

COMP 631: Introduction to Information Retrieval

04/01/2022

XIA (BEN) HU
CS, Rice University

<https://cs.rice.edu/~xh37/index.html>

Today's topic

- Text Classification 1

Introduction to

Standing queries

- The path from IR to text classification:
 - You have an information need to monitor, say:
 - Unrest in the Niger delta region
 - You want to rerun an appropriate query periodically to find new news items on this topic
 - You will be sent new documents that are found
 - I.e., it's not ranking but classification (relevant vs. not relevant)
- Such queries are called standing queries
 - Long used by “information professionals”
 - A modern mass instantiation is Google Alerts
- Standing queries are (hand-written) text classifiers

Spam filtering

text classification task

From: "" <takworlld@hotmail.com>

Subject: real estate is the only way... gem oalvgkay

Anyone can buy real estate with no money down

Stop paying rent TODAY !

There is no need to spend hundreds or even thousands for similar courses

I am 22 years old and I have already purchased 6 properties using the methods outlined in this truly INCREDIBLE ebook.

Change your life NOW !

=====

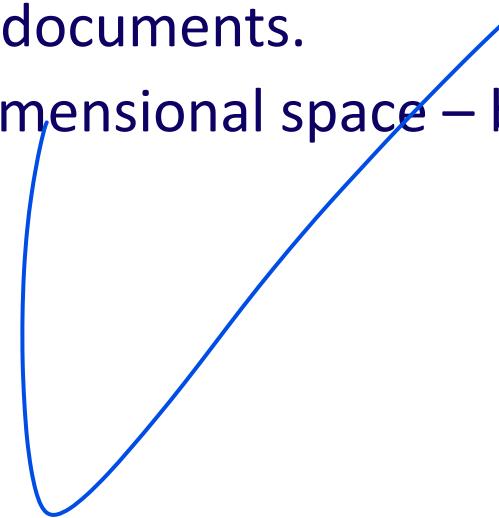
Click Below to order:

<http://www.wholesaledaily.com/sales/nmd.htm>

=====

Categorization/Classification

- Given:
 - A representation of a document d
 - Issue: how to represent text documents.
 - Usually some type of high-dimensional space – bag of words
 - A fixed set of classes:
$$C = \{c_1, c_2, \dots, c_J\}$$
- Determine:
 - The category of d : $\gamma(d) \in C$, where $\gamma(d)$ is a classification function
 - We want to build classification functions (“classifiers”).

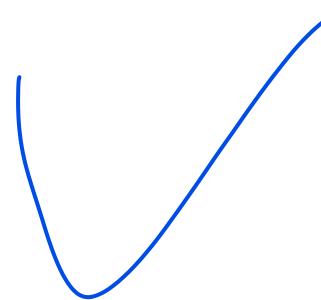


Classification Methods (1)

- Manual classification
 - Used by the original Yahoo! Directory
 - Looksmart, about.com, ODP, PubMed
 - Accurate when job is done by experts
 - Consistent when the problem size and team is small
 - Difficult and expensive to scale
 - Means we need automatic classification methods for big problems

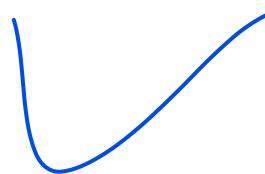
Classification Methods (2)

- Hand-coded rule-based classifiers
 - One technique used by news agencies, intelligence agencies, etc.
 - Widely deployed in government and enterprise
 - Vendors provide “IDE” for writing such rules



Classification Methods (2)

- Hand-coded rule-based classifiers
 - Commercial systems have complex query languages
 - Accuracy is can be high if a rule has been carefully refined over time by a subject expert
 - Building and maintaining these rules is expensive



A Verity topic

A complex classification rule: art

```
comment line      # Beginning of art topic definition
top-level topic   art ACCRUE
topic definition modifiers
    /author = "fsmith"
    /date   = "30-Dec-01"
    /annotation = "Topic created
                  by fsmith"
subtopic topic    * 0.70 performing-arts ACCRUE
evidence topic    ** 0.50 WORD
                   /wordtext = ballet
topic definition modifier
evidence topic    ** 0.50 STEM
                   /wordtext = dance
topic definition modifier
evidence topic    ** 0.50 WORD
                   /wordtext = opera
topic definition modifier
evidence topic    ** 0.30 WORD
                   /wordtext = symphony
topic definition modifier
subtopic          * 0.70 visual-arts ACCRUE
                   ** 0.50 WORD
                   /wordtext = painting
                   ** 0.50 WORD
                   /wordtext = sculpture
subtopic          * 0.70 film ACCRUE
                   ** 0.50 STEM
                   /wordtext = film
subtopic          ** 0.50 motion-picture PHRASE
                   *** 1.00 WORD
                   /wordtext = motion
                   *** 1.00 WORD
                   /wordtext = picture
                   ** 0.50 STEM
                   /wordtext = movie
subtopic          * 0.50 video ACCRUE
                   ** 0.50 STEM
                   /wordtext = video
                   ** 0.50 STEM
                   /wordtext = vcr
# End of art topic
```

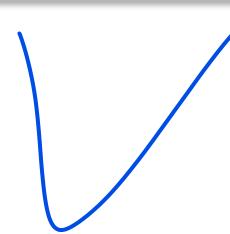
Note:

- maintenance issues
(author, etc.)
- Hand-weighting of terms

[Verity was bought by
Autonomy, which was
bought by HP ...]

Classification Methods (3): Supervised learning

- Given:
 - A document d
 - A fixed set of classes:
$$C = \{c_1, c_2, \dots, c_J\}$$
 - A training set D of documents each with a label in C
- Determine:
 - A learning method or algorithm which will enable us to learn a classifier γ
 - For a test document d , we assign it the class $\gamma(d) \in C$



Classification Methods (3)

- Supervised learning
 - Naive Bayes (simple, common) – see video
 - k-Nearest Neighbors (simple, powerful)
 - Support-vector machines (newer, generally more powerful)
 - ... plus many other methods
 - No free lunch: requires hand-classified training data
 - But data can be built up (and refined) by amateurs
- Many commercial systems use a mixture of methods

Features

- Supervised learning classifiers can use any sort of feature
 - URL, email address, punctuation, capitalization, dictionaries, network features
- In the simplest bag of words view of documents
 - We use only word features
 - we use all of the words in the text (not a subset)

The bag of words representation

Y(

I love this movie! It's sweet, but with satirical humor. The dialogue is great and the adventure scenes are fun... It manages to be whimsical and romantic while laughing at the conventions of the fairy tale genre. I would recommend it to just about anyone. I've seen it several times, and I'm always happy to see it again whenever I have a friend who hasn't seen it yet.

) = C

The bag of words representation

$\gamma($

great	2
love	2
recommend	1
laugh	1
happy	1
...	...

) = C

Feature Selection: Why?

- Text collections have a large number of features
 - 10,000 – 1,000,000 unique words ... and more
- Selection may make a particular classifier feasible
 - Some classifiers can't deal with 1,000,000 features
- Reduces training time
 - Training time for some methods is quadratic or worse in the number of features
- Makes runtime models smaller and faster
- Can improve generalization (performance)
 - Eliminates noise features
 - Avoids overfitting

Feature Selection: Frequency

- The simplest feature selection method:
 - Just use the commonest terms
 - No particular foundation
 - But it make sense why this works
 - They're the words that can be well-estimated and are most often available as evidence
 - In practice, this is often 90% as good as better methods
 - Smarter feature selection

Naïve Bayes

- Classify based on prior weight of class and conditional parameter for what each word says:

$$c_{NB} = \operatorname{argmax}_{c_j \in C} \left[\log P(c_j) + \sum_{i \in positions} \log P(x_i | c_j) \right]$$

- Training is done by counting and dividing:

$$P(c_j) \leftarrow \frac{N_{c_j}}{N} \quad P(x_k | c_j) \leftarrow \frac{T_{c_j x_k} + \alpha}{\sum_{x_i \in V} [T_{c_j x_i} + \alpha]}$$

- Don't forget to smooth

Bayes Theorem

- A probabilistic approach for solving classification problems

- Conditional Probability: $P(Y | X) = \frac{P(X, Y)}{P(X)}$

$$P(X | Y) = \frac{P(X, Y)}{P(Y)}$$

- Bayes theorem:

$$P(Y | X) = \frac{P(X | Y)P(Y)}{P(X)}$$

Example of Bayes Theorem

- Given:
 - A doctor knows that meningitis causes stiff neck 50% of the time
 - The probability that a patient has meningitis is 1/50,000
 - The probability that a patient has stiff neck is 1/20
- If a patient has stiff neck, what's the probability he/she has meningitis?

$$P(M | S) = \frac{P(S | M)P(M)}{P(S)} = \frac{0.5 \times 1/50000}{1/20} = 0.0002$$

Joint Probability Distribution

- An example of two Boolean variables

	Toothache	!Toothache
Cavity	0.04	0.06
!Cavity	0.01	0.89

- Observations: mutually exclusive and collectively exhaustive
- What are $P(\text{Cavity})$, $P(\text{Cavity} \vee \text{Toothache})$, $P(\text{Cavity}|\text{Toothache})$?

Using Bayes Theorem for Classification

- Consider each attribute and class label as random variables
- Given a record with attributes (X_1, X_2, \dots, X_d)
 - Goal is to predict class Y
 - Specifically, we want to find the value of Y that maximizes $P(Y | X_1, X_2, \dots, X_d)$
- Can we estimate $P(Y | X_1, X_2, \dots, X_d)$ directly from data?

Using Bayes Theorem for Classification

- Approach:
 - Compute posterior probability $P(Y | X_1, X_2, \dots, X_d)$ using the Bayes theorem

$$P(Y | X_1 X_2 \dots X_d) = \frac{P(X_1 X_2 \dots X_d | Y) P(Y)}{P(X_1 X_2 \dots X_d)}$$

- Choose Y that maximizes
 $P(Y | X_1, X_2, \dots, X_d)$
- Equivalent to choosing value of Y that maximizes
 $P(X_1, X_2, \dots, X_d | Y) P(Y)$
- How to estimate $P(X_1, X_2, \dots, X_d | Y)$?

Conditional Independence

- X and Y are conditionally independent given Z if $P(X|YZ) = P(X|Z)$
- Example: Arm length and reading skills
 - Young child has shorter arm length and limited reading skills, compared to adults
 - If age is fixed, no apparent relationship between arm length and reading skills
 - Arm length and reading skills are conditionally independent given age

Naïve Bayes Classifier

- Assume independence among attributes X_i when class is given:
 - $P(X_1, X_2, \dots, X_d | Y_j) = P(X_1 | Y_j) P(X_2 | Y_j) \dots P(X_d | Y_j)$
 - $P(X_i | Y_j)$ can be estimated for all X_i and Y_j from data
 - A new point is classified to Y_j if $P(Y_j) \prod P(X_i | Y_j)$ is maximal.

Estimate Probabilities from Data

Tid	Refund	Marital Status	Taxable Income	Evade
1	Yes	Single	125K	No
2	No	Married	100K	No
3	No	Single	70K	No
4	Yes	Married	120K	No
5	No	Divorced	95K	Yes
6	No	Married	60K	No
7	Yes	Divorced	220K	No
8	No	Single	85K	Yes
9	No	Married	75K	No
10	No	Single	90K	Yes

- Class: $P(Y) = N_c/N$
 - e.g., $P(\text{No}) = 7/10$, $P(\text{Yes}) = 3/10$

Let's work out all $P(X|Y)$

- For discrete attributes:

$$P(X_i | Y_k) = |X_{ik}| / N_c$$

- where $|X_{ik}|$ is number of instances having attribute value X_i and belonging to class Y_k
- Examples:

$$P(\text{Status}=\text{Married} | \text{No}) = 4/7$$
$$P(\text{Refund}=\text{Yes} | \text{Yes})=0$$

Estimate Probabilities from Data

- For continuous attributes:
 - Discretization: Partition the range into bins:
 - Replace continuous value with bin value
 - Probability density estimation:
 - Assume an attribute follows a normal distribution
 - Use data to estimate parameters of distribution (e.g., mean and standard deviation)
 - Once probability distribution is known, use it to estimate the conditional probability $P(X_i | Y)$

Estimate Probabilities from Data

Tid	Refund	Marital Status	Taxable Income	Evade
1	Yes	Single	125K	No
2	No	Married	100K	No
3	No	Single	70K	No
4	Yes	Married	120K	No
5	No	Divorced	95K	Yes
6	No	Married	60K	No
7	Yes	Divorced	220K	No
8	No	Single	85K	Yes
9	No	Married	75K	No
10	No	Single	90K	Yes

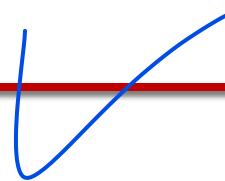
- Normal distribution:

$$P(X_i | Y_j) = \frac{1}{\sqrt{2\pi\sigma_{ij}^2}} e^{-\frac{(X_i - \mu_{ij})^2}{2\sigma_{ij}^2}}$$

- One for each (X_i, Y_i) pair
- For (Income, Class=No):
 - If Class=No
 - sample mean = 110
 - sample variance = 2975

$$P(\text{Income} = 120 | \text{No}) = \frac{1}{\sqrt{2\pi}(54.54)} e^{-\frac{(120-110)^2}{2(2975)}} = 0.0072$$

Example of Naïve Bayes Classifier



Given a Test Record:

$$X = (\text{Refund} = \text{No}, \text{Married}, \text{Income} = 120\text{K})$$

naive Bayes Classifier:

$$P(\text{Refund}=\text{Yes}|\text{No}) = 3/7$$

$$P(\text{Refund}=\text{No}|\text{No}) = 4/7$$

$$P(\text{Refund}=\text{Yes}|\text{Yes}) = 0$$

$$P(\text{Refund}=\text{No}|\text{Yes}) = 1$$

$$P(\text{Marital Status}=\text{Single}|\text{No}) = 2/7$$

$$P(\text{Marital Status}=\text{Divorced}|\text{No}) = 1/7$$

$$P(\text{Marital Status}=\text{Married}|\text{No}) = 4/7$$

$$P(\text{Marital Status}=\text{Single}|\text{Yes}) = 2/3$$

$$P(\text{Marital Status}=\text{Divorced}|\text{Yes}) = 1/3$$

$$P(\text{Marital Status}=\text{Married}|\text{Yes}) = 0$$

For taxable income:

If class=No: sample mean=110
sample variance=2975

If class=Yes: sample mean=90
sample variance=25

- $P(X|\text{Class}=\text{No}) = P(\text{Refund}=\text{No}|\text{Class}=\text{No}) \times P(\text{Married}|\text{ Class}=\text{No}) \times P(\text{Income}=120\text{K}|\text{ Class}=\text{No}) = 4/7 \times 4/7 \times 0.0072 = 0.0024$
- $P(X|\text{Class}=\text{Yes}) = P(\text{Refund}=\text{No}|\text{ Class}=\text{Yes}) \times P(\text{Married}|\text{ Class}=\text{Yes}) \times P(\text{Income}=120\text{K}|\text{ Class}=\text{Yes}) = 1 \times 0 \times 1.2 \times 10^{-9} = 0$

Since $P(X|\text{No})P(\text{No}) > P(X|\text{Yes})P(\text{Yes})$

Therefore $P(\text{No}|X) > P(\text{Yes}|X)$

$\Rightarrow \text{Class} = \text{No}$

SpamAssassin

- Naïve Bayes has found a home in spam filtering
 - Paul Graham’s A Plan for Spam
 - Widely used in spam filters
 - But many features beyond words:
 - black hole lists, etc.
 - particular hand-crafted text patterns

SpamAssassin Features:

- Basic (Naïve) Bayes spam probability
- Mentions: Generic Viagra
- Regex: millions of (dollar) ((dollar) NN,NNN,NNN.NN)
- Phrase: impress ... girl
- Phrase: 'Prestigious Non-Accredited Universities'
- From: starts with many numbers
- Subject is all capitals
- HTML has a low ratio of text to image area
- Relay in RBL, http://www.mail-abuse.com/enduserinfo_rbl.html
- RCVD line looks faked
- http://spamassassin.apache.org/tests_3_3_x.html

Naive Bayes is Not So Naive

- Very fast learning and testing (basically just count words)
- Low storage requirements
- Very good in domains with many equally important features
- More robust to irrelevant features than many learning methods
 - Irrelevant features cancel out without affecting results

Naive Bayes is Not So Naive

- More robust to concept drift (changing class definition over time)
- Naive Bayes won 1st and 2nd place in KDD-CUP 97 competition out of 16 systems
 - Goal: Financial services industry direct mail response prediction: Predict if the recipient of mail will actually respond to the advertisement – 750,000 records.
- A good dependable baseline for text classification (but not the best)!

Evaluating Categorization

- Evaluation must be done on test data that are independent of the training data
 - Sometimes use cross-validation (averaging results over multiple training and test splits of the overall data)
- Easy to get good performance on a test set that was available to the learner during training (e.g., just memorize the test set)

Evaluating Categorization

- Measures: precision, recall, F1, classification accuracy
- Classification accuracy: r/n where n is the total number of test docs and r is the number of test docs correctly classified

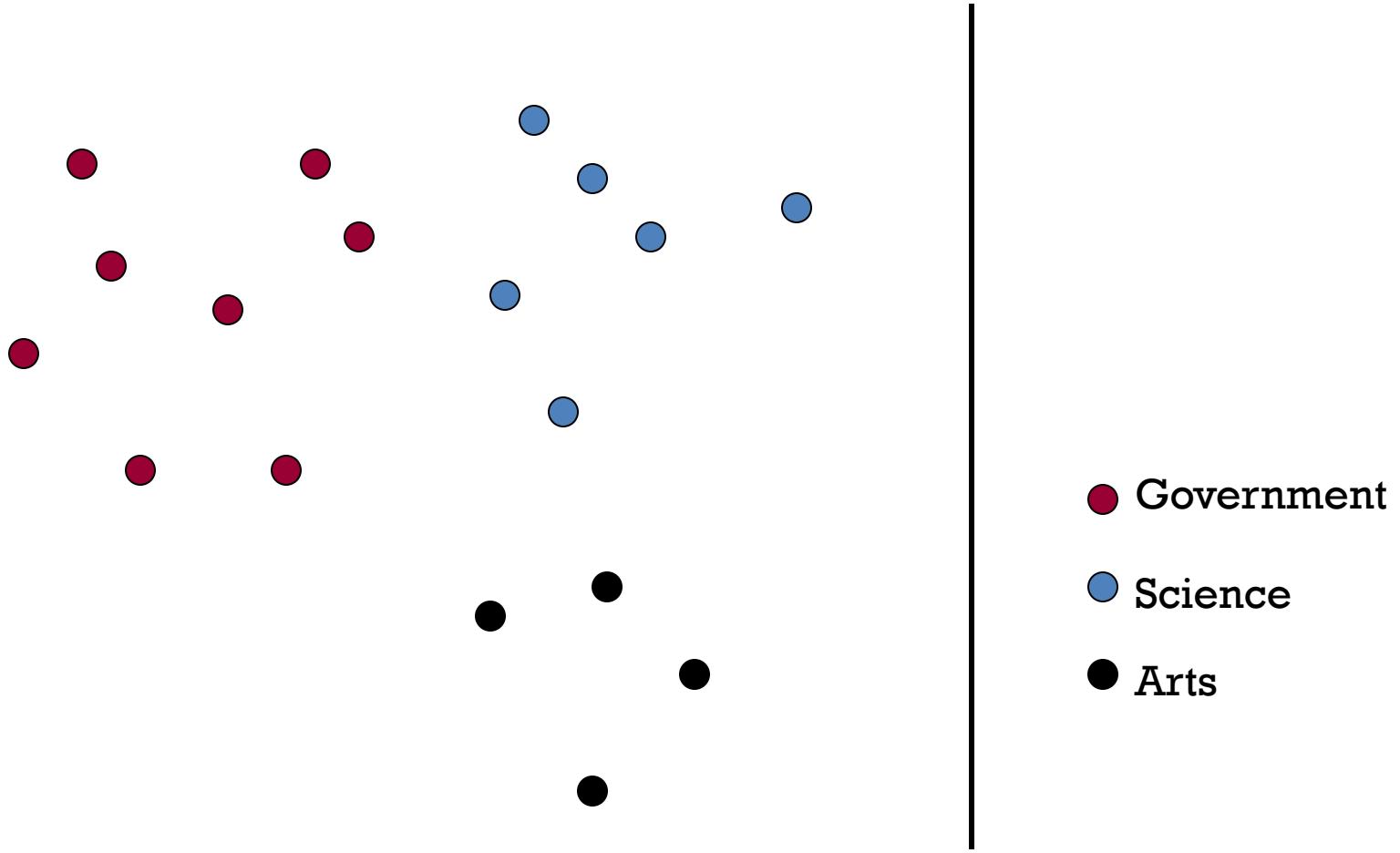
Recall: Vector Space Representation

- Each document is a vector, one component for each term (= word).
- Normally normalize vectors to unit length.
- High-dimensional vector space:
 - Terms are axes
 - 10,000+ dimensions, or even 100,000+
 - Docs are vectors in this space
- How can we do classification in this space?

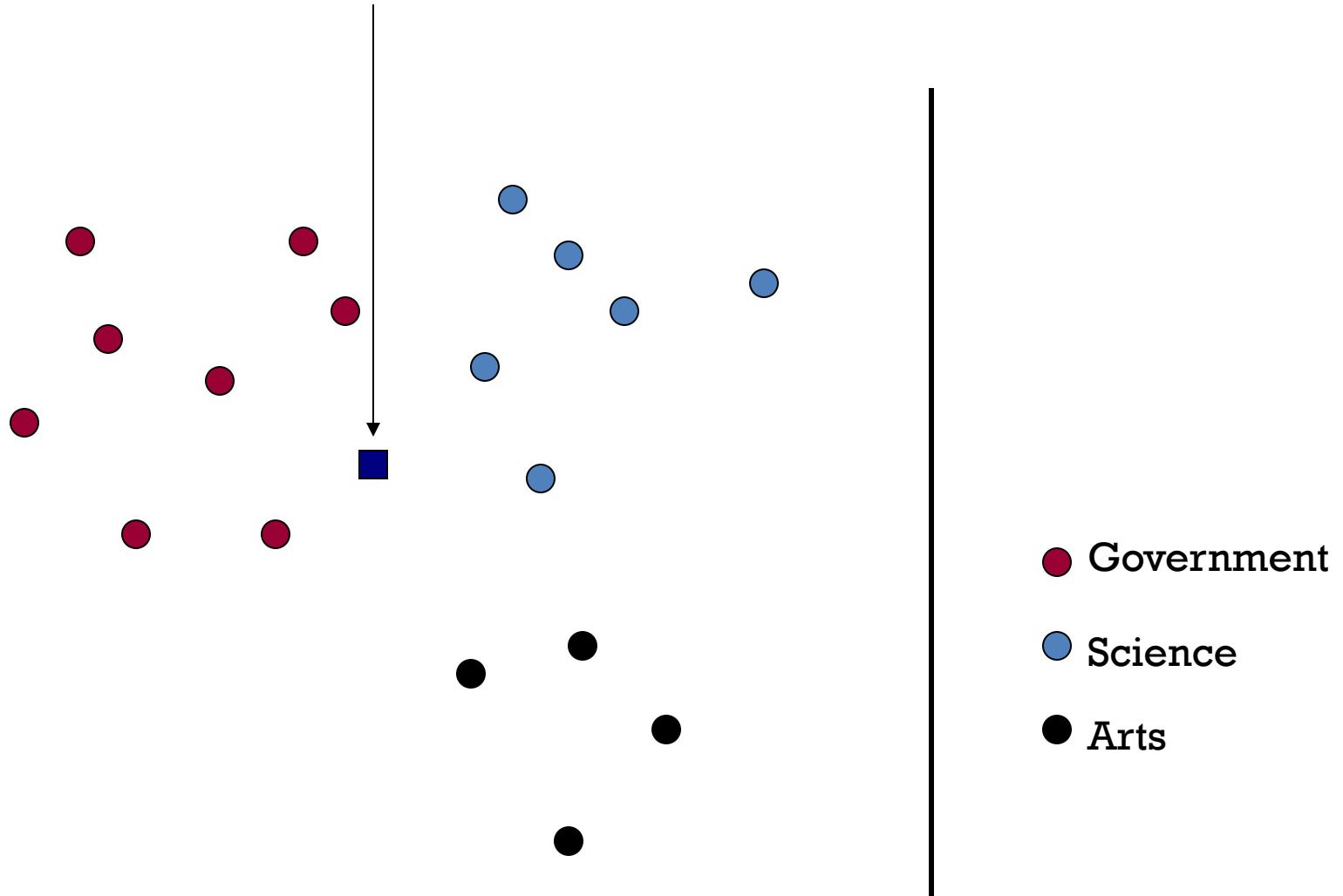
Classification Using Vector Spaces

- In vector space classification, training set corresponds to a labeled set of points (equivalently, vectors)
- Premise 1: Documents in the same class form a contiguous region of space
- Premise 2: Documents from different classes don't overlap (much)
- Learning a classifier: build surfaces to delineate classes in the space

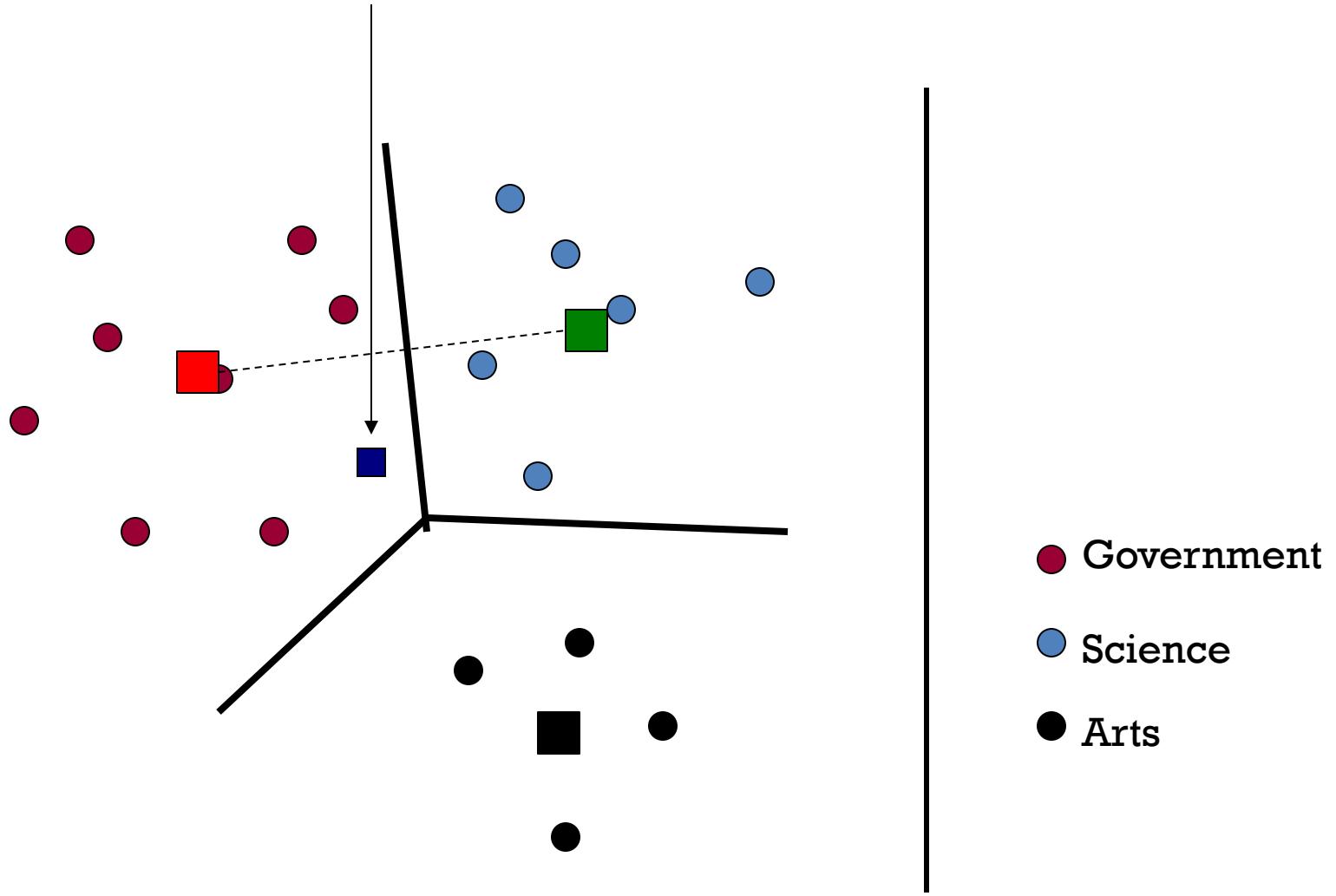
Documents in a Vector Space



Test Document of what class?



Test Document = Government



Our focus: how to find good separators

Definition of centroid

$$\vec{u}(c) = \frac{1}{|D_c|} \sum_{d \in D_c} \vec{v}(d)$$

- Where D_c is the set of all documents that belong to class c and $v(d)$ is the vector space representation of d .
- Note that centroid will in general not be a unit vector even when the inputs are unit vectors.

Rocchio classification

- Rocchio forms a simple representative for each class: the centroid/prototype
- Classification: nearest prototype/centroid
- It does not guarantee that classifications are consistent with the given training data

Why not?

Two-class Rocchio as a linear classifier

- Line or hyperplane defined by:

$$\sum_{i=1}^M w_i d_i = \theta$$

- For Rocchio, set:

$$\vec{w} = \vec{\mu}(c_1) - \vec{\mu}(c_2)$$

$$\theta = 0.5 \times (\|\vec{\mu}(c_1)\|^2 - \|\vec{\mu}(c_2)\|^2)$$

Linear classifier: Example

- Class: “interest” (as in interest rate)
- Example features of a linear classifier
- | w_i | t_i |
|-------------------|---------------|
| • 0.70 prime | • -0.71 dIrs |
| • 0.67 rate | • -0.35 world |
| • 0.63 interest | • -0.33 sees |
| • 0.60 rates | • -0.25 year |
| • 0.46 discount | • -0.24 group |
| • 0.43 bundesbank | • -0.24 dlr |
- To classify, find dot product of feature vector and weights

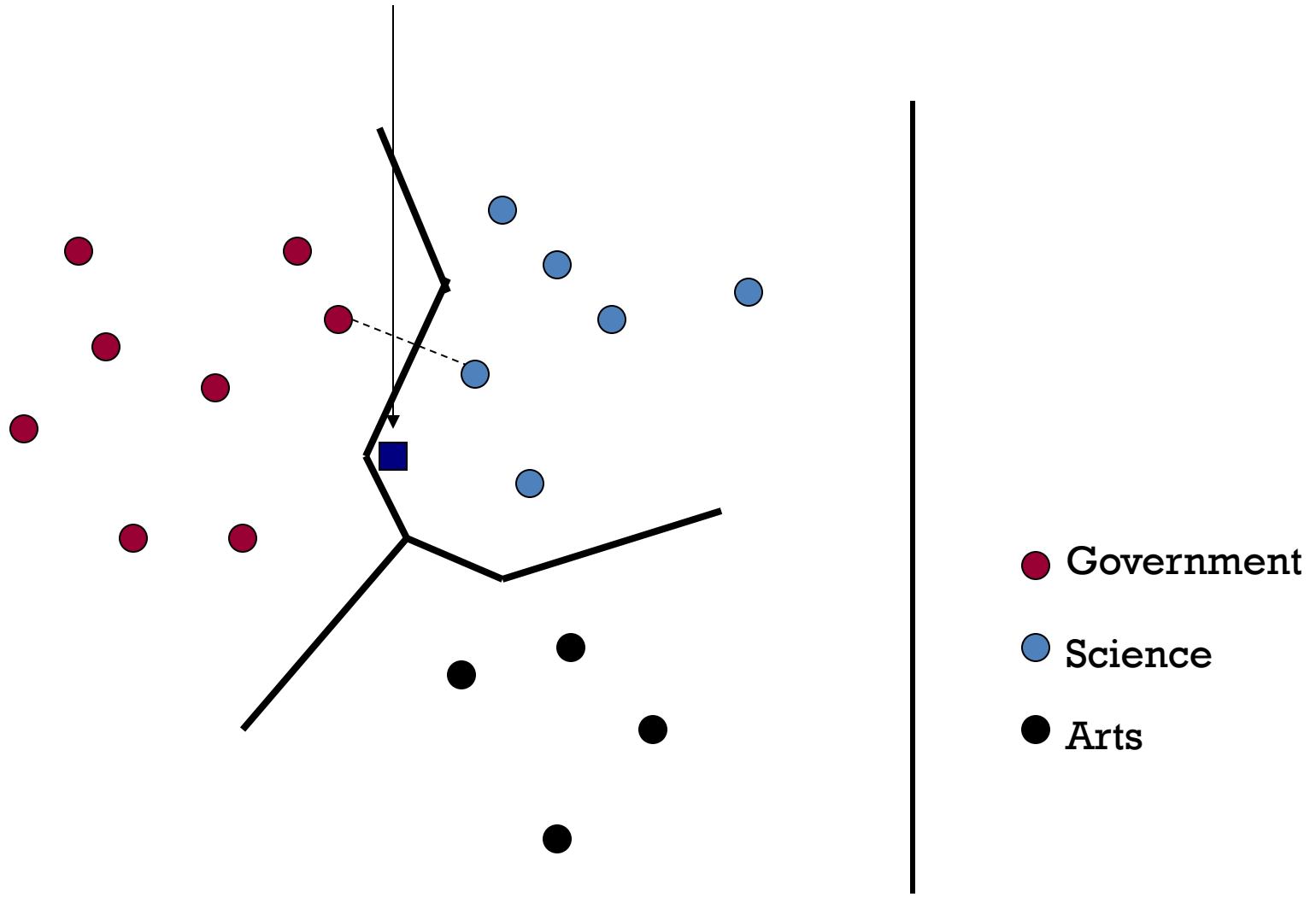
Rocchio classification

- A simple form of Fisher's linear discriminant
- Little used outside text classification
 - It has been used quite effectively for text classification
 - But in general worse than Naïve Bayes
- Again, cheap to train and test documents

k Nearest Neighbor Classification

- $k\text{NN} = k$ Nearest Neighbor
- To classify a document d :
- Define k -neighborhood as the k nearest neighbors of d
- Pick the majority class label in the k -neighborhood
- For larger k can roughly estimate $P(c|d)$ as $\#(c)/k$

Test Document = Science



Voronoi diagram

Nearest-Neighbor Learning

- Learning: just store the labeled training examples D
- Testing instance x (under 1NN):
 - Compute similarity between x and all examples in D .
 - Assign x the category of the most similar example in D .
- Does not compute anything beyond storing the examples
- Also called:
 - Case-based learning
 - Memory-based learning
 - Lazy learning
- Rationale of kNN: contiguity hypothesis

k Nearest Neighbor

- Using only the closest example (1NN) subject to errors due to:
 - A single atypical example.
 - Noise (i.e., an error) in the category label of a single training example.
- More robust: find the k examples and return the majority category of these k
- k is typically odd to avoid ties; 3 and 5 are most common

Nearest Neighbor with Inverted Index

- Naively finding nearest neighbors requires a linear search through $|D|$ documents in collection
- But determining k nearest neighbors is the same as determining the k best retrievals using the test document as a query to a database of training documents.
- Use standard vector space inverted index methods to find the k nearest neighbors.
- Testing Time: $O(B/V_t/)$ where B is the average number of training documents in which a test-document word appears.
 - Typically $B \ll |D|$

kNN: Discussion

- No feature selection necessary
- No training necessary
- Scales well with large number of classes
 - Don't need to train n classifiers for n classes
- Classes can influence each other
 - Small changes to one class can have ripple effect
- Done naively, very expensive at test time
- In most cases it's more accurate than NB or Rocchio

Bias vs. capacity – notions and terminology

- Consider asking a botanist: Is an object a tree?
 - Too much *capacity*, low *bias*
 - Botanist who memorizes
 - Will always say “no” to new object (e.g., different # of leaves)
 - Not enough capacity, high bias
 - Lazy botanist
 - Says “yes” if the object is green
 - You want the middle ground

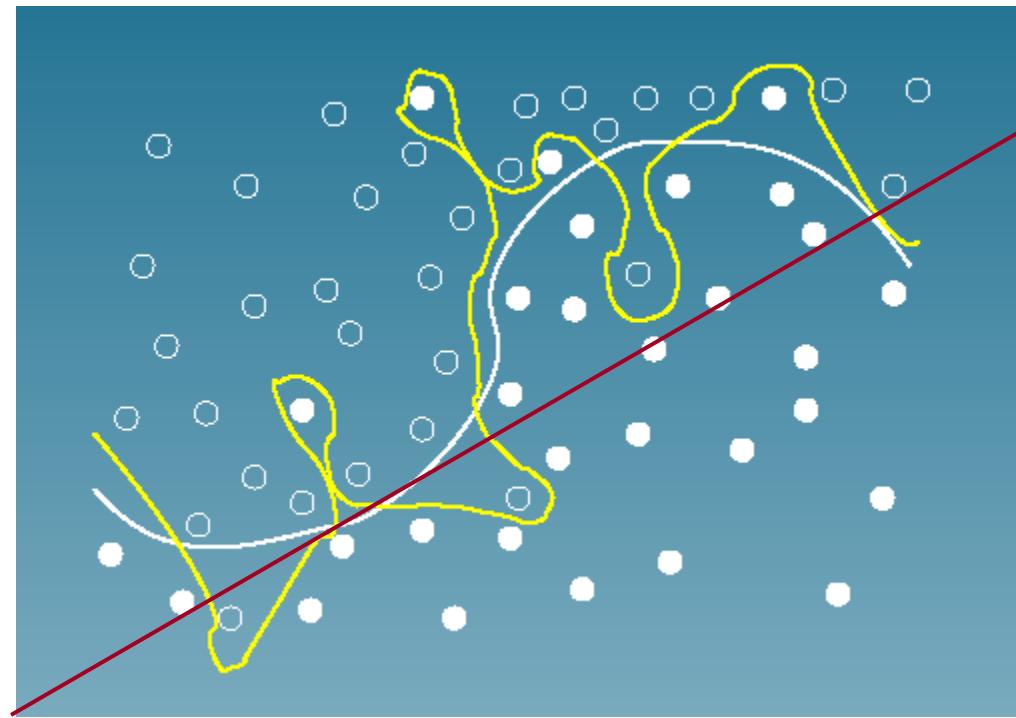
(Example due to C. Burges)

kNN vs. Naive Bayes

- Bias/Variance tradeoff
 - Variance \approx Capacity
- kNN has high variance and low bias.
 - Infinite memory
- Rocchio/NB has low variance and high bias.
 - Linear decision surface between classes

Bias vs. variance:

Choosing the correct model capacity



Summary: Representation of Text Categorization Attributes

- Representations of text are usually very high dimensional
 - “The curse of dimensionality”
- High-bias algorithms should generally work best in high-dimensional space
 - They prevent overfitting
 - They generalize more
- For most text categorization tasks, there are many relevant features and many irrelevant ones

Which classifier do I use for a given text classification problem?

- Is there a learning method that is optimal for all text classification problems?
- No, because there is a tradeoff between bias and variance.
- Factors to take into account:
 - How much training data is available?
 - How simple/complex is the problem? (linear vs. nonlinear decision boundary)
 - How noisy is the data?
 - How stable is the problem over time?
 - For an unstable problem, it's better to use a simple and robust classifier.