

Introduction to Information Retrieval

<http://informationretrieval.org>

IIR 15-1: Support Vector Machines

Wiltrud Kessler & Hinrich Schütze

Institute for Natural Language Processing, Universität Stuttgart

2012-12-04

Overview

- 1 Recap
- 2 Support Vector Machines
- 3 Discussion

Outline

- 1 Recap
- 2 Support Vector Machines
- 3 Discussion

Relevance

- We will evaluate the quality of an information retrieval system and, in particular, its ranking algorithm with respect to **relevance**.
- A document is relevant if it gives the user the information she was looking for.
- To evaluate relevance, we need an **evaluation benchmark** with three elements:
 - A benchmark document collection
 - A benchmark suite of queries
 - An assessment of the relevance of each query-document pair

Relevance: query vs. information need

- The notion of “relevance to the query” is very problematic.
- **Information need i** : You are looking for information on whether drinking red wine is more effective at reducing your risk of heart attacks than white wine.
- **Query q** : WINE RED WHITE HEART ATTACK
- Consider document d' : *He then launched into the heart of his speech and attacked the wine industry lobby for downplaying the role of red and white wine in drunk driving.*
- d' is relevant to the query q , but d' is **not** relevant to the information need i .
- User happiness/satisfaction (i.e., how well our ranking algorithm works) can only be measured **by relevance to information needs, not by relevance to queries.**

Precision and recall

- Precision (P) is the fraction of retrieved documents that are relevant

$$\text{Precision} = \frac{\#(\text{relevant items retrieved})}{\#(\text{retrieved items})} = P(\text{relevant}|\text{retrieved})$$

- Recall (R) is the fraction of relevant documents that are retrieved

$$\text{Recall} = \frac{\#(\text{relevant items retrieved})}{\#(\text{relevant items})} = P(\text{retrieved}|\text{relevant})$$

A combined measure: F

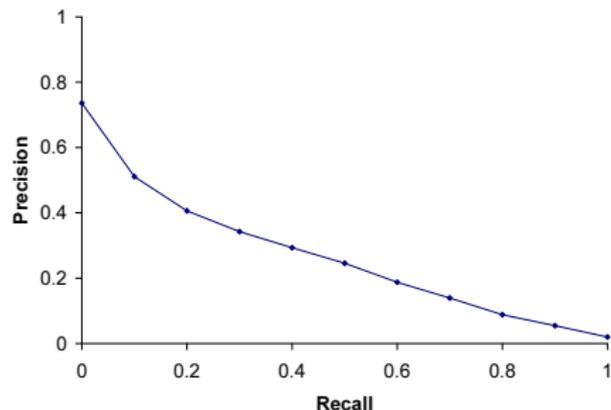
- F allows us to trade off precision against recall.

- Balanced F :

$$F_1 = \frac{2PR}{P + R}$$

- This is a kind of soft minimum of precision and recall.

Averaged 11-point precision/recall graph



- This curve is typical of performance levels for the TREC benchmark.
- 70% chance of getting the first document right (roughly)
- When we want to look at at least 50% of all relevant documents, then for each relevant document we find, we will have to look at about two nonrelevant documents.
- That's not very good.
- High-recall retrieval is an unsolved problem.

Take-away today

- **Support vector machines:** State-of-the-art text classification methods (linear and nonlinear)
- **Discussion:** Which classifier should I use for my problem?

Overview

- 1 Recap
- 2 Support Vector Machines
- 3 Discussion

Outline

- 1 Recap
- 2 Support Vector Machines
- 3 Discussion

Support vector machines

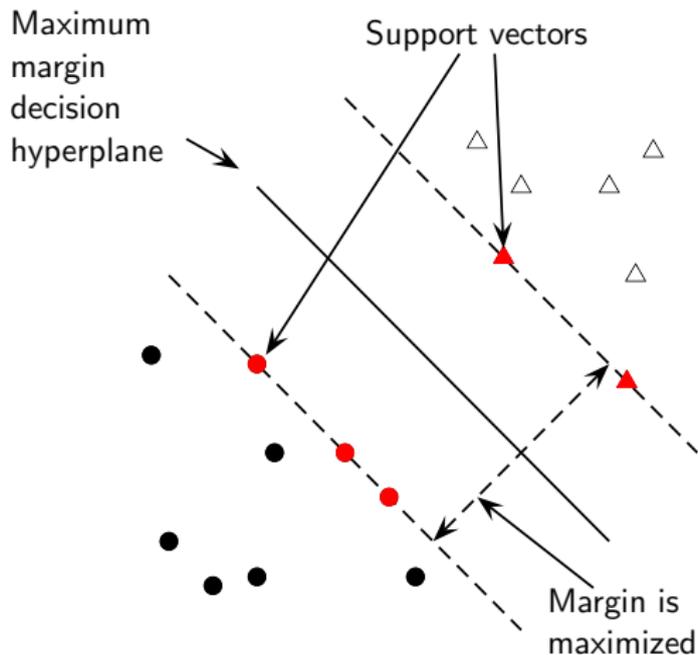
- Machine-learning research in the last two decades has improved classifier effectiveness.
- New generation of state-of-the-art classifiers: support vector machines (SVMs), boosted decision trees, regularized logistic regression, maximum entropy, neural networks, and random forests
- As we saw in IIR: Applications to IR problems, particularly text classification

What is a support vector machine – first take

- Vector space classification (similar to Rocchio, kNN, linear classifiers)
- Difference from previous methods: **large margin** classifier
- We aim to find a separating hyperplane (decision boundary) that is **maximally far** from any point in the training data
- In case of non-linear-separability: We may have to discount some points as outliers or noise.

Support Vector Machines

- 2-class training data
- decision boundary \rightarrow **linear separator**
- criterion: being maximally far away from any data point \rightarrow determines classifier **margin**
- Vectors on margin lines are called **support vectors**
- Set of support vectors are a complete specification of classifier

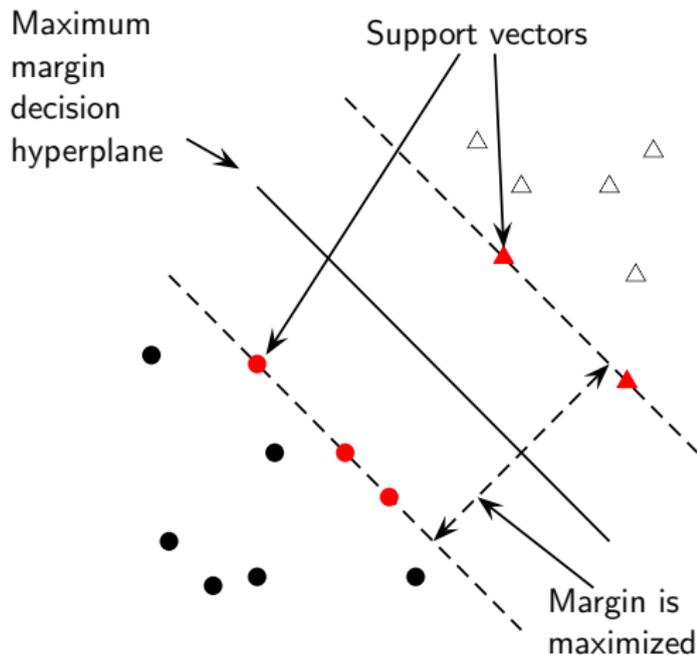


Why maximize the margin?

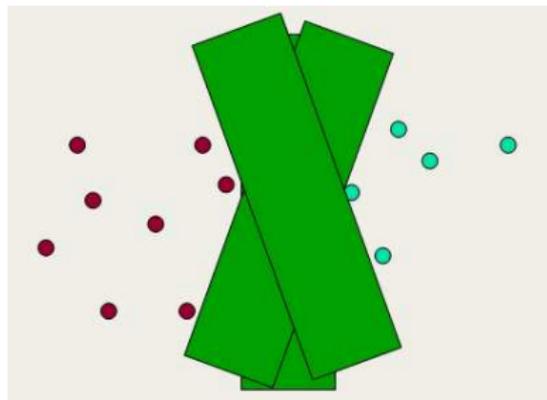
Points near decision surface \rightarrow uncertain classification decisions (50% either way).

A classifier with a large margin makes no low certainty classification decisions.

Gives classification safety margin with respect to errors and random variation



Why maximize the margin?



- SVM classifier: large margin around decision boundary
- compare to decision hyperplane: place fat separator between classes
 - unique solution
- decreased memory capacity
- increased ability to correctly generalize to test data

Separating hyperplane: Recap

Hyperplane

An n -dimensional generalization of a plane (point in 1-D space, line in 2-D space, ordinary plane in 3-D space).

Decision hyperplane

Can be defined by:

- intercept term b (we were calling this θ before)
- normal vector \vec{w} (**weight vector**) which is perpendicular to the hyperplane

All points \vec{x} on the hyperplane satisfy:

$$\vec{w}^T \vec{x} + b = 0$$

Formalization of SVMs

Training set

Consider a binary classification problem:

- \vec{x}_i are the input vectors
- y_i are the labels

For SVMs, the two classes are $y_i = +1$ and $y_i = -1$.

The linear classifier is then:

$$f(\vec{x}) = \text{sign}(\vec{w}^T \vec{x} + b)$$

A value of -1 indicates one class, and a value of $+1$ the other class.

Functional margin of a point

We are confident in the classification of a point if it is far away from the decision boundary.

Functional margin

The functional margin of the vector \vec{x}_i w.r.t the hyperplane $\langle \vec{w}, b \rangle$ is: $y_i(\vec{w}^T \vec{x}_i + b)$

The functional margin of a data set w.r.t a decision surface is twice the functional margin of any of the points in the data set with minimal functional margin

- factor 2 comes from measuring across the whole width of the margin

But we can increase functional margin by scaling \vec{w} and b . We need to place some constraint on the size of the \vec{w} vector.

Geometric margin

Geometric margin of the classifier: maximum width of the band that can be drawn separating the support vectors of the two classes.

$$r = y \frac{\vec{w}^T \vec{x} + b}{|\vec{w}|}$$

The geometric margin is clearly invariant to scaling of parameters: if we replace \vec{w} by $5\vec{w}$ and b by $5b$, then the geometric margin is the same, because it is normalized by the length of \vec{w} .

Optimization problem solved by SVMs

Assume canonical distance

Assume that all data is at least distance 1 from the hyperplane, then:

$$y_i(\vec{w}^T \vec{x}_i + b) \geq 1$$

Since each example's distance from the hyperplane is $r_i = y_i(\vec{w}^T \vec{x}_i + b)/|\vec{w}|$, the margin is $\rho = 2/|\vec{w}|$.

We want to maximize this margin.

That is, we want to find \vec{w} and b such that:

- $\rho = 2/|\vec{w}|$ is maximized
- For all $(\vec{x}_i, y_i) \in \mathbb{D}$, $y_i(\vec{w}^T \vec{x}_i + b) \geq 1$

Optimization problem solved by SVMs (2)

Maximizing $2/|\vec{w}|$ is the same as minimizing $|\vec{w}|/2$. This gives the final standard formulation of an SVM as a minimization problem:

Example

Find \vec{w} and b such that:

- $\frac{1}{2}\vec{w}^T\vec{w}$ is minimized (because $|\vec{w}| = \sqrt{\vec{w}^T\vec{w}}$), and
- for all $\{(\vec{x}_i, y_i)\}$, $y_i(\vec{w}^T\vec{x}_i + b) \geq 1$

We are now optimizing a **quadratic function** subject to linear constraints. Quadratic optimization problems are standard mathematical optimization problems, and many algorithms exist for solving them (e.g. Quadratic Programming libraries).

Recap

- We start with a training set.
- The data set defines the maximum-margin separating hyperplane (if it is separable).
- We use quadratic optimization to find this plane.
- Given a new point \vec{x} to classify, the classification function $f(\vec{x})$ computes the projection of the point onto the hyperplane normal.
- The sign of this function determines the class to assign to the point.
- If the point is within the margin of the classifier, the classifier can return “don’t know” rather than one of the two classes.
- The value of $f(\vec{x})$ may also be transformed into a probability of classification

Soft margin classification

What happens if data is not linearly separable?

- Standard approach: allow the fat decision margin to make a few mistakes
 - some points, outliers, noisy examples are inside or on the wrong side of the margin
- Pay cost for each misclassified example, depending on how far it is from meeting the margin requirement

Slack variable ξ_i : A non-zero value for ξ_i allows \vec{x}_i to not meet the margin requirement at a cost proportional to the value of ξ_i .

Optimization problem: trading off how fat it can make the margin vs. how many points have to be moved around to allow this margin.

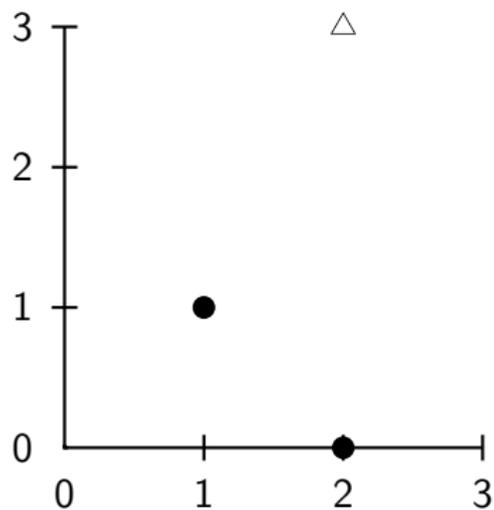
The sum of the ξ_i gives an upper bound on the number of training errors.

Soft-margin SVMs minimize training error traded off against margin.

Using SVM for one-of classification

- Recall how to use binary linear classifiers (k classes) for one-of: train and run k classifiers and then select the class with the highest confidence
- Another strategy used with SVMs: build $k(k - 1)/2$ one-versus-one classifiers, and choose the class that is selected by the most classifiers. While this involves building a very large number of classifiers, the time for training classifiers may actually decrease, since the training data set for each classifier is much smaller.
- Yet another possibility: structured prediction. Generalization of classification where the classes are not just a set of independent, categorical labels, but may be arbitrary structured objects with relationships defined between them

Exercise



Which vectors are the support

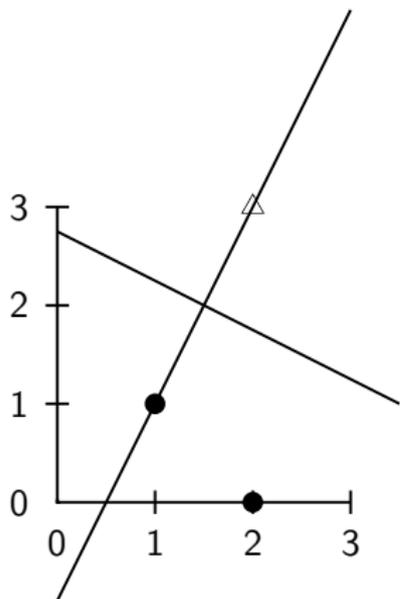
vectors? Draw the maximum margin separator. What values of w_1 , w_2 and b (for $w_1x + w_2y + b = 0$) describe this separator? Recall that we must have $w_1x + w_2y + b \in \{1, -1\}$ for the support vectors.

Walkthrough example: building an SVM over the data set shown in the figure

Working geometrically:

- The maximum margin weight vector will be parallel to the shortest line connecting points of the two classes, that is, the line between $(1, 1)$ and $(2, 3)$, giving a weight vector of $(1, 2)$.
- The optimal decision surface is orthogonal to that line and intersects it at the halfway point. Therefore, it passes through $(1.5, 2)$.
- The SVM decision boundary is:

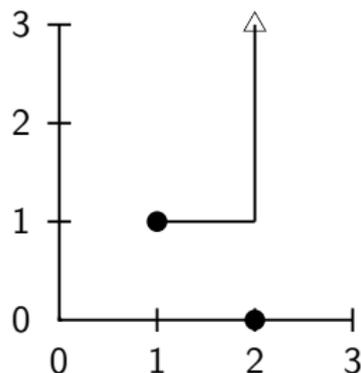
$$0 = x + 2y - (1 \cdot 1.5 + 2 \cdot 2) \Leftrightarrow 0 = \frac{2}{5}x + \frac{4}{5}y - \frac{11}{5}$$



Walkthrough example: building an SVM over the data set shown in the figure

Working algebraically:

- With the constraint $\text{sign}(y_i(\vec{w}^T \vec{x}_i + b)) \geq 1$, we seek to minimize $|\vec{w}|$.
- We know that the solution is $\vec{w} = (a, 2a)$ for some a . So:
 $a + 2a + b = -1$, $2a + 6a + b = 1$
- Hence, $a = 2/5$ and $b = -11/5$. So the optimal hyperplane is given by $\vec{w} = (2/5, 4/5)$ and $b = -11/5$.
- The margin ρ is $2/|\vec{w}| = 2/\sqrt{4/25 + 16/25} = 2/(2\sqrt{5}/5) = \sqrt{5} = \sqrt{(1-2)^2 + (1-3)^2}$.



Outline

- 1 Recap
- 2 Support Vector Machines
- 3 Discussion

Text classification

- Many commercial applications
- There are many applications of text classification for corporate Intranets, government departments, and Internet publishers.
- Often greater performance gains from exploiting domain-specific text features than from changing from one machine learning method to another.
- Understanding the data is one of the keys to successful categorization, yet this is an area in which many categorization tool vendors are weak.

Choosing what kind of classifier to use

When building a text classifier, first question: **how much training data is there currently available?**

Practical challenge: creating or obtaining enough training data

Hundreds or thousands of examples from each class are required to produce a high performance classifier and many real world contexts involve large sets of categories.

- None?
- Very little?
- Quite a lot?
- A huge amount, growing every day?

If you have no labeled training data

Use hand-written rules

Example

IF (wheat OR grain) AND NOT (whole OR bread) THEN
 $c = \text{grain}$

In practice, rules get a lot bigger than this, and can be phrased using more sophisticated query languages than just Boolean expressions, including the use of numeric scores. With careful crafting, the accuracy of such rules can become very high (high 90% precision, high 80% recall). Nevertheless the amount of work to create such well-tuned rules is very large. A reasonable estimate is 2 days per class, and extra time has to go into maintenance of rules, as the content of documents in classes drifts over time.

A Verity topic (a complex classification rule)

```
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          * 0.70 film ACCRUE
                  ** 0.50 STEM
                  /wordtext = film
evidencetopic    ** 0.50 WORD
topic definition modifier /wordtext = ballet
subtopic          ** 0.50 motion-picture PHRAS
evidencetopic    *** 1.00 WORD
                  /wordtext = motion
topic definition modifier /wordtext = dance
evidencetopic    *** 1.00 WORD
                  /wordtext = picture
topic definition modifier /wordtext = opera
evidencetopic    ** 0.50 STEM
                  /wordtext = movie
topic definition modifier /wordtext = symphony
subtopic          * 0.50 video ACCRUE
evidencetopic    ** 0.50 STEM
                  /wordtext = video
subtopic          ** 0.50 WORD
                  /wordtext = painting
                  ** 0.50 STEM
                  /wordtext = vcr
                  # End of art topic
```

Westlaw: Example queries

Information need: Information on the legal theories involved in preventing the disclosure of trade secrets by employees formerly employed by a competing company *Query:* "trade secret" /s disclos! /s prevent /s employe! *Information need:* Requirements

for disabled people to be able to access a workplace *Query:* disab! /p access! /s work-site work-place (employment /3 place)

Information need: Cases about a host's responsibility for drunk guests *Query:* host! /p (responsib! liab!) /p (intoxicat! drunk!) /p guest

If you have fairly little data and you are going to train a supervised classifier

Work out how to get more labeled data as quickly as you can.

- Best way: insert yourself into a process where humans will be willing to label data for you as part of their natural tasks.

Example

Often humans will sort or route email for their own purposes, and these actions give information about classes.

Active Learning

A system is built which decides which documents a human should label.

Usually these are the ones on which a classifier is uncertain of the correct classification.

If you have labeled data

Good amount of labeled data, but not huge

Use everything that we have presented about text classification.
Consider hybrid approach (overlay Boolean classifier)

Huge amount of labeled data

Choice of classifier probably has little effect on your results.
Choose classifier based on the scalability of training or runtime efficiency. **Rule of thumb: each doubling of the training data size produces a linear increase in classifier performance, but with very large amounts of data, the improvement becomes sub-linear.**

Large and difficult category taxonomies

If you have a small number of well-separated categories, then many classification algorithms are likely to work well. But often: very large number of very similar categories.

Example

Web directories (e.g. the Yahoo! Directory consists of over 200,000 categories or the Open Directory Project), library classification schemes (Dewey Decimal or Library of Congress), the classification schemes used in legal or medical applications.

Accurate classification over large sets of closely related classes is **inherently difficult**. – No general high-accuracy solution.

Recap

- 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 problem?
 - How stable is the problem over time?
 - For an unstable problem, it's better to use a simple and robust classifier.

Take-away today

- **Support vector machines:** State-of-the-art text classification methods (linear and nonlinear)
- **Discussion:** Which classifier should I use for my problem?

Resources

- Chapter 14 of IIR (basic vector space classification)
- Chapter 15 of IIR (SVMs)
- Resources at <http://ifnlp.org/ir>