Information Retrieval WS 2016 / 2017

Lecture 12, Tuesday January 24th, 2017 (Linear Classifiers, Perceptrons)

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

REIBURG

Overview of this lecture

- Organizational
 - Experiences with ES 10 Naïve Bayes
 - Online evaluation for this course until January 29

- Contents
 - Linear classifiers definition
 - Naïve Bayes (again) is a linear classifier
 - Perceptrons another linear classifier
 - ES12: prove that Naïve Bayes is a linear classifier

Experiences with ES11 1/2

Results

- Overall precision: 73% on genres, 65% on ratings
- Works well for some classes, not so well for others
 Quite well: Documentary (F1 = 86%), G-Rating (F1 = 79%)
 Not so well: Science Fiction (F1 = 46%), PG-13 (F1 = 38%)

Words with largest p_{wc} are mostly stopwords:
Comedy: the and a film by is of in directed comedy was to ...
R: the and a film by is of in was ... horror ... thriller
However: they don't really hurt, because they occur for all classes ... indeed, removing them change results only little

Experiences with ES11 2/2

Which precision is considered "good" ?

- Baseline 1: guess label uniformly at random
 Overall precision: 20% for genres, 25% for ratings
 Pr(label predicted correctly) = 1 / #classes
- Baseline 2: always pick label most frequent in training set
 Overall precision: 41% for genres, 51% for ratings
 Works quite well when one class is very frequent

- Baseline 3: pick label c with probability $p_c = |T_c| / |T|$ This is not better than Baseline 2, and generally worse Think about why ... this might be an exam question

Official course evaluation

Instructions

 You should have received an email from EvaSys Admin on Monday, January 16 with a link to an evaluation form

- We are **very** interested in your feedback
- Please take your time for this
- Please be honest and concrete
- The free text comments are most interesting for us

Please complete by Sunday, January 29

The evaluation is centralized, and will be closed after that date, and there is nothing we can do about that

- Objects are (as usual now) vectors in d dimensions
- Exactly two classes ... often denoted +1 and -1

See slide 10 for how to generalize to more classes

- A linear classifier tries to separate the data points by a (d-1)-dimensional hyperplane, as defined on next slide
 For d=2 this means: try to separate by a straight line
 Note that the points may not be fully separable
- Predictions are made based on which side of the hyperplane (for d=2: straight line) the object lies on

Linear Classifiers 2/6

Hyperplane, definitions

– Two common ways to define a hyperplane H in R^d

The two definitions are equivalent ... proof on next slide

- **Definition 1** (by anchor point and basis):

There is an anchor point $a \in \mathbb{R}^d$ and pairwise orthogonal $h_1, ..., h_{d-1} \in \mathbb{R}^d$ such that H consists of all linear combinat. $a + \Sigma_i \alpha_i h_i$ for arbitrary $\alpha_1, ..., \alpha_{d-1} \in \mathbb{R}$

- **Definition 2** (by normal vector and offset):

There is a normal vector $w \in \mathbb{R}^d$ and an offset $b \in \mathbb{R}$ such that H consists of all points $x \in \mathbb{R}^d$ with $w \bullet x = b$



Linear Classifiers 4/6

Distance from a point to a hyperplane

- Let $H = \{ x \in \mathbb{R}^d : w \bullet x = b \}$ be a hyperplane in \mathbb{R}^d
- Then the distance of a point $x \in \mathbb{R}^d$ to H is $|w \bullet x b| / |w|$

- The sign of $w \cdot x - b$ says on which side of H lies x $w_0 = w/[w] \Rightarrow |w_0| = 1$ $w = w_0 \cdot |w|$ $x = x' + x \cdot w_0 \Rightarrow w \cdot x = w \cdot x' + x \cdot w \cdot w_0$ = b $= x \cdot |w| \cdot |w_0|^2$ $because x'_{GH} = 1$ $\Rightarrow w \cdot x = b + x \cdot |w|$ $\Rightarrow x = (w \cdot x - b)/|w|$ This was CASE 1 (x on "night" side of H) $to v He other case, x = x' - x \cdot w_0$ $\Rightarrow x = -(w \cdot x - b)/|w|$

- Generalization to more classes
 - 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 $k \cdot (k - 1) / 2$ classifiers, one for each subset of two classes

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

 Option 3: Extend theory of the respective approach to deal with more than two classes directly
 Drawback: Sometimes hard (not for Naïve Bayes though)

For ES11 we work with Option 1, but only for one class

When the data is not linearly separable

- Option 1: extend the method to accept "outliers"

Naïve Bayes does this by definition: it always finds some hyperplane, whether the data is linearly separable or not

 Option 2: suitably transform the data to some higherdimensional space, where it becomes (better) separable



Naïve Bayes 1/4

Recap from last lecture

- Let the vocabulary of all words be $V = \{v_1, ..., v_{|V|}\}$

Beware: in todays lecture, **w** is reserved for normal vectors of hyperplanes, which is why we denote a word by **v** here

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

 $Pr(C=c \mid D=d) = \prod_{v_i \text{ in } D} p_{ic} \cdot p_c / Pr(D=d)$

where $p_{ic} = Pr(W=v_i | C=c)$ and the factor is taken multiple times for multiple occurrence of v_i in document d

- We can equivalently write this as

 $Pr(C=c \mid D=d) = \prod_{i=1,...,|V|} p_{ic}^{tf_i} \cdot p_c / Pr(D=d)$ where tf_i is the number of occurrences of v_i in D Two-class NB is a linear classifier

- Assume our two classes are called A and B, and define $b \in \mathbf{R}$ and $w \in \mathbf{R}^{|V|}$ as follows:

 $b = -\log_2(p_A / p_B), w_i = \log_2(p_{iA} / p_{iB})$

Then NB predicts A if and only if $w \bullet x - b > 0$

You should prove this yourself in Exercise 11.1

This is a good exercise for understanding the linear algebra behind linear classifiers. It's not hard, but you have to understand the basic concepts, so perfect exercise :-)

Naïve Bayes 3/4

Our toy example from Lecture 10

- Let us recap the math for a general document



INI REII

14



UNI FREIBURG

Intuition

- A perceptron is a linear classifier that iteratively computes the w and b that define the hyperplane used for prediction
- The w and b are greedily improved in each iteration
- If the training set is linearly separable, the algorithm provably terminates

Perceptrons 2/6

ane round = one pass over all training objects (in random order)

, NI

Algorithm

- Initialization: set w = 0 (all-zero vector) and b = 0
- Then iterate over the objects from the training set in random order, and for each object x do the following:

If $w \bullet x - b$ gives the right prediction, do nothing

If $w \bullet x - b$ gives the wrong prediction, update w and b:

• if
$$w \bullet x - b \le 0$$
: $w \leftarrow w + x$ and $b \leftarrow b - 1$

• if $w \bullet x - b \ge 0$: $w \leftarrow w - x$ and $b \leftarrow b + 1$

Repeat until no more change in w or fixed no. of rounds
 For ES11, find out a good termination condition yourself

First few iterations on our toy example

$x_{1} = (2, 1)$	aba	A
$x_{2} = (5,2)$	baabaaa	A
$x_{3} = (3,5)$	bbaabbab	В
$x_{q} = (3_{l}2)$	abbaa	A
×5= (113)	abbb	В
× ₆ = (2,4)	bbbaab	В

Perceptrons 3/6

 $w = \begin{pmatrix} 0 \\ 0 \end{pmatrix}, \quad b = 0$ $w \cdot x_{1} - b = \begin{pmatrix} 0 \\ 0 \end{pmatrix} \begin{pmatrix} 2 \\ x \end{pmatrix} - 0 = 0$ $w \cdot x_{1} - b = \begin{pmatrix} 0 \\ x \end{pmatrix} \begin{pmatrix} 2 \\ x \end{pmatrix} + x = \begin{pmatrix} 2 \\ x \end{pmatrix}, \quad b \in b - 1 = -1$ $(towork) \land h$ $w \cdot x_{2} - b = \begin{pmatrix} 2 \\ 1 \end{pmatrix} \begin{pmatrix} 5 \\ 2 \end{pmatrix} + 1 = 13 > 0$ $w \cdot x_{3} - b = \begin{pmatrix} 2 \\ 1 \end{pmatrix} \begin{pmatrix} 3 \\ 8 \end{pmatrix} + 1 = -12 > 0$ $w \cdot x_{3} - b = \begin{pmatrix} 2 \\ 1 \end{pmatrix} \begin{pmatrix} 3 \\ 8 \end{pmatrix} + 1 = -12 > 0$ $w \cdot x_{3} - b = \begin{pmatrix} 2 \\ 1 \end{pmatrix} \begin{pmatrix} 3 \\ 8 \end{pmatrix} + 1 = -12 > 0$ $w \cdot x_{3} - b = \begin{pmatrix} 2 \\ 1 \end{pmatrix} \begin{pmatrix} 3 \\ 8 \end{pmatrix} + 1 = -12 > 0$ $w \cdot x_{3} - b = \begin{pmatrix} 2 \\ 1 \end{pmatrix} \begin{pmatrix} 3 \\ 8 \end{pmatrix} + 1 = -12 > 0$ $w \cdot x_{3} - b = \begin{pmatrix} 2 \\ 1 \end{pmatrix} \begin{pmatrix} 3 \\ 8 \end{pmatrix} + 1 = -12 > 0$ $w \cdot x_{3} - b = \begin{pmatrix} 2 \\ 1 \end{pmatrix} \begin{pmatrix} 3 \\ 8 \end{pmatrix} + 1 = -12 > 0$ $w \cdot x_{3} - b = \begin{pmatrix} 2 \\ 1 \end{pmatrix} \begin{pmatrix} 3 \\ 8 \end{pmatrix} + 1 = -12 > 0$ $w \cdot x_{3} - b = \begin{pmatrix} 2 \\ 1 \end{pmatrix} \begin{pmatrix} 3 \\ 8 \end{pmatrix} + 1 = -12 > 0$ $w \cdot x_{3} - b = \begin{pmatrix} 2 \\ 1 \end{pmatrix} \begin{pmatrix} 3 \\ 8 \end{pmatrix} + 1 = -12 > 0$ (towords B)

BURG

INI

and so am ...

w·x-b>0 =1 A

w•x-6<0 =>B

Perceptrons 4/6

$\times \bullet \times = 1 \times l^2$

REIN

Convergence, intuition

 In each iteration, w and b are fixed "towards" the right prediction for the x under consideration:

Assume $w \bullet x - b \le 0$ when it should be > 0

Then we update to w' = w + x and b' = b - 1 $= (w + x) \cdot x$ Then w' • x - b' = w • x - b + $|x|^2 + 1$ This pushes w' • x - b' in the right direction (towards > 0)

The hope is that the various updates (for the various x) do not cancel each other out, and eventually all predictions are correct on the training set

That is, a separating hyperplane is found ... if it exists

Convergence, proof sketch

– Without loss of generality, we can omit the b

Perceptrons $5/6 \frac{d \dim 2}{1 \mod d \dim 2} \begin{pmatrix} w \\ b \end{pmatrix} \cdot \begin{pmatrix} \times \\ -1 \end{pmatrix} = w \cdot x - b$

Just add another dimension d+1, and let the value of all objects be -1 in that dimension ... then w_{d+1} is like the b

b is after called "Intercent" term

INI

- Let $w^{(k)}$ be the w after the k-th correction of w
- Then we can prove that $|w^{(k+1)}| \ge \varepsilon \cdot k$ $= \varepsilon^2 \cdot 2^2 - \varepsilon^2 \cdot 2^2 = c^2 \cdot 2^2 \cdot 2^2 = c^2 \cdot 2^2 = c^2 \cdot 2^2 = c^2 \cdot 2^2 = c^2 \cdot 2^2 \cdot 2^2 \cdot 2^2 \cdot 2^2 = c^2 \cdot 2^2 \cdot 2^2 \cdot 2^2 \cdot 2^2 = c^2 \cdot 2^2 \cdot 2^2 \cdot 2^2 \cdot 2^2 \cdot 2^2 = c^2 \cdot 2^2 \cdot 2^2$

That is, |w| increases by a fixed amount for each correction $= 2 = \frac{2}{\sqrt{2}}$

– We can also prove that $|w^{(k+1)}| \leq C \cdot \sqrt{k}$

That is, |w| increases only sublinearly with k

– This implies that $k \leq C^2 / \epsilon^2 \dots$ a fixed #iterations suffice

Perceptrons 6/6

 H_1^2

Problems

- If the training data is linearly separable, the perceptron algorithm finds a separating hyperplane
- However, there are many options for that hyperplane, some more reasonable than others H_3 ²

+

Hz 2

Hy 2

The perceptron algorithm finds any of these

Perceptron Refinements 1/4

UNI FREIBURG

- Refinements we already discussed
 - Change the (pre-determined) number of iterations
 - Terminate when change in precision (on training set) drops below a certain threshold
 - Remove frequent words
 - Use tf.idf instead of tf to represent documents
 - Use different / additional features, e.g. word bigrams

Averaging

 Take the average of all w from all iterations ... including all the iterations where w did not change That is, if you have 10 iterations and a training set of size 100, you take the average of 1000 w vectors

 Intuition 1: the final changes to w are due to relative few documents (which are still misclassified)

Averaging de-emphasizes the w vectors from the end

 Intuition 2: good values of w are not changed for many iterations (where they classify elements correctly)

Averaging emphasizes those "good" w vectors

– This gives the following refined update step:

Class of x is +1 : $w \leftarrow w + \alpha \cdot a \cdot x$

Class of x is -1: $w \leftarrow w - \alpha \cdot a' \cdot x$

where $a = 1 - S(w \bullet x)$ and $a' = S(w \bullet x)$ and α is a tuning parameter (the so-called learning rate)

Batching

Given a w, consider a whole batch B of training elements
 The size of the batch is a parameter to play around with

- For each x_i ∈ B compute update term with respect to w
 Simple Perceptron: + x if class is +1, -x otherwise
 "Logistic" Perceptron: + α ⋅ a ⋅ x ... or ... α ⋅ a' ⋅ x
- Then add all the update terms to w to obtain a new w
 Batching mainly improves performance (a lot), but it also affects the precision (since it leads to a different w)

References

Wikipedia

- http://en.wikipedia.org/wiki/Linear classifier
- <u>http://en.wikipedia.org/wiki/Perceptron</u>
- <u>https://en.wikipedia.org/wiki/Logistic regression</u>

UNI FREIBURG