for machine learning

1 Machine learning - Wikipedia
tag | edit links cache | chatter | spam
Machine learning is a scientific discipline that is concerned with the design and development of algorithms that allow computers to evolve behaviors based on empirical data.
en.wikipedia.org/wiki/Machine_learning

2 UCI Machine Learning Repository
tag | edit links cache | chatter | spam
A repository of databases, domain theories and data generators that are used by the machine learning community for the empirical analysis of machine learning algorithms at ics.uci.edu/~mlearn/ Repository.html

3 Machine Learning textbook
tag | edit links cache | chatter | spam
A textbook by Tom Mitchell. McGraw Hill, 1997. Machine Learning, Tom Mitchell, McGraw Hill, 1997. Machine Learning is the study of computer algorithms that improve automatically cs.cmu.edu/~tom/mlbook.html

4 Journal of Machine Learning Research Homepage
tag | edit links cache | chatter | spam
Publishes scholarly articles in all areas of machine learning electronically. Paper volume release 6 times annually. Website contains content in different file formats, information and jmlr.csail.mit.edu

5 The International Machine Learning Society - About

Slashtags

blekko | machine learning

1 to 20 of 30M web results for machine learning

1 Machine learning - Wikipedia
tag | edit links cache | chatter | spam
Machine learning is a scientific discipline that is concerned with the design and development of algorithms that allow computers to evolve behaviors based on empirical data.
en.wikipedia.org/wiki/Machine_learning

2 UCI Machine Learning Repository
tag | edit links cache | chatter | spam
A repository of databases, domain theories and data generators that are used by the machine learning community for the empirical analysis of machine learning algorithms at ics.uci.edu/~mlearn/ Repository.html

3 Machine Learning textbook
tag | edit links cache | chatter | spam
A textbook by Tom Mitchell. McGraw Hill, 1997. Machine Learning, Tom Mitchell, McGraw Hill, 1997. Machine Learning is the study of computer algorithms that improve automatically cs.cmu.edu/~tom/mlbook.html

4 Journal of Machine Learning Research Homepage
tag | edit links cache | chatter | spam
Publishes scholarly articles in all areas of machine learning electronically. Paper volume release 6 times annually. Website contains content in different file formats, information and jmlr.csail.mit.edu

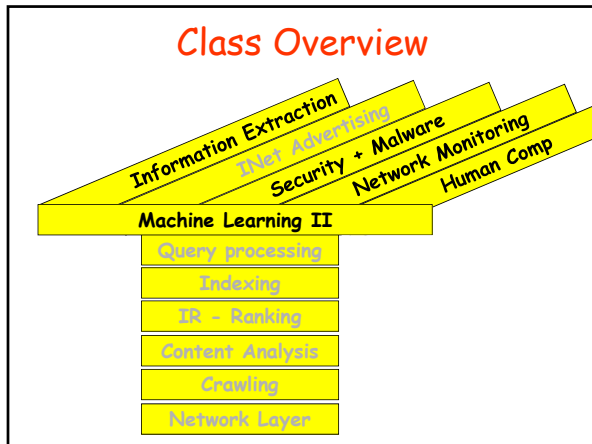
5 The International Machine Learning Society - About

Slashtag these results

Bundles the websites for these search results into a single slashtag

Google Local Search

- Place search
 - Finding local business w/ maps, reviews, etc
- Boost
 - Ads for local businesses
- 20% searches are related to location



- ## Next few classes
- Machine Learning
 - Malware (Arvind Krishnamurthy)
 - Information Extraction
 - " " Continued
 - NLP Basics: Parsing & POS Tagging
 - Internet-Enabled Human Computation

- ## 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

Sample Category Learning Problem

- Instance language: $\langle \text{size, color, shape} \rangle$
 - size $\in \{\text{small, medium, large}\}$
 - color $\in \{\text{red, blue, green}\}$
 - shape $\in \{\text{square, circle, triangle}\}$
- $C = \{\text{positive, negative}\}$
- D :

Example	Size	Color	Shape	Category
1	small	red	circle	positive
2	large	red	circle	positive
3	small	red	triangle	negative
4	large	blue	circle	negative

Example: County vs. Country?

- **Given:**
 - A description of an instance, $x \in X$, where X is the instance language or instance space.
 - A fixed set of categories $C = \{c_1, c_2, \dots, c_n\}$
- **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 .

Bag of words representation

Learning for Categorization

- A **training example** is an instance $x \in X$, paired with its correct category $c(x)$: $\langle x, c(x) \rangle$ for an unknown categorization function, c .
- Given a set of training examples, D .

$\{ \langle \text{img}, \text{county} \rangle, \langle \text{img}, \text{country} \rangle, \dots \}$

- Find a hypothesized categorization function, $h(x)$, such that $\forall x \in D: h(x) = c(x)$

Consistency

Generalization

- Hypotheses must **generalize** to correctly classify instances not in the training data.
- Simply memorizing training examples is a consistent hypothesis **that does not generalize**.

19

Why is Learning Possible?

Experience alone never justifies any conclusion about any unseen instance.

Learning occurs when
PREJUDICE meets **DATA!**

Learning a "Frobnitz"

© Daniel S. Weld

20

Bias

- Which hypotheses **will you consider?**
- Which hypotheses do you **prefer?**

© Daniel S. Weld

21

Some Typical Biases

Occam's razor

"It is needless to do more when less will suffice"

- William of Occam,

died 1349 of the Black plague

MDL - Minimum description length

Concepts can be approximated by

... **conjunctions** of predicates

... by **linear** functions

... by **short** decision trees

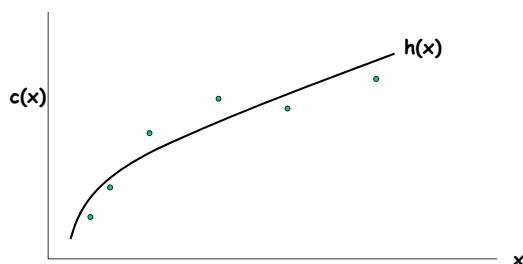
Frobnitz?

© Daniel S. Weld

22

ML = Function Approximation

May not be any perfect fit
Classification ~ discrete functions



23

Supervised Learning

- **Inductive learning** or "Prediction":
Given examples of a function $(X, F(X))$
Predict function $F(X)$ for new examples X
- **Classification**
 $F(X) = \text{Discrete}$
- **Regression**
 $F(X) = \text{Continuous}$
- **Probability estimation**
 $F(X) = \text{Probability}(X)$

© Daniel S. Weld

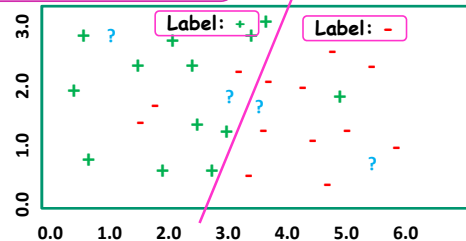
24

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 *doesn't* include desired outputs
- **Reinforcement learning**
Rewards from sequence of actions

Classifier

Hypothesis:
Function for labeling
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

© Daniel S. Weld

27

Experimental Evaluation

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

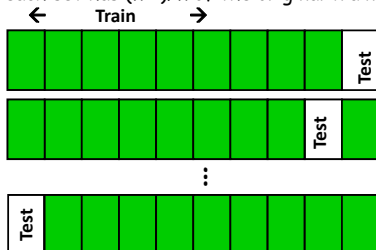
- Can't just at accuracy on training data - this will yield an over optimistic estimate of performance
- Solution: Cross-validation
- Note: this is sometimes called estimating how well the classifier will generalize

© Daniel S. Weld

28

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

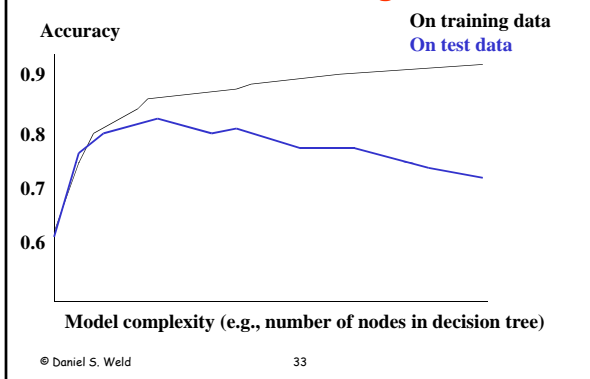
© Daniel S. Weld

31

Overfitting Definition

- 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
 - Large number of features
- Big problem in machine learning
- One solution: Validation set

Overfitting



33

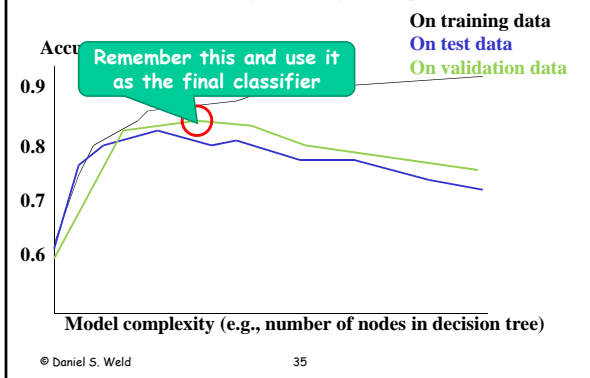
Validation/Tuning Set

- Split data into train and validation set



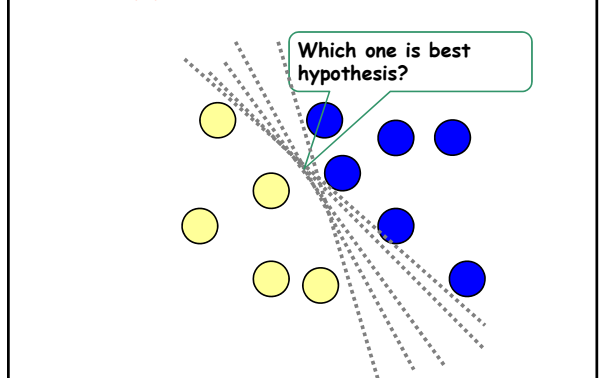
- Score each model on the tuning set, use it to pick the 'best' model

Early Stopping

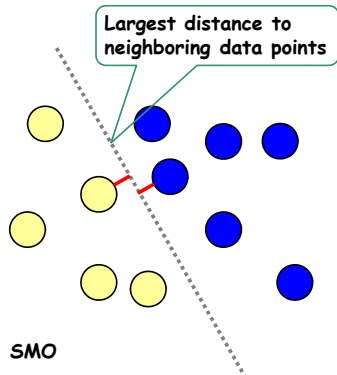


35

Support Vector Machines



Support Vector Machines



SVMs in Weka: SMO

Construct Better Features

- Key to machine learning is having good features
- In industrial data mining, large effort devoted to constructing appropriate features
- Ideas??

© Daniel S. Weld

38

Possible Feature Ideas

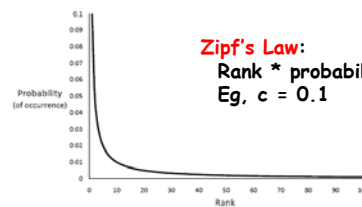
- Look at capitalization (may indicated a proper noun)
- Look for commonly occurring sequences
 - E.g. New York, New York City
 - Limit to 2-3 consecutive words
 - Keep all that meet minimum threshold (e.g. occur at least 5 or 10 times in corpus)

© Daniel S. Weld

39

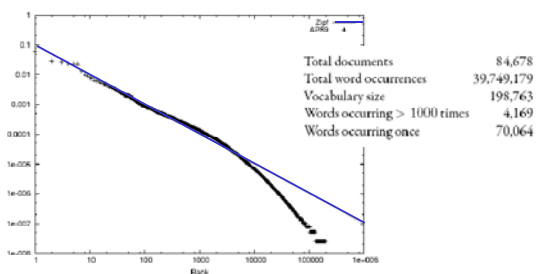
Properties of Text

- Word frequencies - skewed distribution
- 'The' and 'of' account for 10% of all words
- Six most common words account for 40%



From [Croft, Metzler & Strohman 2010]

Associate Press Corpus `AP89'



From [Croft, Metzler & Strohman 2010]

Middle Ground

- Very common words → bad features
- Language-based **stop list**:
 - words that bear little meaning
 - 20-500 words
 - http://www.dcs.gla.ac.uk/idom/ir_resources/linguistic_utils/stop_words
- Subject-dependent stop lists
- Very rare words *also* bad features
 - Drop words appearing less than k times / corpus

Stop lists

- **Language-based stop list:**
words that bear little meaning
20-500 words
http://www.dcs.gla.ac.uk/idom/ir_resources/linguistic_utils/stop_words
- **Subject-dependent stop lists**

From Peter Brusilovsky Univ Pittsburg INFSCI 2140

43

Stemming

- **Are there different index terms?**
retrieve, retrieving, retrieval, retrieved, retrieves...
- **Stemming algorithm:**
(retrieve, retrieving, retrieval, retrieved, retrieves) \Rightarrow **retriev**
Strips prefixes of suffixes (-s, -ed, -ly, -ness)
Morphological stemming

Copyright © Weld 2002-2007

44

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

© Daniel S. Weld

45

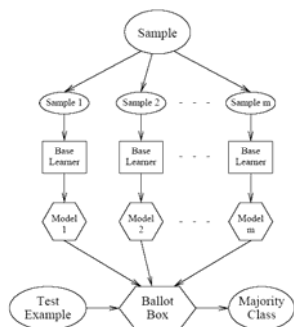
Ensembles of Classifiers

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

© Daniel S. Weld

46

Voting



© Daniel S. Weld

47

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

48

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

49

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

50

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"

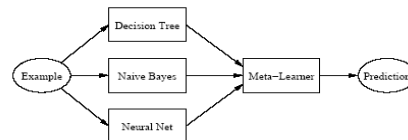
© Daniel S. Weld

51

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

52

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

© Daniel S. Weld

53

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 *doesn't* include desired outputs
- **Reinforcement learning**
 - Rewards from sequence of actions

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??????
- Semi-supervised learning

© Daniel S. Weld

55

Co-training

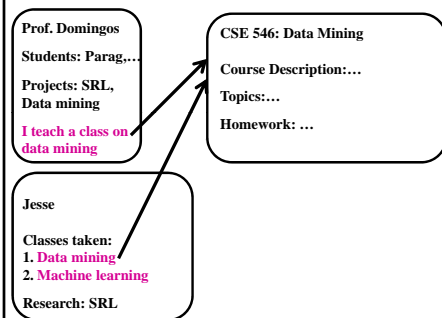
Suppose

- 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

56

Co-training Example



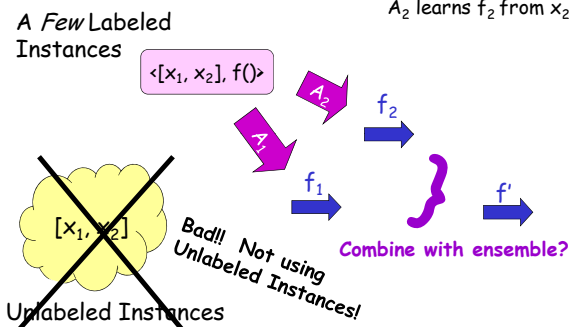
© Daniel S. Weld

57

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



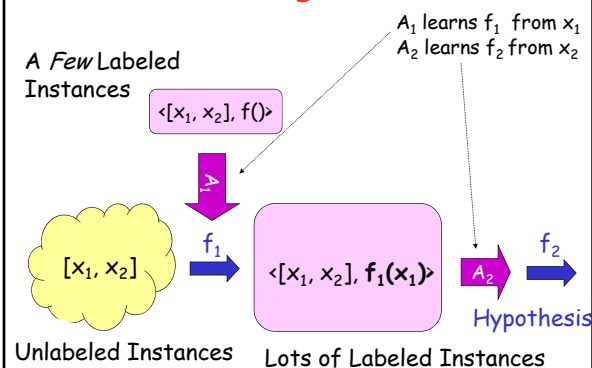
© Daniel S. Weld

58

Co-training

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

A Few Labeled Instances



© Daniel S. Weld

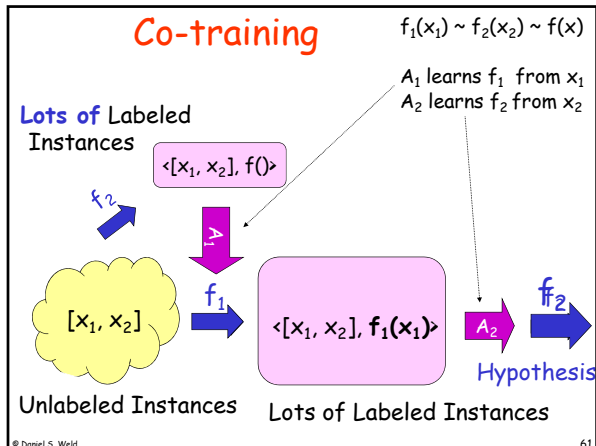
59

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

60



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	11.5	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 62

- ### 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 *doesn't* include desired outputs
 - **Reinforcement learning**
Rewards from sequence of actions
- © Daniel S. Weld 63

Learning with Hidden Labels

- Expectation Maximization Algorithm

© Daniel S. Weld 64

