

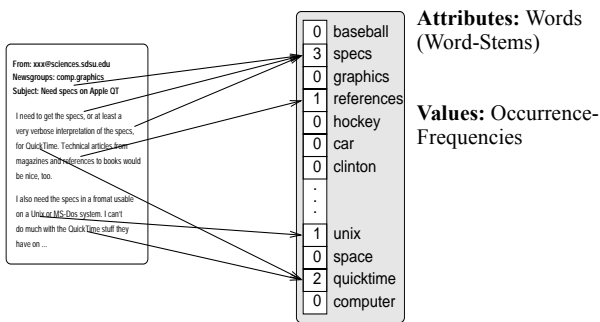
# CS630 Representing and Accessing Digital Information

## Text Classification: Support Vector Machines

Thorsten Joachims  
Cornell University

0

## Representing Text as Attribute Vectors



==> The ordering of words is ignored!

1

## Generative vs. Discriminative Training

### Process:

- Generator: Generates descriptions  $\vec{x}$  according to distribution  $P(\vec{x})$ .
- Teacher: Assigns a value  $y$  to each description  $\vec{x}$  based on  $P(y|\vec{x})$ .

==> Training examples  $(\vec{x}_1, y_1), \dots, (\vec{x}_n, y_n) \sim P(\vec{x}, y) \quad \vec{x}_i \in \mathcal{X}^N \quad y_i \in \{1, -1\}$

### Generative Training

- make assumptions about the parametric form of  $P(\vec{x}, y)$ .
- estimate the parameters of  $P(\vec{x}, y)$  from the training data
- derive optimal classifier using Bayes' rule
- example: naive Bayes

### Discriminative Training

- make assumptions about the set  $H$  of classifiers
- estimate error of classifiers in  $H$  from the training data
- select classifier with lowest error rate
- example: SVM, decision tree

3

## Principle: Empirical Risk Minimization (ERM)

### Learning Principle:

Find the decision rule  $h^\circ \in H$  for which the training error is minimal:

$$h^\circ = \arg \min_{h \in H} \{Err_S(h)\}$$

### Training Error:

$$Err_S(h) = \frac{1}{n} \sum_{i=1}^n \Delta(y_i \neq h(\vec{x}_i))$$

==> Number of misclassifications on training examples.

### Central Problem of Statistical Learning Theory:

When does a low training error lead to a low generalization error?  
(i.e. what is the probability of error after n training examples)

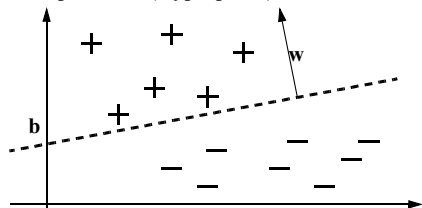
3

## Linear Classifiers

Rules of the Form: weight vector  $\vec{w}$ , threshold  $b$

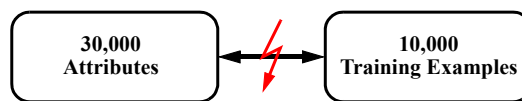
$$h(\vec{x}) = \text{sign} \left[ \sum_{i=1}^N w_i x_i + b \right] = \begin{cases} 1 & \text{if } \sum_{i=1}^N w_i x_i + b > 0 \\ -1 & \text{else} \end{cases}$$

Geometric Interpretation (Hyperplane):



4

## Paradoxon of Text Classification



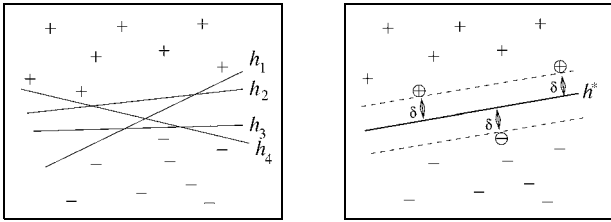
Good News: SVMs can overcome this problem!

Bad News: This does not hold for all high-dimensional problems!

5

## Optimal Hyperplane (SVM Type 1)

**Assumption:** The training examples are linearly separable.



**Support Vectors:** Examples with minimal distance (circled).

## SVM Primal Optimization Problem

**Training Examples:**  $(x_1, y_1), \dots, (x_n, y_n)$   $x_i \in \mathfrak{R}^N$   $y_i \in \{1, -1\}$

**Training:**

• minimize  $P(\vec{w}, b) = \frac{1}{2} \vec{w} \cdot \vec{w}$  over  $\vec{w}$  and  $b$

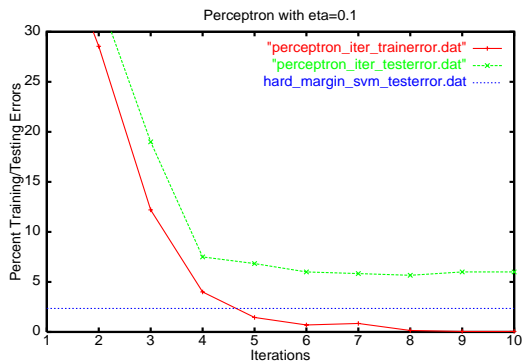
subject to fulfilling the constraints 
$$\begin{cases} y_1 [\vec{w} \cdot \vec{x}_1 + b] \geq 1 \\ \dots \\ y_n [\vec{w} \cdot \vec{x}_n + b] \geq 1 \end{cases}$$

• typically single solution (i. e.  $\vec{w}, b$  is unique)

**Prediction Rule:**

$$h(\vec{x}) = \text{sign} \left[ \sum_{i=1}^N w_i x_i + b \right] = \text{sign} [\vec{w} \cdot \vec{x} + b]$$

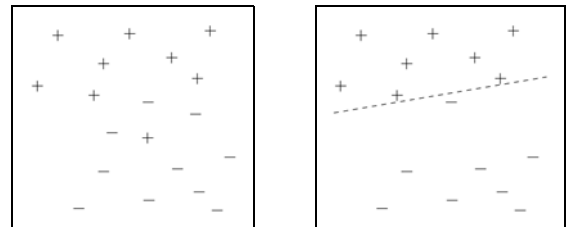
## Example: Optimal Hyperplane vs. Perceptron



Train on 1000 pos / 1000 neg examples for "acq" (Reuters-21578).

## Non-Separable Training Samples

- For some training samples there is no separating hyperplane!
- Complete separation is suboptimal for many training samples!



=> minimize trade-off between margin and training error.

## Soft-Margin Separation

**Idea:** Maximize margin and minimize training error simultaneously.

**Hard Margin:**

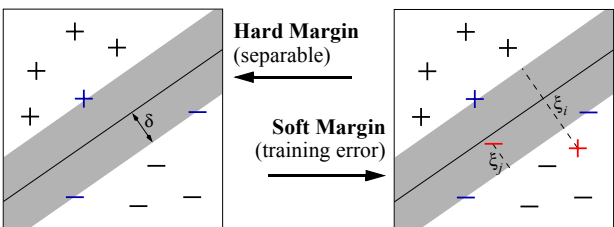
$$\text{minimize } P(\vec{w}, b) = \frac{1}{2} \vec{w} \cdot \vec{w}$$

$$\text{s. t. } y_i [\vec{w} \cdot \vec{x}_i + b] \geq 1$$

**Soft Margin:**

$$\text{minimize } P(\vec{w}, b, \xi) = \frac{1}{2} \vec{w} \cdot \vec{w} + C \sum_{i=1}^n \xi_i$$

$$\text{s. t. } y_i [\vec{w} \cdot \vec{x}_i + b] \geq 1 - \xi_i \text{ and } \xi_i \geq 0$$

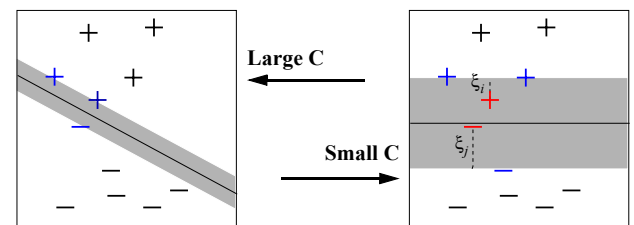


## Controlling Soft-Margin Separation

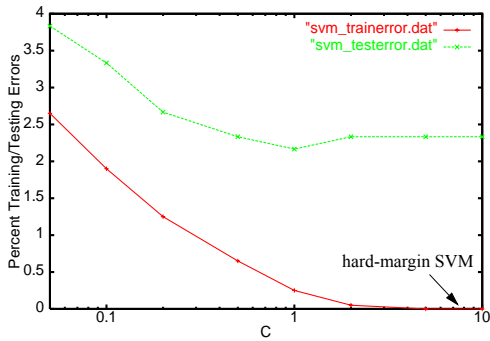
$$\text{Soft Margin: minimize } P(\vec{w}, b, \xi) = \frac{1}{2} \vec{w} \cdot \vec{w} + C \sum_{i=1}^n \xi_i$$

$$\text{s. t. } y_i [\vec{w} \cdot \vec{x}_i + b] \geq 1 - \xi_i \text{ and } \xi_i \geq 0$$

- $\xi_i$  is an upper bound on the number of training errors.
- $C$  is a parameter that controls trade-off between margin and error.



### Example Reuters "acq" : Varying C



**Observation:** Typically no local optima, but not necessarily...

### Solution as Linear Combination

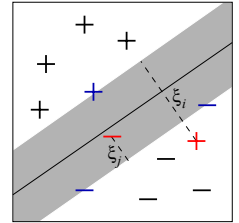
**Primal OP:** minimize  $P(\vec{w}, b, \xi) = \frac{1}{2} \vec{w} \cdot \vec{w} + C \sum_{i=1}^n \xi_i$   
 s. t.  $y_i [\vec{w} \cdot \vec{x}_i + b] \geq 1 - \xi_i$  and  $\xi_i \geq 0$

**Lemma:** The solution  $w^o$  can always be written as a linear combination

$$\vec{w}^o = \sum_{i=1}^n \alpha_i y_i \vec{x}_i \quad 0 \leq \alpha_i \leq C$$

of the training data.

- one factor  $\alpha_i$  for each training example
- “influence” of single training example limited by C
- $0 < \alpha_i < C \iff$  SV with  $\xi_i = 0$
- $\alpha_i = C \iff$  SV with  $\xi_i > 0$
- $\alpha_i = 0$  else
- SVM-light outputs  $\alpha_i$  with option “-a”



### Summary What is a (linear) SVM?

**Given:**

- Training examples  $(\vec{x}_1, y_1), \dots, (\vec{x}_n, y_n)$   $\vec{x}_i \in \mathbb{R}^N$   $y_i \in \{1, -1\}$
- Parameter C for trading-off training error and margin size (often C=1 is a good choice for normalized document vectors)

**Training:**

- Finds hyperplane.
- The hyperplane has maximum margin with minimal training error (upper bound  $\sum \xi_i$ ) given C.
- The result of training are  $\alpha_1, \dots, \alpha_n$ . They determine  $\vec{w}, b$ .

**Classification:** For new example  $h(\vec{x}) = \text{sign}_{TM} \left( \sum_{x_i \in SV} \alpha_i y_i \vec{x}_i \cdot \vec{x} + b \right)$

### Leave-One-Out

**Training set:**  $(\vec{x}_1, y_1), (\vec{x}_2, y_2), (\vec{x}_3, y_3), \dots, (\vec{x}_n, y_n)$

**Approach:** Repeatedly leave one example out for testing.

train on	test on
$(\vec{x}_2, y_2), (\vec{x}_3, y_3), (\vec{x}_4, y_4), \dots, (\vec{x}_n, y_n)$	$(\vec{x}_1, y_1)$
$(\vec{x}_1, y_1), (\vec{x}_3, y_3), (\vec{x}_4, y_4), \dots, (\vec{x}_n, y_n)$	$(\vec{x}_2, y_2)$
$(\vec{x}_1, y_1), (\vec{x}_2, y_2), (\vec{x}_4, y_4), \dots, (\vec{x}_n, y_n)$	$(\vec{x}_3, y_3)$
...	...
$(\vec{x}_1, y_1), (\vec{x}_2, y_2), (\vec{x}_3, y_3), \dots, (\vec{x}_{n-1}, y_{n-1})$	$(\vec{x}_n, y_n)$

**Error estimate:**  
 $Err_{loo}(h) = \frac{1}{n} \sum_{i=1}^n |h(\vec{x}_i) - y_i|$

**Question:** Is there a connection between margin and the estimate?

### Necessary Cond. for Leave-One-Out Error of SVM

**Lemma:** SVM  $[h_i(\vec{x}_i) \neq y_i] \iff [2\alpha_i R^2 + \xi_i \geq 1]$  [Joachims, 2000] [Jaakkola & Haussler, 1999] [Vapnik & Chapelle, 2000]

**Input:**

- $\alpha_i$  dual variable of example i
- $\xi_i$  slack variable of example i
- $\|\vec{x}\| \leq R$  bound on length

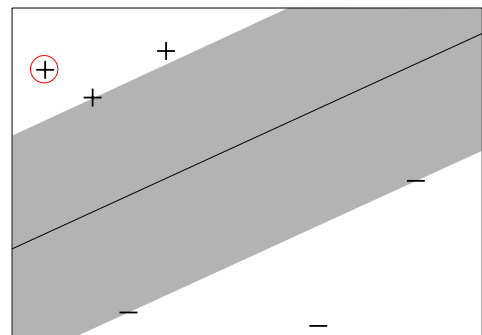
**Example:**

$\rho \alpha_i R^2 + \xi_i$	leave-one-out error
0.0	OK
0.7	OK
3.5	ERROR
0.1	OK
1.3	OK
0.0	OK
0.0	OK
...	...

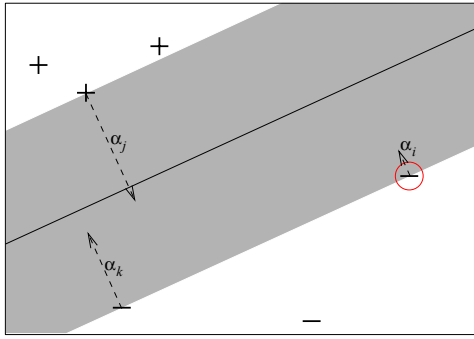
Available after training SVM on the full training data

### Case 1: Example is no SV

$(\alpha_i = 0) \iff (\xi_i = 0) \iff (2\alpha_i R^2 + \xi_i < 1) \iff$  no leave-one-out error



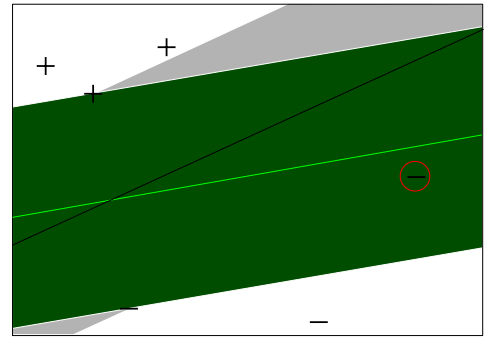
### Case 2: Example is SV with Low Influence



$$\alpha_i < \frac{0.5}{R^2} < C \quad \vee \quad (\xi_i = 0) \quad \vee \quad (2\alpha_i R^2 + \xi_i < 1) \quad \vee \quad \text{no leave-one-out error}$$

18

### Case 2: Example is SV with Low Influence

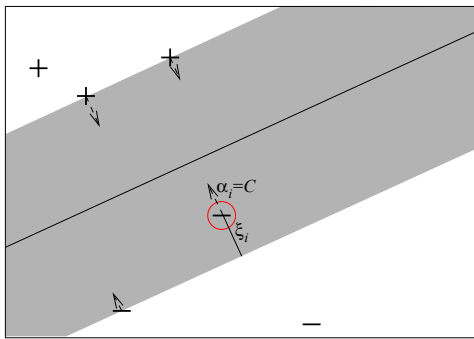


$$\alpha_i < \frac{0.5}{R^2} < C \quad \vee \quad (\xi_i = 0) \quad \vee \quad (2\alpha_i R^2 + \xi_i < 1) \quad \vee \quad \text{no leave-one-out error}$$

19

### Case 3: Example has Small Training Error

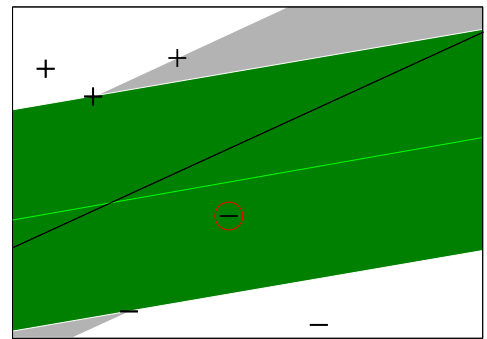
$$(\alpha_i = C) \wedge (\xi_i < 1 - 2CR^2) \quad \vee \quad (2\alpha_i R^2 + \xi_i < 1) \quad \vee \quad \text{no leave-one-out error}$$



20

### Case 3: Example has Small Training Error

$$(\alpha_i = C) \wedge (\xi_i < 1 - 2CR^2) \quad \vee \quad (2\alpha_i R^2 + \xi_i < 1) \quad \vee \quad \text{no leave-one-out error}$$

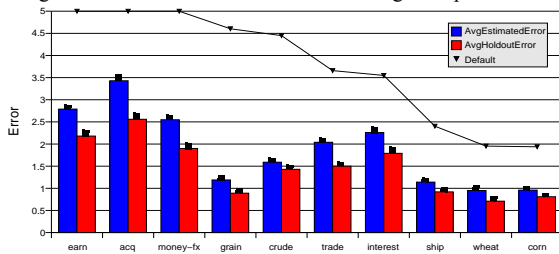


21

### Experiment: Reuters-21578

- 6451 training examples
- 6451 test examples for holdout testing
- ~27,000 features

Average error estimate over 10 random training/test splits:



=> small bias, variance of estimators is approximately equal

22

### Fast Leave-One-Out Estimation for SVMs

**Lemma:** Training errors are always leave-out-out errors.

- Algorithm:**
- $(R, \alpha, \xi) = \text{train\_SVM}(X, 0, 0)$ ;
  - for all training examples, do
    - if  $\xi_i > 1$  then  $\text{loo}++$ ;
    - else if  $(\rho\alpha_i R^2 + \xi_i < 1)$  then  $\text{loo}=\text{loo}$ ;
    - else  $\text{train\_SVM}(X_i, \alpha, \xi)$ ;

**Experiment:**

	Training Examples	Retraining Steps (%)		CPU-Time (sec)	
		$\rho = 1$	$\rho = 2$	$\rho = 1$	$\rho = 2$
Reuters	6451	0.20%	0.58%	11.1	32.3
WebKB	2092	6.78%	20.42%	78.5	235.4
Ohsumed	10000	1.07%	2.56%	433.0	1132.3

23

## Expected Error Rate of SVM

**Leave-One-Out Error Estimate:**  $Err_{loo}(h) = \frac{1}{n} \sum_{i=1}^n |h_i(\tilde{x}_i) - y_i|$

For separable problems:

$$\left[ h_i(\tilde{x}_i) \neq y_i \right] \vee \left[ \alpha_i R^2 \geq 1 \right] \quad \left\| \tilde{x} \right\| \leq R$$

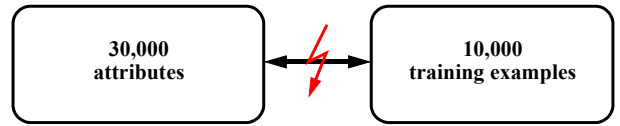
$$\Rightarrow Err_{loo}(h_{SVM}) \leq \frac{1}{n} \sum_{i=1}^n \left[ \alpha_i R^2 \geq 1 \right] \leq \frac{1}{n} \sum_{i=1}^n \alpha_i R^2 \leq \frac{1}{n} \frac{R^2}{\delta^2}$$

**Bound on the expected error rate [Vapnik, 98]:**

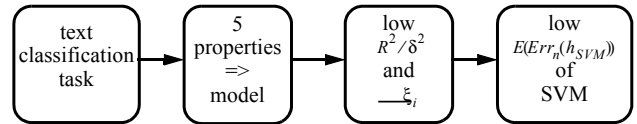
$$E[Err_{true(n)}(h_{SVM})] = E[Err_{loo(n+1)}(h_{SVM})] \leq \frac{1}{n+1} E \left[ \frac{R^2}{\delta^2} \right]$$

24

## Why Do SVMs Work Well for Text Classification?



A statistical learning model of text classification with SVMs:



25

## Margin/Loss Based Bound on the Expected Error

**Theorem:** The expected error of a soft margin SVM is bounded by

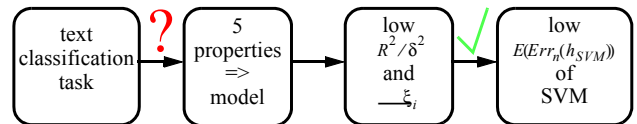
$$E(Err_n(h_{SVM})) \leq \frac{\rho E \left[ \frac{R^2}{\delta^2} \right] + \rho C R^2 E \left[ \frac{\xi_i}{\delta} \right]}{n+1} \quad C \geq \frac{1}{\rho R^2}$$

$$E(Err_n(h_{SVM})) \leq \frac{\rho E \left[ \frac{R^2}{\delta^2} \right] + \rho(CR^2 + 1) E \left[ \frac{\xi_i}{\delta} \right]}{n+1} \quad C < \frac{1}{\rho R^2}$$

Where  $E \left[ \frac{R^2}{\delta^2} \right]$  is the expected soft margin and  $E \left[ \frac{\xi_i}{\delta} \right]$  is the expected training loss on training sets of size  $n+1$ .

26

## First Step Completed



27

## Properties 1+2: Sparse Examples in High Dimension

- High dimensional feature vectors (30,000 features)
- Sparse document vectors: only a few words of the whole language occur in each document

	Training Examples	Number of Features	Distinct Words (Sparsity)
Reuters Newswire Articles	9,603	27,658	74 (0.27%)
Ohsumed MeSH Abstracts	10,000	38,679	100 (0.26%)
WebKB WWW-Pages	3,957	38,359	130 (0.34%)

28

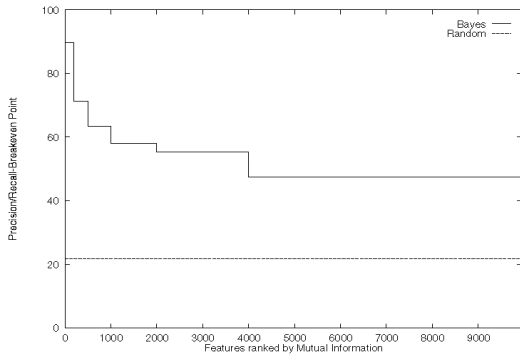
## Property 3: Heterogeneous Use Of Words

<b>MODULAIRE BUYS BOISE HOMES PROPERTY</b> Modulaire Industries said it acquired the design library and manufacturing rights of privately-owned Boise Homes for an undisclosed amount of cash. Boise Homes sold commercial and residential prefabricated structures, Modulaire said.	<b>USX, CONS. NATURAL END TALKS</b> USX Corp' s Texas Oil and Gas Corp subsidiary and Consolidated Natural Gas Co have mutually agreed not to pursue further their talks on Consolidated' s possible purchase of Apollo Gas Co from Texas Oil. No details were given.
<b>JUSTICE ASKS U.S. DISMISSAL OF TWA FILING</b> The Justice Department told the Transportation Department it supported a request by USAir Group that the DOT dismiss an application by Trans World Airlines Inc for approval to take control of USAir. "Our rationale is that we reviewed the application for control filed by TWA with the DOT and ascertained that it did not contain sufficient information upon which to base a competitive review." James Weiss, an official in Justice' s Antitrust Division, told Reuters.	<b>E.D. And F. MAN TO BUY INTO HONG KONG FIRM</b> The U.K. Based commodity house E.D. And F. Man Ltd and Singapore' s Yeo Hiap Seng Ltd jointly announced that Man will buy a substantial stake in Yeo' s 71.1 pct held unit, Yeo Hiap Seng Enterprises Ltd. Man will develop the locally listed soft drinks manufacturer into a securities and commodities brokerage arm and will rename the firm Man Pacific (Holdings) Ltd.

No pair of documents shares any words, but "it", "the", "and", "of", "for", "an", "a", "not", "that", "in".

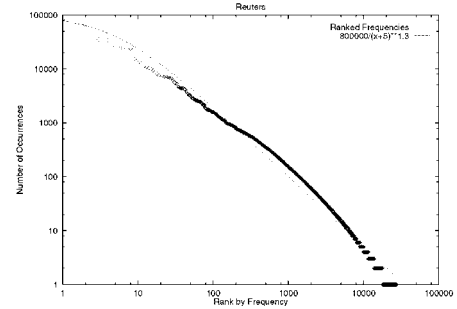
29

### Property 4: High Level Of Redundancy



=> Few features are irrelevant!

### Property 5: "Zipf's Law"



Zipf's Law: In text, the  $i$ -th frequent word occurs  $f_i = \frac{k}{(c+i)^\theta}$  times.

=> Most words occur very infrequently!

### Text Classification Model

**Definition:** For the TCat-concept there are  $s$  disjoint sets of features.

$$TCat([p_1|n_1|f_1], \dots, [p_s|n_s|f_s])$$

Each positive (negative) example contains  $p_i$  ( $n_i$ ) occurrences from the  $f_i$  features in set  $i$ .

**Example:**  $TCat([20|20|100], [4|1|200], [1|4|200], [5|5|600], [9|1|3000], [1|9|3000], [10|10|4000])$

### TCat-Concept for WebKB "Course"

$$TCat([77|29|98], [4|21|52], [16|2|431], [1|12|341], [9|1|5045], [1|21|24276], [169|191|8116])$$

	high frequency	medium frequency	low frequency
pos	98 words	431 words	5045 words
neg	52 words	341 words	24276 words

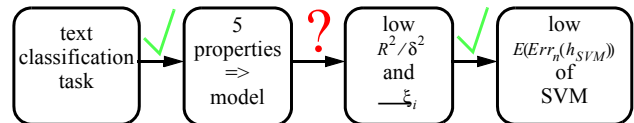
### Real Text Classification Tasks as TCat-Concepts

$$\text{Reuters "Earn": } TCat([33|2|65], [32|65|152], [2|1|171], [3|21|974], [3|1|3455], [1|10|17020], [78|52|5821])$$

$$\text{Webkb "Course": } TCat([77|29|98], [4|21|52], [16|2|431], [1|12|341], [9|1|5045], [1|21|24276], [169|191|8116])$$

$$\text{Ohsumed "Pathology": } TCat([2|1|10], [1|4|22], [2|1|92], [1|2|94], [5|1|4080], [1|10|20922], [197|190|13459])$$

### Second Step Completed



## The Margin $\delta^2$ of TCat-Concepts

**Lemma 1:** For  $TCat([p_1|n_1|f_1], \dots, [p_s|n_s|f_s])$  -concepts there is always a hyperplane passing through the origin with margin  $\delta^2$  at least

$$a = \frac{\sum_{i=1}^s p_i^2}{f_i}$$

$$\delta^2 \geq \frac{ad - b^2}{a + 2b + d} \quad \text{with} \quad d = \frac{\sum_{i=1}^s n_i^2}{f_i}$$

$$b = \frac{\sum_{i=1}^s n_i p_i}{f_i}$$

**Example:** The previous example WebKB “course” has a margin of at least

$$\delta^2 \geq 0.23$$

36

## The Length $R^2$ of Document Vectors

**Lemma 2:** If the ranked term frequencies  $f_i$  in a document with  $l$  words have the form of the generalized Zipf’s Law

$$f_i = \frac{k}{(c+i)^\Theta}$$

based on their frequency rank  $i$ , then the Euclidean length of the document vector  $x$  is bounded by

$$\|x\| \leq \sqrt{\frac{d}{\sum_{i=1}^l \frac{k}{(c+i)^\Theta}}}$$

with  $\frac{d}{\sum_{i=1}^l \frac{k}{(c+i)^\Theta}} = l$

**Example:** For WebKB “course” with

$$f_i = \frac{470000}{(5+i)^{1.25}}$$

follows that  $R^2 \leq 1900$ .

37

## $R^2$ , $\delta^2$ , and $\xi_i$ for Text Classification

### Reuters Newswire Stories

- 10 most frequent categories
- 9603 training examples
- 27658 attributes

$$E(\text{Err}_n(h_{SVM})) \leq \frac{E\left(\frac{R^2}{\delta^2}\right) + CR^2 E\left(\frac{\xi_i}{n}\right)}{n+1}$$

	$R^2/\delta^2$	$\xi_i$
earn	1143	0
acq	1848	0
money-fx	1489	27
grain	585	0
crude	810	4

	$R^2/\delta^2$	$\xi_i$
trade	869	9
interest	2082	33
ship	458	0
wheat	405	2
corn	378	0

38

## Learnability of TCat-Concepts

**Theorem:** For  $TCat([p_1|n_1|f_1], \dots, [p_s|n_s|f_s])$  -concepts and documents with  $l$  words that follow the generalized Zipf’s Law  $f_i = k/(c+i)^\Theta$  the expected generalization error of an unbiased SVM after training on  $n$  examples is bounded by

$$E(\text{Err}_n(h_{SVM})) \leq \frac{R^2}{n+1} \frac{ad - b^2}{a + 2b + d} \quad \text{with}$$

$$a = \frac{\sum_{i=1}^s p_i^2}{f_i}$$

$$d = \frac{\sum_{i=1}^s n_i^2}{f_i}$$

$$b = \frac{\sum_{i=1}^s n_i p_i}{f_i}$$

$$R^2 \leq \frac{\sum_{i=1}^l \frac{k}{(c+i)^\Theta}}{l}$$

39

## Comparison Theory vs. Experiments

	Learning Curve Bound	Predicted Bound on Error Rate	Error Rate in Experiment
Reuters “earn”	$E(\text{Err}_n(h_{SVM})) \leq \frac{138}{n+1}$	1.5%	1.3%
WebKB “course”	$E(\text{Err}_n(h_{SVM})) \leq \frac{443}{n+1}$	11.2%	4.4%
Ohsumed “pathology”	$E(\text{Err}_n(h_{SVM})) \leq \frac{9457}{n+1}$	94.6%	23.1%

- Model can differentiate between “difficult” and “easy” tasks
- Predicts and reproduces the effect of information retrieval heuristics (e.g. TFIDF-weighting)

40

## Sensitivity Analysis

What makes a text classification problem suitable for a linear SVM?

### High Redundancy:

$$TCat_{TM}^{(R)} \left[ \begin{array}{l} [40|40|50] \\ [5|1000], [5|25|1000] \\ [30|30|30000] \end{array} \right] \left. \begin{array}{l} \text{high frequency} \\ \text{medium frequency} \\ \text{low frequency} \end{array} \right\}$$

### High Discriminatory Power:

$$TCat_{TM}^{(R)} \left[ \begin{array}{l} [40|40|50] \\ [15|0|500], [0|15|500], [15|15|1000] \\ [30|30|30000] \end{array} \right] \left. \begin{array}{l} \text{high frequency} \\ \text{medium frequency} \\ \text{low frequency} \end{array} \right\}$$

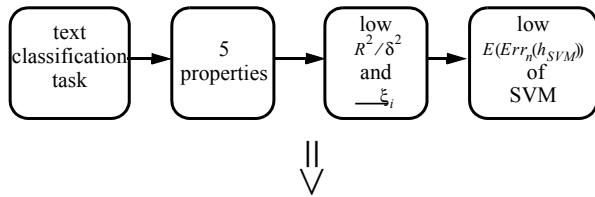
### High Frequency:

$$TCat_{TM}^{(R)} \left[ \begin{array}{l} [16|4|10], [4|16|10], [20|20|30] \\ [30|30|2000] \\ [30|30|30000] \end{array} \right] \left. \begin{array}{l} \text{high frequency} \\ \text{medium frequency} \\ \text{low frequency} \end{array} \right\}$$

41

### What does this Model Provide?

Connects the statistical properties of text classification tasks with generalization error of SVM!



- Explains the behavior of (linear) SVMs on text classification tasks
- Gives guideline for when to apply (linear) SVMs
- Provides formal basis for developing new methods

42

### Summary When do (Linear) SVMs Work Well?



**Intuition:** If the problem can be cast as a TCat-concept with

- high redundancy,
- strongly discriminating features
- particularly in the high frequency region

then linear SVMs achieve a low generalization error.

**Assumptions and Restrictions:**

- no noise (attribute and classification)
- no variance (only “average” examples)
- only upper bounds, no lower bounds

43