

Transductive Learning

**CS4780 – Machine Learning
Fall 2009**

**Thorsten Joachims
Cornell University**

Outline

- **Transductive Learning Setting**
- **Transduction via Graph Cuts**
 - Minimum s-t-cuts
 - Minimum ratio cuts
- **Transductive Support Vector Machines**
- **Co-Training**

Transductive Learning Process

Sampling Training data

- select random subset of l examples from DB of size n

$$\Rightarrow Z = [x_1, \dots, x_l]$$

- receive labels for these examples positive (+1) / negative (-1)

$$\Rightarrow Z = [(x_1, y_1), \dots, (x_l, y_l)]$$

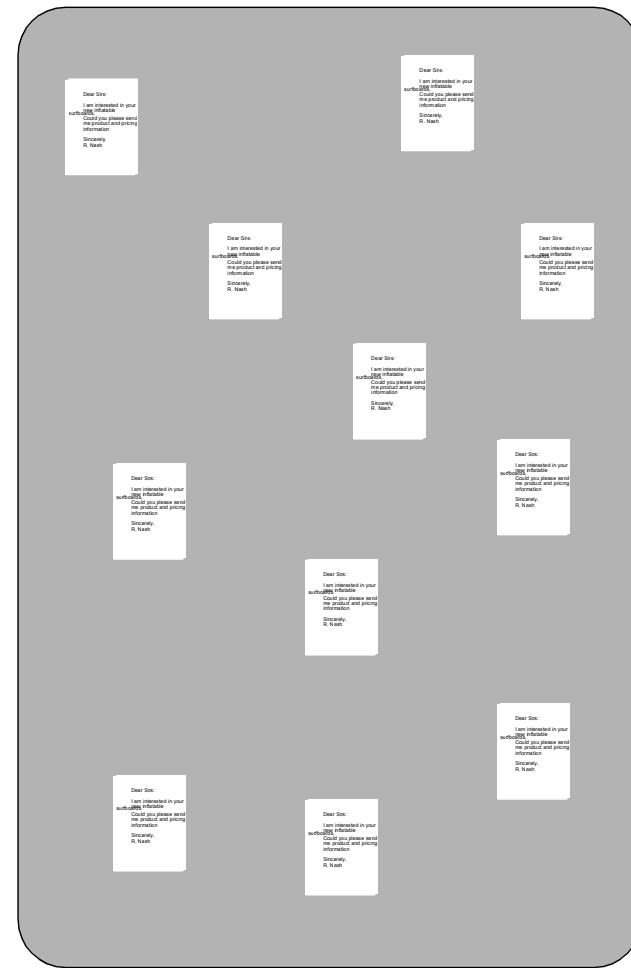
Goal of Learner

- predict the labels of the remaining examples $Z_x^* = [x_1^*, \dots, x_k^*]$

Opportunity

- Learning algorithm can study the test examples $Z_x^* = [x_1^*, \dots, x_k^*]$

Document DB



Transductive Learning Process

Sampling Training data

- select random subset of l examples from DB of size n

$$\Rightarrow Z = [x_1, \dots, x_l]$$

- receive labels for these examples positive (+1) / negative (-1)

$$\Rightarrow Z = [(x_1, y_1), \dots, (x_l, y_l)]$$

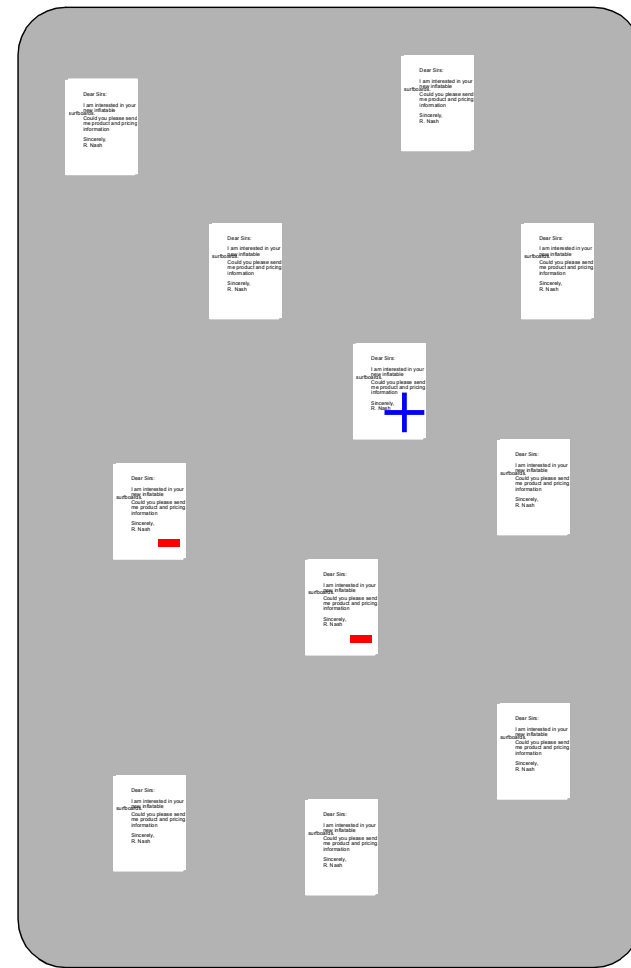
Goal of Learner

- predict the labels of the remaining examples $Z_x^* = [x_1^*, \dots, x_k^*]$

Opportunity

- Learning algorithm can study the test examples $Z_x^* = [x_1^*, \dots, x_k^*]$

Document DB



Transductive Learning Process

Sampling Training data

- select random subset of l examples from DB of size n

$$\Rightarrow Z = [x_1, \dots, x_l]$$

- receive labels for these examples positive (+1) / negative (-1)

$$\Rightarrow Z = [(x_1, y_1), \dots, (x_l, y_l)]$$

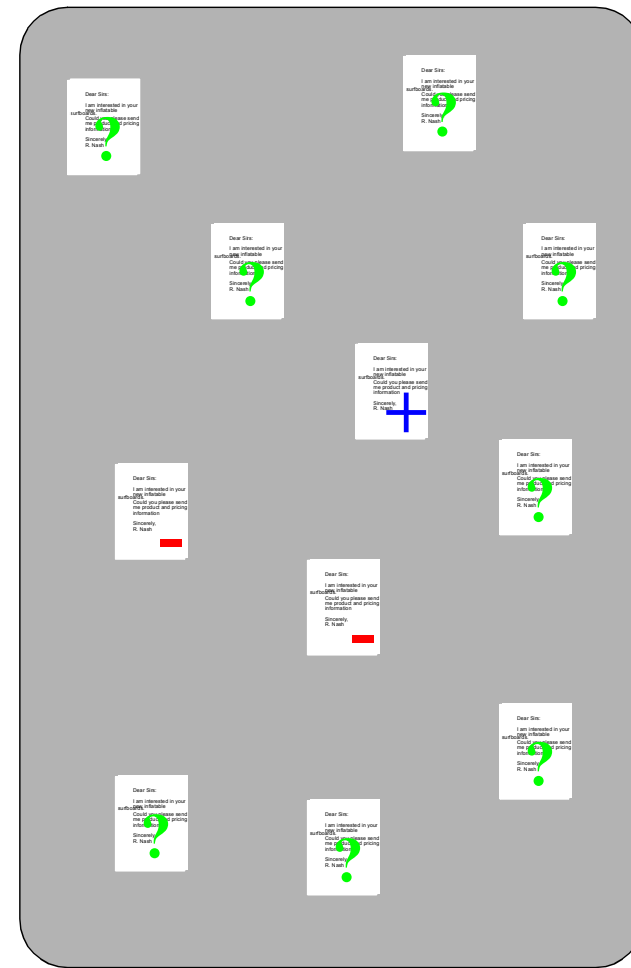
Goal of Learner

- predict the labels of the remaining examples $Z_x^* = [x_1^*, \dots, x_k^*]$

Opportunity

- Learning algorithm can study the test examples $Z_x^* = [x_1^*, \dots, x_k^*]$

Document DB



Example: Exploiting the Test Set

How would you classify the test set for Term/document matrix A_i

	nuclear	physics	atom	pepper	basil	salt	and
+	D1	1					1
?	D2	1	1	1			1
?	D3			1			1
?	D4			1	1		1
?	D5			1		1	1
-	D6				1	1	1

[Joachims, 1999]

- training set {D1, D6}
- test set {D2, D3, D4, D5}

Example: Exploiting the Test Set

How would you classify the test set for Term/document matrix A_i ?

	nuclear	physics	atom	pepper	basil	salt	and
+	D1	1					1
?	D2	1	1	1			1
?	D3			1			1
?	D4				1	1	1
?	D5				1	1	1
-	D6					1	1

[Joachims, 1999]

- training set {D1, D6}
- test set {D2, D3, D4, D5}

Transductive Support Vector Machines [Vapnik]

Objective: maximize margin δ on both training and test examples

Training sample: $Z = [(x_1, y_1), \dots, (x_l, y_l)]$

Test sample: $Z_x^* = [x_1^*, \dots, x_k^*]$

Find solution $W^*(y^*, w) = \frac{1}{\delta^2}$ of

$$\min_{y_1^* \dots y_k^* \in \{-1, 1\}} \min_{w \in \mathfrak{R}^d} w \cdot w$$

subject to

$$y_1[w \cdot x_1 + b] \geq 1$$

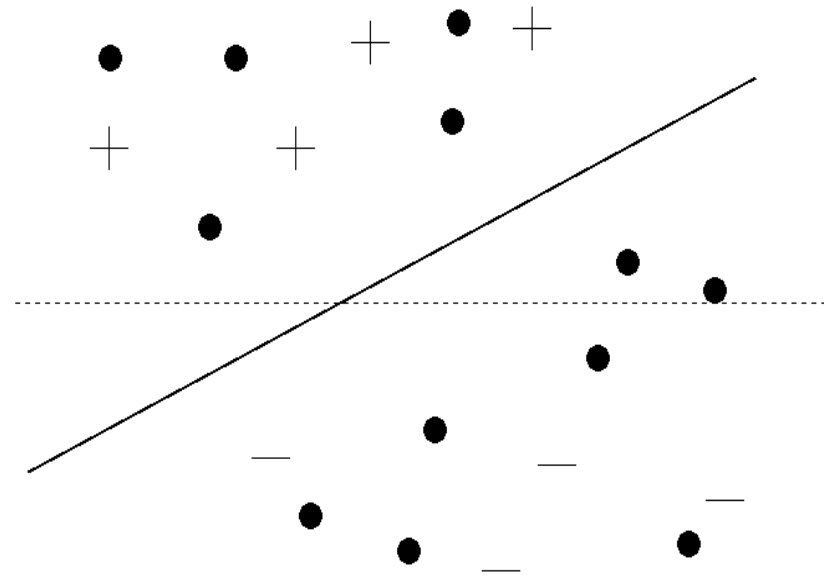
...

$$y_l[w \cdot x_l + b] \geq 1$$

$$y_1^*[w \cdot x_1^* + b] \geq 1$$

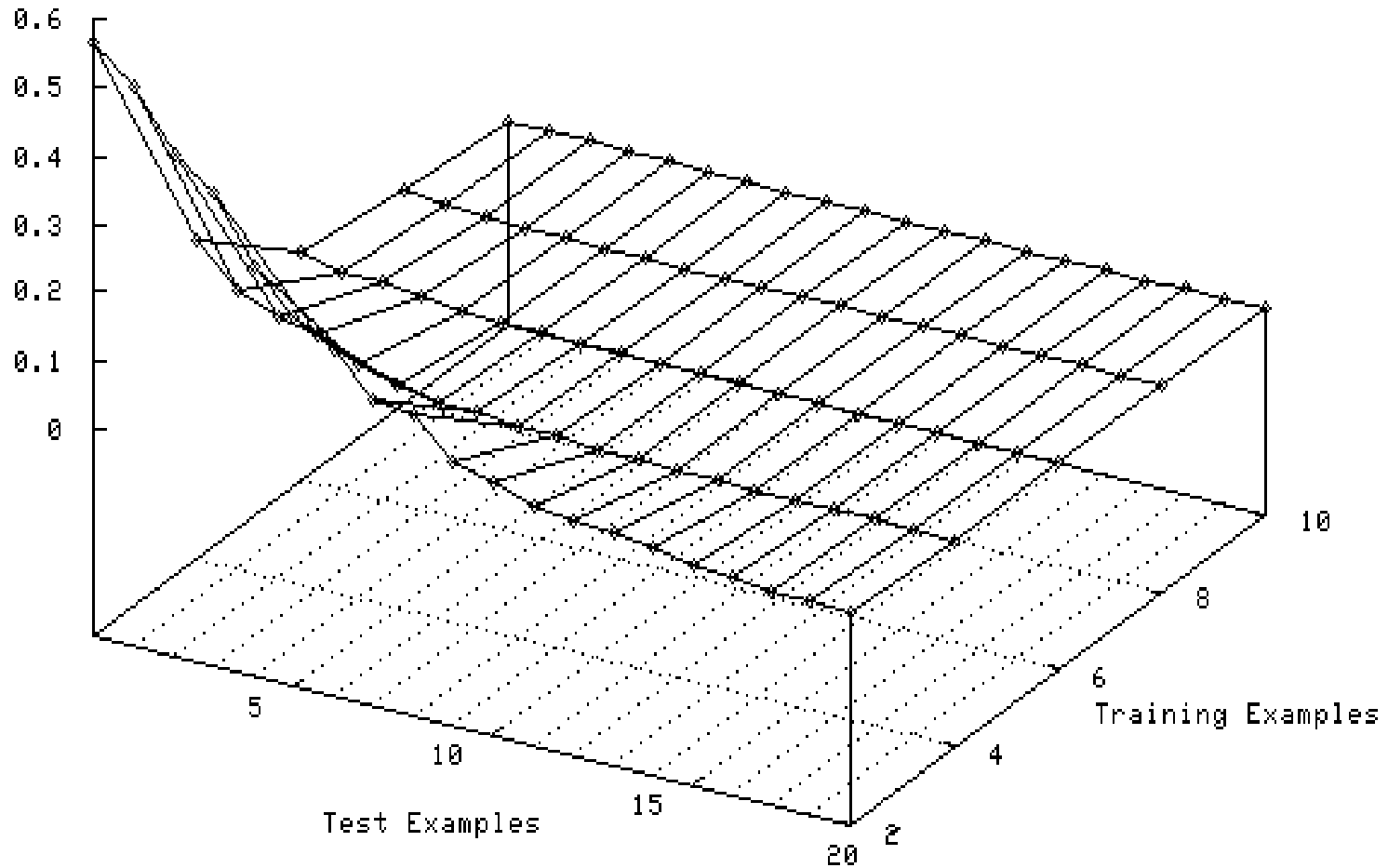
...

$$y_k^*[w \cdot x_k^* + b] \geq 1$$

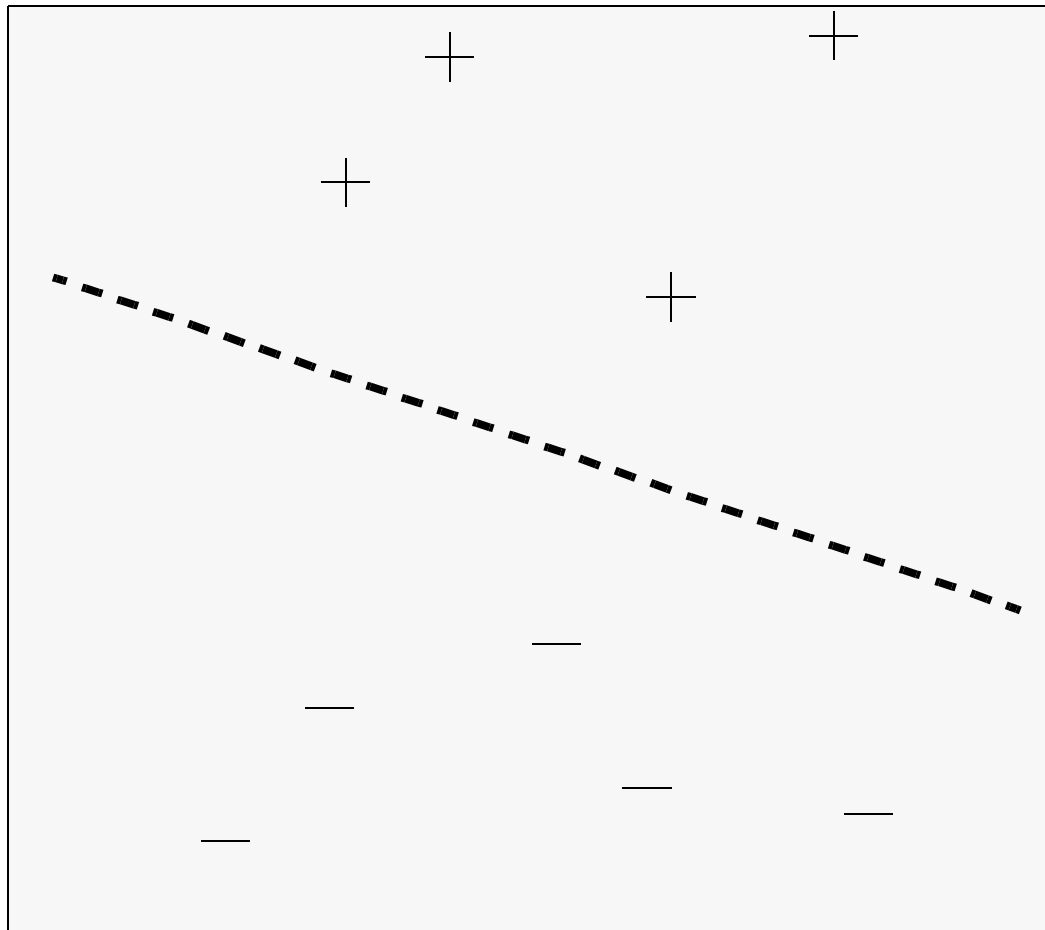


Simulation

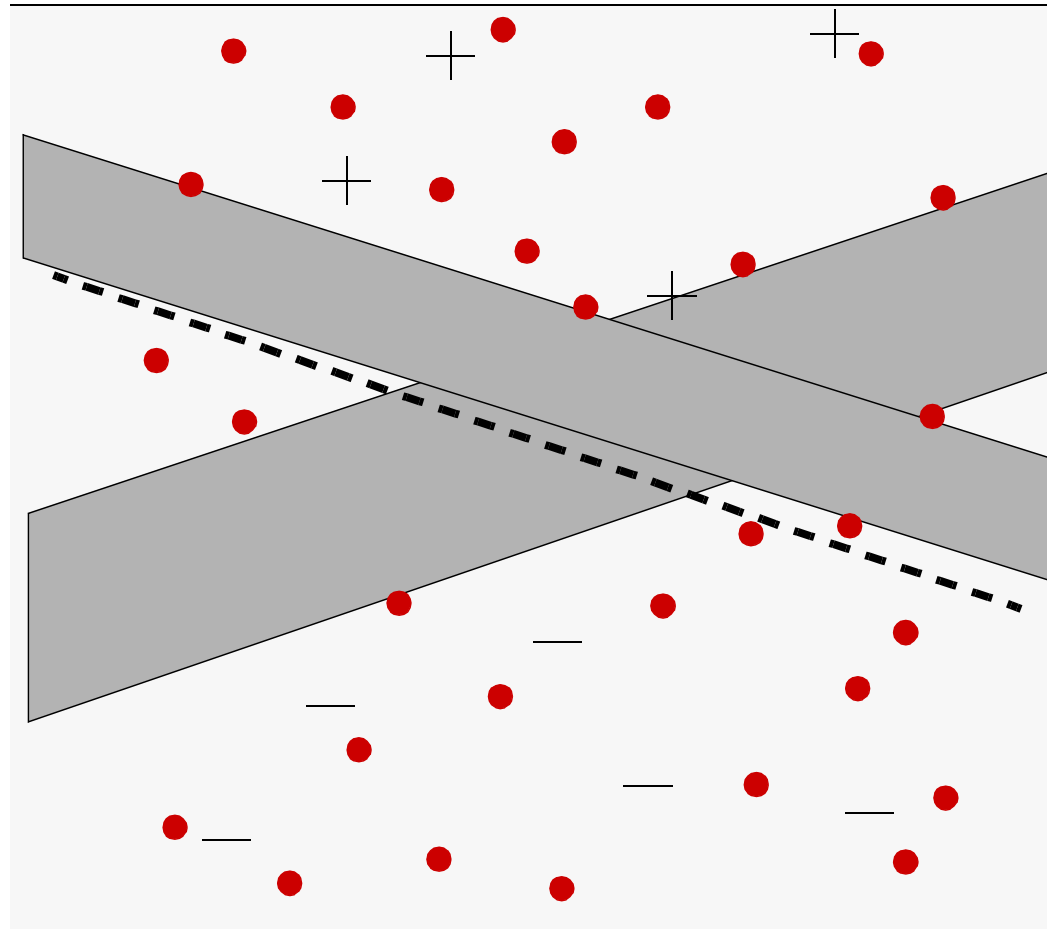
Target concept: $TCat([1:0:1],[0:1:1],[4:4:8])$



Why Does Adding Test Examples Reduce Error?

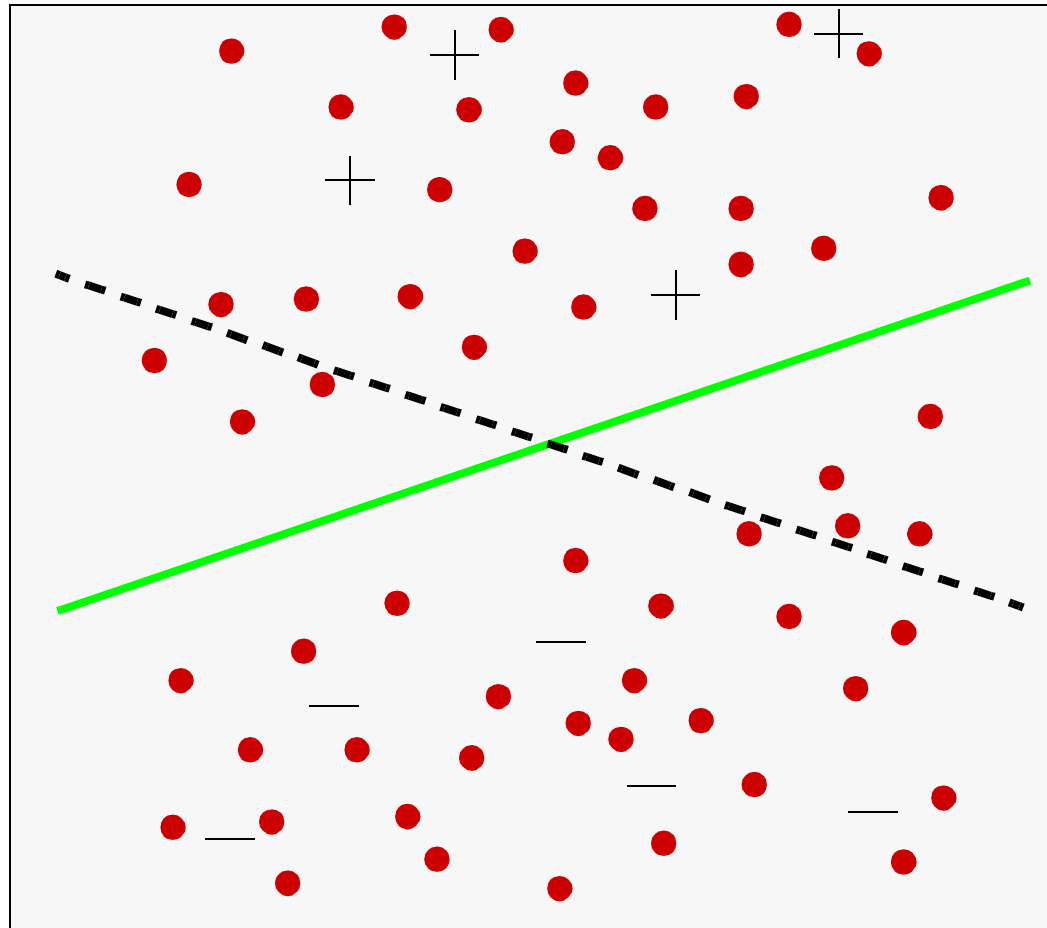


Why Does Adding Test Examples Reduce Error?



$$\text{Margin } \delta \geq \frac{1}{\sqrt{2}}$$

Why Does Adding Test Examples Reduce Error?



$$\text{Margin } \delta \geq \frac{1}{\sqrt{2}}$$

Experiment: Reuters-21578

Reuters Newswire Stories

- 90 categories
- 9603 training documents
- 3299 test documents

Experiment

- 10 most frequent categories
- 17 training documents
- 3299 test documents
- ca. 700-12000 features

	Bayes	SVM	TSVM
earn	78.8	91.3	95.4
acq	57.4	67.8	76.6
money-fx	43.9	41.3	60.0
grain	40.1	56.2	68.5
crude	24.8	40.9	83.6

	Bayes	SVM	TSVM
trade	22.1	29.5	34.0
interest	24.5	35.6	50.8
ship	33.2	32.5	46.3
wheat	19.5	47.9	54.4
corn	14.5	41.3	43.7

=> avg. TSVM run-time: ~ 1 minute 40 seconds

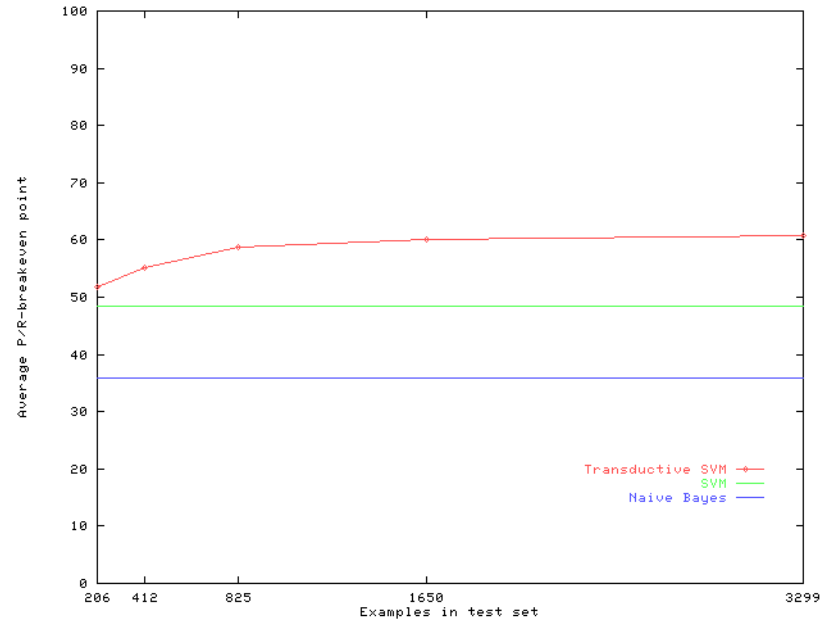
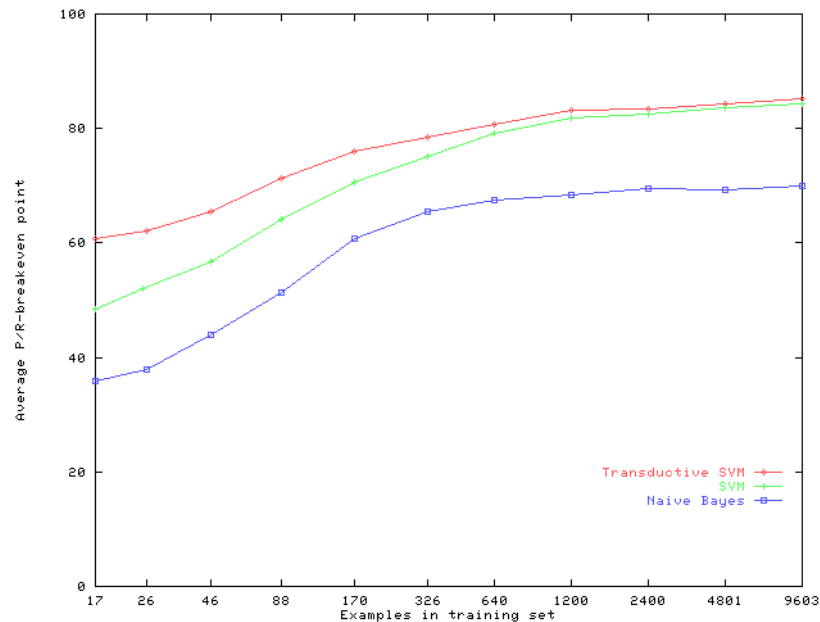
Training Set vs. Test Set

Increasing training set size:

- avg. over 10 Reuters categories
- 3299 test documents
- feature selection: MI with local dictionaries of 1000 for Bayes

Increasing test set size:

- avg. over 10 categories
- 17 training documents



Co-Training (Blum & Mitchell)

Idea: Exploit two sufficiently redundant representations $X = A \times B$.

Scenarios:

- Web-page body text / Hyperlinks pointing to page
- sound of person saying “hello” / image of lip movements

Training example: $(\langle a_i, b_j \rangle, y)$

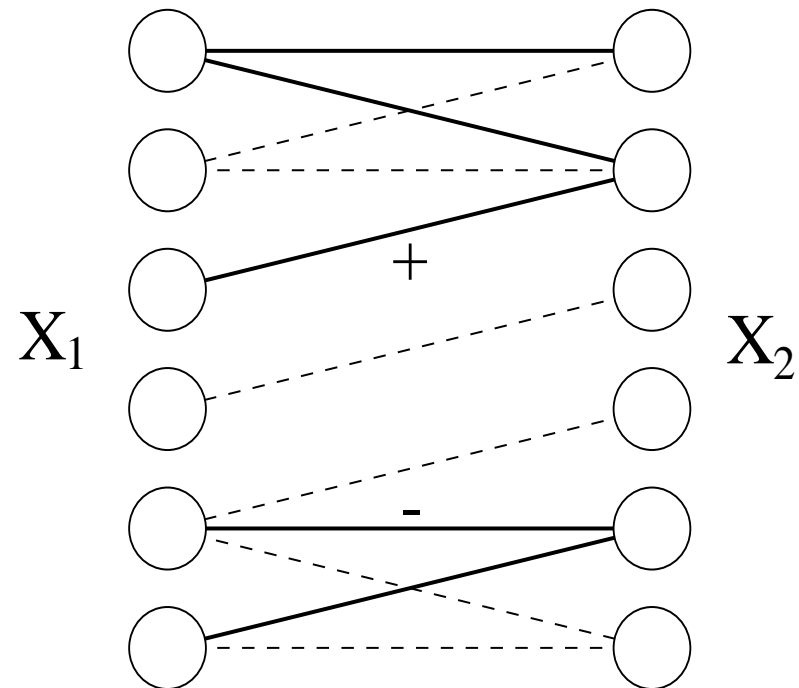
Test example: $\langle a_i, b_j \rangle$

Composition: $\langle a_i, b_j \rangle \in X_1 \times X_2$

Hypotheses: $H_1 \times H_2$

Compatible:

$$\langle a_i, b_j \rangle \Rightarrow h_1(a_i) = h_2(b_j)$$



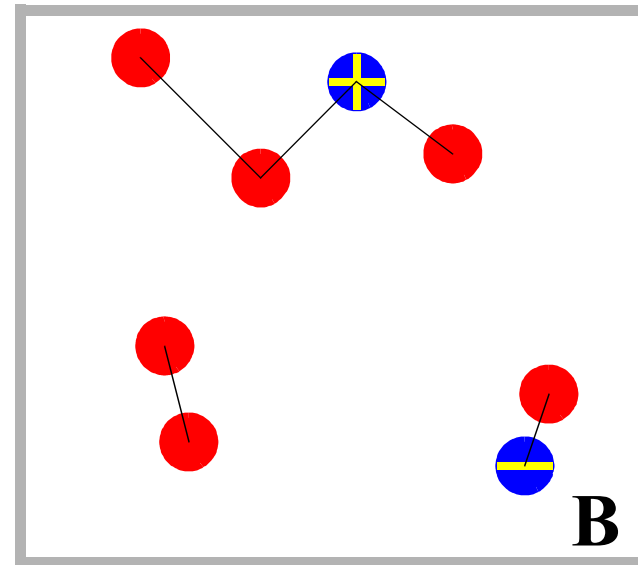
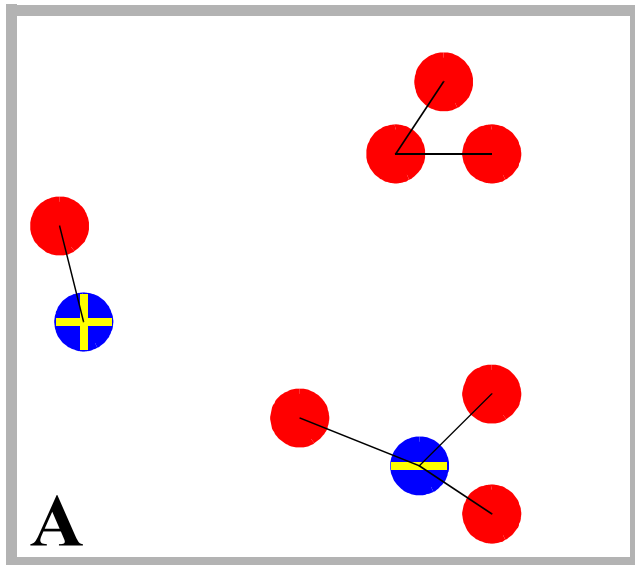
Co-Training [Blum & Mitchell]

Idea: Exploit two sufficiently redundant representations $X = A \times B$.

Scenario:

- Web-page body text (A) / Hyperlinks pointing to page (B)

Compatible: Perfect classifiers on A and B do not disagree!



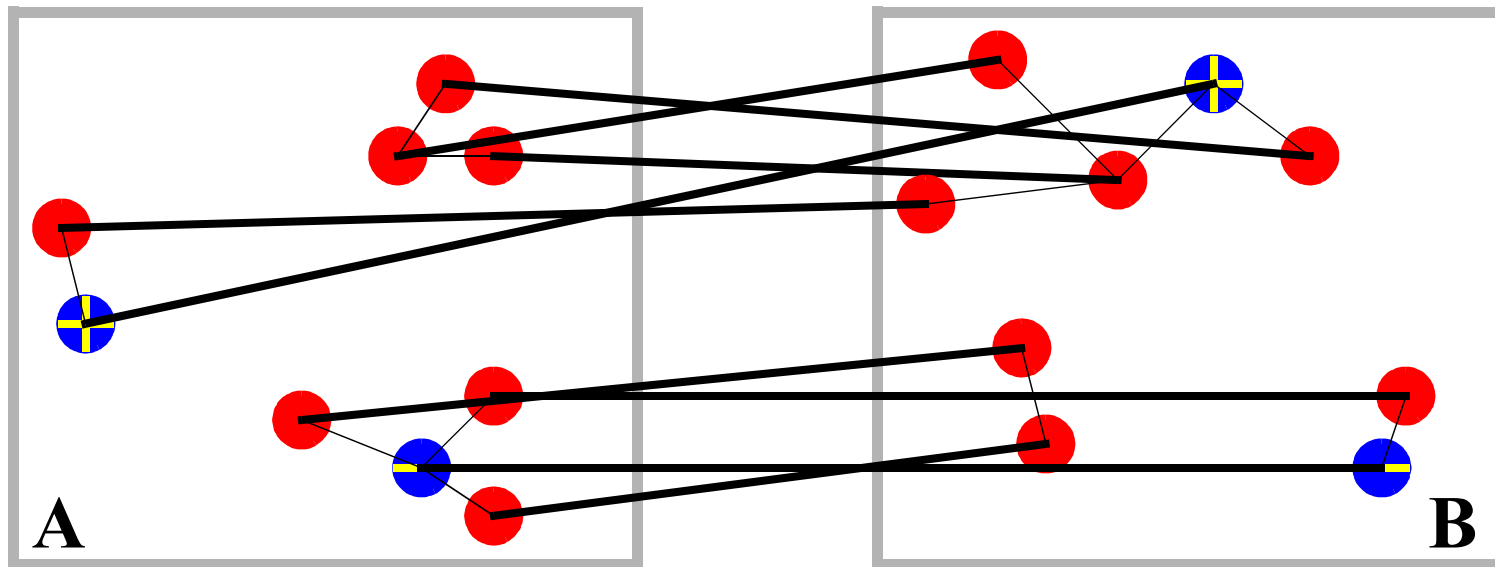
Co-Training [Blum & Mitchell]

Idea: Exploit two sufficiently redundant representations $X = A \times B$.

Scenario:

- Web-page body text (A) / Hyperlinks pointing to page (B)

Compatible: Perfect classifiers on A and B do not disagree!



=> SGT maximizes consistency between two k-NN classifiers.

Co-Training Experiment

	SGT	KNN	TSVM	SVM	B&M
cotrain	3.3	-	-	-	5.0
page+link	5.9	10.1	4.3	20.3	-
page	6.2	13.3	4.6	21.6	12.9
link	22.1	13.1	8.9	18.5	12.4

- Dataset: classifying course homepages from Blum and Mitchell
- 12 training examples, 1039 test examples
- Error on test set averaged over 100 random test/training splits
- Parameters:
 - SGT: cosine similarity, $c = 3200$, $d = 80$, 200NN in each view
 - others: optimized on the test set