# Information Retrieval
## WS 2012 / 2013

Lecture 11, Wednesday January 23rd, 2013
(SVMs = Support Vector Machines)

Prof. Dr. Hannah Bast
Chair of Algorithms and Data Structures
Department of Computer Science
University of Freiburg

# Overview of this lecture

- **Organizational**

    - Your results + experiences with Ex. Sheet 10 (Naïve Bayes)

    - Oral exams are on **March 5 + 6** (Tuesday + Wednesday)

- **Support Vector Machines (SVMs)**

    - Another linear classifier, just like Naïve Bayes

    - But different objective function + harder to optimize

    - Some more linear algebra … you will love it

    - Play around with SVM Light software

    - **Exercise Sheet 11:** Compare SVMs with Naïve Bayes with respect to linear separability and classification accuracy

■ Summary / excerpts        last checked January 23, 14:45

   – Theory was clear and not too hard to implement

   – Confusion about multiple labels and choice of training set

   – Great observation: better compare the **log** $\Pr(C = c \mid doc)$

    Reason: They easily become $\leq -1000$ , and the $\exp(\ldots)$ of any such value is $0$ on a typical machine, and all such classes then become indistinguishable  → hurts prediction quality badly

   – Ignoring stop-words helps a bit, but not much

   – Another promising idea from you: consider only words that strongly discriminate between classes in the training set

# Your results for ES#10　(Naive Bayes)

- **For our dataset**　(38.115 docs, 18 classes)

  - Reading time:　on the order of 10 seconds

  - Training time:　on the order of 1 second

  - Prediction time:　on the order of a few seconds

  - **Bottom line 1:**　Naive Bayes is definitely efficient !

  - Quality around 50%

  - With non-exp-trick 60% and more

  - **Bottom line 2:**　Without having seen other methods, it's hard to tell whether this is good or bad or so-so
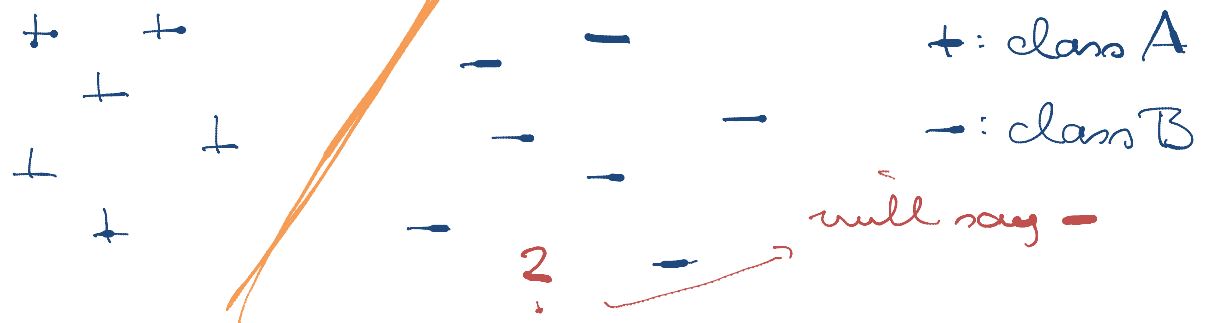
- ■ Informally

    - – Assume the objects are points in d dimensions

    - – Let's assume we have only two classes for now

    - – A linear classifier tries to separate the data points by a (d-1)-dimensional **hyperplane** … definition on next slide

        For d = 2 this means:  try to separate by a **straight line**

    - – Predictions are made based on which side of the hyperplane / straight line the object lies on

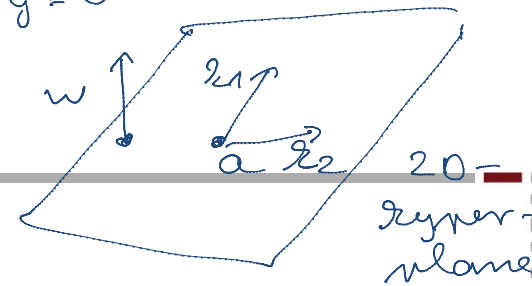    - – Note: points in the training set may not be separable

$d = 2$

$+ : $ class A

$- : $ class B

will say $-$

2

$x \perp y := x$ orthogonal to $y \iff x \cdot y = 0$

$d = 3$

$w_1, z_1, \ldots, z_{d-1}$ form a base of $\mathbb{R}^d$

hyper-plane

- **Formal definition of a hyperplane**

  - A hyperplane H in $\mathbf{R}^d$ if defined by an anchor point $a \in \mathbf{R}^d$, and linearly independent $h_1, \ldots, h_{d-1}$ and consists of all linear combinations $a + \Sigma_i \, a_i \, h_i$ for arbitrary $a_1, \ldots, a_{d-1} \in \mathbf{R}$

  - **Lemma:**  For each such H, there exists $w \in \mathbf{R}^d$ and $b \in \mathbf{R}$ such that  $H = \{ \, x \in \mathbf{R}^d : w \bullet x = b \, \}$   $, w \perp z_1, \ldots, z_{d-1}$

We love to show   $x \in H \iff w \cdot x = b$

"$\Rightarrow$":  $x \in H \Rightarrow x = a + \Sigma \alpha_i \cdot z_i$   $\underset{=0 \text{ because } w \perp z_i}{\underbrace{}}$

$\Rightarrow w \bullet x = w \bullet a + \Sigma \alpha_i \, w \bullet z_i = w \cdot a =: b$

"$\Leftarrow$":  $x \in \mathbb{R}^d, \, w \bullet x = b$, write $x - a = \alpha \cdot w + \Sigma \alpha_i \cdot z_i$

$\Rightarrow w \bullet (a + \alpha \cdot w + \Sigma \alpha_i z_i) = b$   (because $w, z_1, \ldots, z_{d-1}$ form a base)

$\underset{w \bullet a}{\underbrace{w \bullet a}} + \alpha \cdot |w|^2 + \underset{=0}{\underline{\Sigma \alpha_i z_i w}}$   $\Rightarrow \alpha \cdot |w|^2 = 0$

$\Rightarrow \alpha = 0$

■ Distance from a point to a hyperplane

  – Let $H = \{ x \in \mathbf{R}^d : w \bullet x = b \}$ be a hyperplane in $\mathbf{R}^d$

  – Then the distance of a point $x \in \mathbf{R}^d$ to $H$ is $|w \bullet x - b| \,/\, |w|$

  – The sign of $w \bullet x - b$ says on which side of $H$ lies $x$



Handwritten notes:

$$\text{Let } w_o = w / |w| \implies |w_o| = 1$$

$r = \text{dist}(x, H)$

[The picture is for CASE 1]

CASE 1: $x' = x + r \cdot w_o$, $r > 0$
($w$ points from $x$ towards $H$)

$$x' \in H \implies w \bullet x' = w \bullet x + r \cdot w \bullet w_o = b$$
$$\qquad = |w| \cdot w_o \bullet w_o = |w|$$

$$\implies r = -\frac{w \bullet x - b}{|w|}$$

CASE 2: $x' = x - r \cdot w_o$, $r > 0$
($w$ points from $x$ away from $H$)

$$\implies r = \frac{w \bullet x - b}{|w|}$$

# Linear Classifiers   4/6

■ Two-class Naïve Bayes (**NB**) is a linear classifier

  – Recall how **NB** predicts the probability of a class C for d

    $Pr(C \mid d) = Pr(C) \cdot \Pi_{i=1,\ldots,|d|} \, Pr(w_i \mid C)$,  $|d|$ = #words in d

    where $w_i$ is the i-th word of d

  – We can equivalently write this as

    $Pr(C \mid d) = Pr(C) \cdot \Pi_{i=1,\ldots,|V|} \, Pr(w_i \mid C)^{fi}$,  $V$ = vocabulary

    where $w_i$ is the i-th word in V, and fi = #occ of $w_i$ in d

  – **Lemma:**  For two classes A and B, define $b \in \mathbf{R}$ and $w \in \mathbf{R}^{|V|}$

    $b = -\log(Pr(A) / Pr(B))$,  $w_i = \log(Pr(w_i \mid A) / Pr(w_i \mid B))$

    Then **NB** predicts A for x if $w \bullet x - b > 0$, and B otherwise

8

■ **Proof of Lemma**

$x = (f_1, \ldots, f_{|V|})^T$
$= $ feature vector for doc

– **NB** predicts A for x if w • x – b > 0, and B otherwise

$b = -\log(Pr(A) / Pr(B)), \quad w_i = \log(Pr(w_i \mid A) / Pr(w_i \mid B))$

$$Pr(A \mid doc) = Pr(A) \cdot \prod_{i=1}^{|V|} Pr(w_i \mid A)^{f_i} / P$$

$$Pr(B \mid doc) = Pr(B) \cdot \prod_{i=1}^{|V|} Pr(w_i \mid B)^{f_i} / P$$

$$\log \frac{Pr(A \mid doc)}{Pr(B \mid doc)} = \underbrace{\log \frac{Pr(A)}{Pr(B)}}_{=-b} + \sum_{i=1}^{|V|} f_i \cdot \underbrace{\log \frac{Pr(w_i \mid A)}{Pr(w_i \mid B)}}_{= w_i}$$

$$= w \bullet x - b$$

$$Pr(A \mid doc) > Pr(B \mid doc)$$
$$\iff \log \frac{Pr(A \mid doc)}{Pr(B \mid doc)} > 0$$
$$\iff w \bullet x - b > 0$$

$$x > y \iff \frac{x}{y} > 1$$
$$\iff \log \frac{x}{y} > 0$$

log

# Linear Classifiers   6/6

- The toy example from our last lecture again:

| Doc 1: | aba | class A |
|--------|-----|---------|
| Doc 2: | baabaaa | class A |
| Doc 3: | bbaabbab | class B |
| Doc 4: | abbaa | class A |
| Doc 5: | abbb | class B |
| Doc 6: | bbbaab | class B |

$$m_A = 3 \; , \; m_B = 3$$
$$m_{aA} = 10 \; , \; m_{aB} = 6$$
$$m_{bA} = 5 \; , \; m_{bB} = 12$$
$$Pr(A) = \tfrac{1}{2} \quad Pr(B) = \tfrac{1}{2}$$
$$Pr(a|A) = \tfrac{2}{3} \; ; \; Pr(a|B) = \tfrac{1}{3}$$
$$Pr(b|A) = \tfrac{1}{3} \; ; \; Pr(b|B) = \tfrac{2}{3}$$

new doc $m_a \times a \, , \, m_b \times b$

$$Pr(A|doc) = \tfrac{1}{2} \cdot \left(\tfrac{2}{3}\right)^{m_a} \cdot \left(\tfrac{1}{3}\right)^{m_b} / P$$

$$Pr(B|doc) = \tfrac{1}{2} \cdot \left(\tfrac{1}{3}\right)^{m_a} \cdot \left(\tfrac{2}{3}\right)^{m_b} / P$$

$$\frac{Pr(A|doc)}{Pr(B|doc)} = \left(\tfrac{2}{3}\right)^{m_a - m_b} \cdot \underbrace{\left(\tfrac{1}{3}\right)^{m_b - m_a}}_{=3^{\,m_a - m_b}}$$

$$= 2^{\,m_a - m_b} > 1 \iff m_a > m_b$$

w and b from the proof

"a"
$$w_1 = \log_2 \frac{Pr(a|A)}{Pr(b|A)} = 1$$

"b"
$$w_2 = \log_2 \frac{Pr(a|B)}{Pr(b|B)} = -1$$

$$b = -\log_2 \frac{Pr(A)}{Pr(B)} = 0$$

$$\binom{m_a}{m_b} \cdot \binom{1}{-1} = m_a - m_b$$

$$\underbrace{\phantom{\binom{m_a}{m_b}}}_{=x} \quad \underbrace{\phantom{\binom{1}{-1}}}_{=w} \quad \lessgtr 0$$

10

Doc 1: aba      class A      $(2,1)$

Doc 2: baabaaa      class A      $(5,2)$

Doc 3: bbaabbab      class B      $(3,5)$

Doc 4: abbaa      class A      $(3,2)$

Doc 5: abbb      class B      $(1,3)$

Doc 6: bbbaab      class B      $(2,4)$

let's consider these as points in the plane

■ : class A

△ : class B

$$H = \left\{ x : \begin{pmatrix} 1 \\ -1 \end{pmatrix} \cdot x = 0 \right\}$$

$\underbrace{\quad}_{w}$    $\underbrace{\quad}_{b}$

$x_1 - x_2$

$|w| = \frac{4}{3} \cdot \frac{1}{\sqrt{2}}$

$= \frac{2}{3}\sqrt{2}$

$= 0.9428$

$\frac{2}{|w|} = \frac{3}{2}\sqrt{2}$

#b's

H

$\sqrt{2}$

$\frac{1}{2\sqrt{2}}$

#a's

1   2   3   4   5

11

- **Intuition**

  - Place the separating hyperplane H such that on both sides, there is a **margin** r > 0 as large as possible to the points

  - In $\mathbf{R}^2$ this means: try to separate the point sets with not just a line, but a "band" of width 2r, with r > 0 as large as possible

  - Points on the margin boundary are called **support vectors**

# Support Vector Machines  2/7

- **Derivation of formal optimization problem**

  - Let $x_1, ..., x_m \in \mathbf{R}^d$ be the objects from the training set

  - Let $y_i = +1$ if $x_i$ is in class A,  $y_i = -1$ if $x_i$ is in class B

  - Let $H = \{ \, x$ in $\mathbf{R}^d : w \bullet x = b \, \}$ be a separating hyperplane, such that $w \bullet x_i - b > 0$ for $x_i$ from A, and $< 0$ for $x_i$ from B

  - Then  $\text{dist}(x_i, H) = y_i \cdot (w \bullet x_i - b) \, / \, |w|$     (see slide 7)

  - This gives rise to the following maximization problem:

    Maximize $2r$, such that  $y_i \cdot (w \bullet x_i - b) \, / \, |w| \geq r$  for all i

  - We can equivalently formulate this as ... proof on next slide

    Minimize $|w|^2$, such that $y_i \cdot (w \bullet x_i - b) \geq 1$  for all i

  - This is a well-known kind of optimization problem ... slide 14

■ Proof of equivalence of

$\frac{2}{|w|} = 2r$ is the width of the margin / band (you need that for the exercise)

  – Maximize $2r$, such that $y_i \cdot (w \bullet x_i - b) / |w| \geq r$ for all $i$

  – Minimize $|w|^2$, such that $y_i \cdot (w \bullet x_i - b) \geq 1$ for all $i$

$$\max 2r \quad \text{s.t.} \quad y_i (w \bullet x_i - b) \geq r \cdot |w| \quad \forall i$$

$\alpha = r \cdot |w|$

$$\iff \max \frac{2\alpha}{|w|} \quad \text{s.t.} \quad y_i (w \cdot x_i - b) \geq \alpha \quad \forall i$$

Observe: if $w, b, \alpha$ is optimum, so is $\frac{w}{\alpha}, \frac{b}{\alpha}, 1$

$$y_i(w x_i - b) \geq \alpha \iff y_i\left(\frac{w}{\alpha} x_i - \frac{b}{\alpha}\right) \geq 1$$

and objective $\frac{2 \cdot 1}{|w/\alpha|} = \frac{2\alpha}{|w|}$

so just take $\alpha = 1$ and just optimize over $w$ and $b$

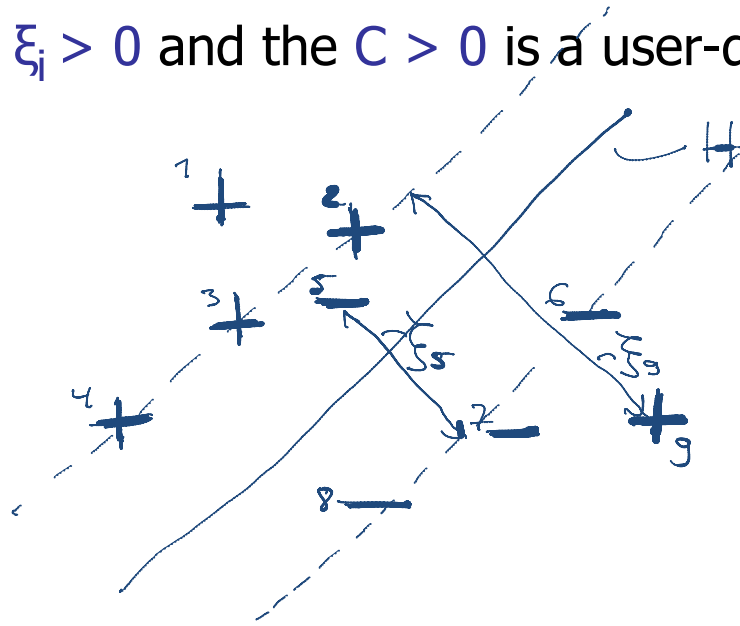$$\max \frac{2}{|w|} \iff \min \frac{|w|}{2} \iff \min |w|^2$$

14

- ■ We now have a quadratic optimization problem

  - The $|w|^2 = w \bullet w$ is a **quadratic** objective function

  - The $y_i \cdot (w \bullet x_i - b) \geq 1$ are **linear** constraints

  - There are established numerical methods for this kind of problem, but the details are beyond the scope of this course

  - Similar as for the SVD, we will use third-party software

- ■ SVM Light Software

  - Solve this optimization problem

  - Download from http://svmlight.joachims.org

  - I will show to download and install it, then let's apply it to our toy example (the 6 documents, with words a and b)

■ So far complete linear separation or nothing

  – The optimization problem can be easily extended to incorporate **outliers** = objects in the training set that lie inside of the margin or even on the wrong side of it:

  Minimize $|w| / 2 + C \cdot \Sigma_i \xi_i$

  such that $y_i \cdot (w \bullet x_i - b) / |w| \geq 1 - \xi_i$  for all $i$

  where $\xi_i > 0$ and the $C > 0$ is a user-defined parameter



In SVM Light,
the C can be set
with the -c
option.
Set to something
very large to
disallow outliers

■ **Multi-Class Support Vector Machines**

– Assume we have an arbitrary number of k classes again

– **Option 1:**  Build k classifiers, one for each class, with the i-th one doing the classification:  Class i  OR  not Class i

Drawback:  Need to "vote" when more than one class wins

– **Option 2:**  Build k · (k − 1) / 2 classifiers, one for each subset of two classes

Drawback:  For large k, that's a lot of classifiers !

– **Option 3:**  Extend the SVM theory to be able to deal with more than two classes directly

Drawback:  optimization problem becomes more complex

- What if the data is not at all linearly separable

  – … even when allowing for a few outliers

  – Standard trick: map objects to a different vector space, where they become (almost) linearly separable again

  – For SVMs, this can be done particularly efficiently, with the so-called "kernel" trick … see machine learning lecture

# References

- **Further reading**

  – Textbook Chapter 15: Support vector machines

    http://nlp.stanford.edu/IR-book/pdf/15svm.pdf

- **Wikipedia**

  – http://en.wikipedia.org/wiki/Linear_classifier

  – http://en.wikipedia.org/wiki/Support_vector_machine

- **SVM Light Software**

  – http://svmlight.joachims.org