

Syntactic N-grams as Machine Learning Features for Natural Language Processing

Marvin Gülzow

Table of Contents

- 1 Basics
 - Syntactic N-Grams
 - Support Vector Machines
 - Naive Bayes
 - Tree Learners (J48)

- 2 Approach

- 3 Results

- 4 Assessment

- 5 Own implementation

- 6 Sources

Syntactic
N-grams as
Machine
Learning
Features for
Natural
Language
Processing

Marvin
Gülzow

Basics

Syntactic
N-Grams
Support Vector
Machines
Naive Bayes
Tree Learners
(J48)

Approach

Results

Assessment

Own imple-
mentation

Sources

- Introduce Syntactic n-grams
- Use them for authorship attribution
- Compare machine learning approaches
 - Support Vector Machines
 - Naive Bayes
 - J48 (decision tree)
- \Rightarrow SVM + SN-Grams work well

Syntactic
N-grams as
Machine
Learning
Features for
Natural
Language
Processing

Marvin
Gülzow

Basics

Syntactic
N-Grams
Support Vector
Machines
Naive Bayes
Tree Learners
(J48)

Approach

Results

Assessment

Own imple-
mentation

Sources

Section 1

Basics

N-Grams

Syntactic
N-grams as
Machine
Learning
Features for
Natural
Language
Processing

Marvin
Gülzow

Basics

Syntactic
N-Grams
Support Vector
Machines
Naive Bayes
Tree Learners
(J48)

Approach

Results

Assessment

Own imple-
mentation

Sources

Definition

n-gram:

$$w = (w_1, \dots, w_n) \in \Sigma$$

- n sequential items from a text
- “item”: characters, words, phonetic units, linguistic features, ...
- “sequential”: Neighborhood relation required
- \Rightarrow Text fragments
- \Rightarrow Probabilistic features

[2]

Syntactic N-Gram:

Syntactic
N-grams as
Machine
Learning
Features for
Natural
Language
Processing

Marvin
Gülzow

Basics

Syntactic
N-Grams
Support Vector
Machines
Naive Bayes
Tree Learners
(J48)

Approach

Results

Assessment

Own imple-
mentation

Sources

Definition

Syntactic N-Gram: “An n-gram obtained based on the order in which the elements appear in syntactic trees”

- Items: SR-Tag (syntactic-relation tag)
- Neighborhood relation: Lie on same path
- Syntactic tree: Parse result according to *formal* grammar
- Issue: *Natural* language processing?
- Stanford NLP suite
- “SN-Grams of SR-tags”

[1], [2]

SN-Grams Example

“Cars with wheels can move”

```
1 -> move/VB (root)
2   -> Cars/NNS (nsubj)
3     -> wheels/NNS (nmod:with)
4       -> with/IN (case)
5 -> can/MD (aux)
```

“Ships with hulls can move”

```
1 -> move/VB (root)
2   -> Ships/NNS (nsubj)
3     -> hulls/NNS (nmod:with)
4       -> with/IN (case)
5 -> can/MD (aux)
```

SN-Grams Example

Syntactic
N-grams as
Machine
Learning
Features for
Natural
Language
Processing

Marvin
Gülzow

Basics

Syntactic
N-Grams
Support Vector
Machines
Naive Bayes
Tree Learners
(J48)

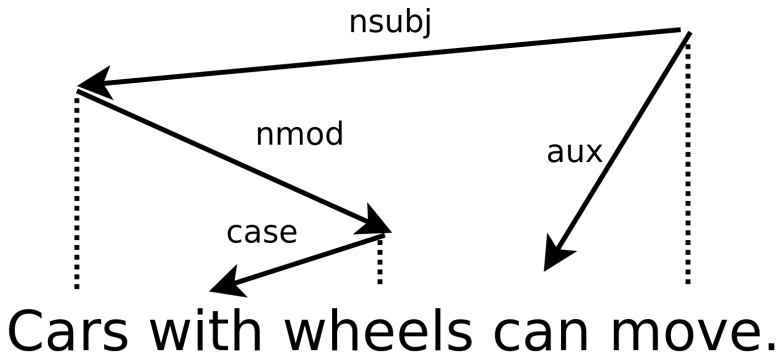
Approach

Results

Assessment

Own imple-
mentation

Sources



SN-Grams Example

Syntactic
N-grams as
Machine
Learning
Features for
Natural
Language
Processing

Marvin
Gülzow

Basics

Syntactic
N-Grams
Support Vector
Machines
Naive Bayes
Tree Learners
(J48)

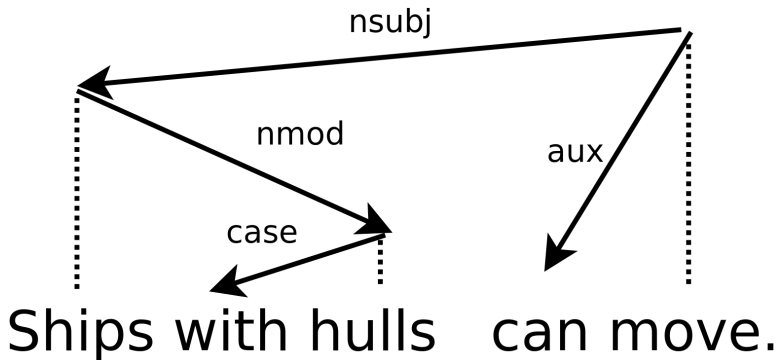
Approach

Results

Assessment

Own imple-
mentation

Sources



Resulting SN-Grams

Syntactic
N-grams as
Machine
Learning
Features for
Natural
Language
Processing

Marvin
Gülzow

Basics

Syntactic
N-Grams

Support Vector
Machines

Naive Bayes

Tree Learners
(J48)

Approach

Results

Assessment

Own imple-
mentation

Sources

- (aux)
- (nsubj, nmod)
- (nmod, case)
- (nsubj, nmod, case)
- \Rightarrow Independent of content.

Syntactic N-Grams

Syntactic
N-grams as
Machine
Learning
Features for
Natural
Language
Processing

Marvin
Gülzow

Basics

Syntactic
N-Grams
Support Vector
Machines
Naive Bayes
Tree Learners
(J48)

Approach

Results

Assessment

Own imple-
mentation

Sources

- Advantages
 - “real” neighbors: No arbitrary influence from content
 - Assumption: Captures author’s writing style
- Disadvantages
 - Preprocessing is expensive (only once though).
 - Parser Quality determines results
 - Good parsers not available for every language

Support Vector Machine (SVM)

Syntactic
N-grams as
Machine
Learning
Features for
Natural
Language
Processing

Marvin
Gülzow

Basics

Syntactic
N-Grams
Support Vector
Machines
Naive Bayes
Tree Learners
(J48)

Approach

Results

Assessment

Own imple-
mentation

Sources

- Deterministic binary classifier
- **linear** separation of classes
- Separator: Hyperplane
- → Gap between classes has maximum width
- Non-linearly separable Data?

[4]

Kernel trick

Syntactic
N-grams as
Machine
Learning
Features for
Natural
Language
Processing

Marvin
Gülzow

Basics

Syntactic
N-Grams

Support Vector
Machines

Naive Bayes

Tree Learners
(J48)

Approach

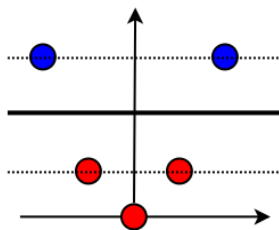
Results

Assessment

Own imple-
mentation

Sources

$$f(x) = (x, x^2)$$



Kernel trick

Syntactic
N-grams as
Machine
Learning
Features for
Natural
Language
Processing

Marvin
Gülzow

Basics

Syntactic
N-Grams
Support Vector
Machines
Naive Bayes
Tree Learners
(J48)

Approach

Results

Assessment

Own imple-
mentation

Sources

- Add dimensions
- Warp data
- \Rightarrow Transformation via kernel-function
- \Rightarrow Restricted to numerical data
- \Rightarrow Multiclass-classification via multiple Binary classification

SVM learning

Syntactic
N-grams as
Machine
Learning
Features for
Natural
Language
Processing

Marvin
Gülzow

Basics

Syntactic
N-Grams
Support Vector
Machines
Naive Bayes
Tree Learners
(J48)

Approach

Results

Assessment

Own imple-
mentation

Sources

- 1 Choose appropriate kernel (human)
- 2 Project data into target vector space
- 3 Find optimum separator
 - Maximize distance of each object to separator
 - \Rightarrow Items defining border are **support vectors**

[6]

Support Vector Machine (SVM)

Syntactic
N-grams as
Machine
Learning
Features for
Natural
Language
Processing

Marvin
Gülzow

Basics

Syntactic
N-Grams
Support Vector
Machines
Naive Bayes
Tree Learners
(J48)

Approach

Results

Assessment

Own imple-
mentation

Sources

- Advantages
 - Non-linear separation
 - “Tunable” to noise
 - More robust against biased data
 - Unique, global solution exists
 - \Rightarrow High accuracy
- Disadvantages
 - Only work on numerical data
 - Learned model not interpretable
 - Training in $O(n^2)$

[5]

Support Vector Machine (SVM)

Syntactic
N-grams as
Machine
Learning
Features for
Natural
Language
Processing

Marvin
Gülzow

Basics

Syntactic
N-Grams
Support Vector
Machines
Naive Bayes
Tree Learners
(J48)

Approach

Results

Assessment

Own imple-
mentation

Sources

- WEKA/LibSVM
- SVMs work on numerical Data
- We have: Nominal data
- \Rightarrow Map semantic relation to numbers

Naive Bayes

Syntactic
N-grams as
Machine
Learning
Features for
Natural
Language
Processing

Marvin
Gülzow

Basics

Syntactic
N-Grams
Support Vector
Machines
Naive Bayes
Tree Learners
(J48)

Approach

Results

Assessment

Own imple-
mentation

Sources

- How many times does an attribute appear in a class?
- \Rightarrow Look at each attribute of item to classify
- \Rightarrow Probabilities determine class
- Each classified object contributes to training set
- Used as a reference for other learners

[4]

Naive Bayes

- Bayes' Theorem

$$P(A|B) = \frac{P(B|A) \cdot P(A)}{P(B)}$$

- Naïve assumption: all attributes are *independent*



$$P(E = (a_1, \dots, a_n) | h) = \prod_{a_i \in E} P(a_i | h)$$



$$P(a_i | h) = \frac{\# \text{data from class } h \text{ with } A_i = a_i}{\# \text{data from class } h}$$

- Object class \Rightarrow most probable

[4]

Naive Bayes

Syntactic
N-grams as
Machine
Learning
Features for
Natural
Language
Processing

Marvin
Gülzow

Basics

Syntactic
N-Grams
Support Vector
Machines
Naive Bayes
Tree Learners
(J48)

Approach

Results

Assessment

Own imple-
mentation

Sources

Advantages

- Easy to implement
- Fast implementation possible
- Learns with each example
- Somewhat accurate
- Standard comparison for other classifiers

Disadvantages

- Attributes are usually **not** independant
- Probabilities may be unavailable

- Decision tree builder
- Entropy based
- \Rightarrow Which attribute yields the highest **information gain**?
- Builds optimum decision tree
- \Rightarrow Human-interpretable model

[6]

J48 - Information Gain

- Given: Labelled dataset
- Find: Attribute which is optimal for discriminating between classes
- Calculate *entropy* of training set T

$$e(T) = - \sum_{i=1}^k p_i \cdot \log_2 p_i$$

- Calculate information gain for attribute A

$$IG(T, A) = e(T) - \sum_{i=1}^m \frac{|T_i|}{|T|} \cdot e(T_i)$$

- \Rightarrow Tree splits data on this attribute
- Repeat
- Other split criteria: Gini-Index, χ^2 , Randomly, ...

[6]

Advantages

- Model can be interpreted for other uses
- Fast classification (precomputed model)
- Can fix missing values (parser errors)

Disadvantages

- Require pruning
- Sensitive to noise
- Greedy approach can get stuck

Section 2

Approach

Dataset

Syntactic
N-grams as
Machine
Learning
Features for
Natural
Language
Processing

Marvin
Gülzow

Basics

Syntactic
N-Grams
Support Vector
Machines
Naive Bayes
Tree Learners
(J48)

Approach

Results

Assessment

Own imple-
mentation

Sources

- English novels
 - Booth Tarkington (13)
 - George Vaizey (13)
 - Louis Tracy (13)
- 24 for Training, 11 for classification
- Total of 6.1 MB

Algorithm

Syntactic
N-grams as
Machine
Learning
Features for
Natural
Language
Processing

Marvin
Gülzow

Basics

Syntactic
N-Grams
Support Vector
Machines
Naive Bayes
Tree Learners
(J48)

Approach

Results

Assessment

Own imple-
mentation

Sources

- 1 Parse Corpus using StanfordNLP
- 2 Extract syntactic relations (SR-tags)
- 3 Construct SN-grams \Rightarrow Profile
- 4 Classify as usual
- 5 Establish baseline using other classifiers

Experiments

Syntactic
N-grams as
Machine
Learning
Features for
Natural
Language
Processing

Marvin
Gülzow

Basics

Syntactic
N-Grams
Support Vector
Machines
Naive Bayes
Tree Learners
(J48)

Approach

Results

Assessment

Own imple-
mentation

Sources

- N-Grams
 - Word based
 - POS (Part Of Speech)
 - Character based
 - SR-Tags
- Vary n-gram size from 2 to 5
- Profile sizes from 400 to 11000
- Use J48 and NB as baseline

Section 3

Results

Results - In brief

Syntactic
N-grams as
Machine
Learning
Features for
Natural
Language
Processing

Marvin
Gülzow

Basics

Syntactic
N-Grams
Support Vector
Machines
Naive Bayes
Tree Learners
(J48)

Approach

Results

Assessment

Own imple-
mentation

Sources

- All classifiers better than 50% accuracy
- SVMs outperform other classifiers
- SR-tags yield better results than other tags
- Bigrams and trigrams better than 4- and 5-grams
- 100% accuracy in some cases

[1], [3]

Results

Syntactic
N-grams as
Machine
Learning
Features for
Natural
Language
Processing

Marvin
Gülzow

Basics

Syntactic
N-Grams
Support Vector
Machines
Naive Bayes
Tree Learners
(J48)

Approach

Results

Assessment

Own imple-
mentation

Sources

```
1 // Show tables from the paper now.  
2 goto PAPER_RESULTS;
```

Section 4

Assessment

Positive

Syntactic
N-grams as
Machine
Learning
Features for
Natural
Language
Processing

Marvin
Gülzow

Basics

Syntactic
N-Grams
Support Vector
Machines
Naive Bayes
Tree Learners
(J48)

Approach

Results

Assessment

Own imple-
mentation

Sources

- SN-Grams provably more accurate than other approaches
- Able to reliably identify author in a small pool of possible authors
- Solid theoretical basis (SVM and parsing)
- Hard to hide author's grammatical habits

Negative

Syntactic
N-grams as
Machine
Learning
Features for
Natural
Language
Processing

Marvin
Gülzow

Basics

Syntactic
N-Grams
Support Vector
Machines
Naive Bayes
Tree Learners
(J48)

Approach

Results

Assessment

Own imple-
mentation

Sources

- Parsing takes “considerable time” on 39 novels
⇒ Mentioned in paper, as expected
- Parser has extreme influence on result
- ⇒ What about “wierd” texts?
 - Non-natives with the speaking of bad grammatics
 - Fantasy/Scifi “bogus” words
- SVM models not interpretable

Paper quality

Syntactic
N-grams as
Machine
Learning
Features for
Natural
Language
Processing

Marvin
Gülzow

Basics

Syntactic
N-Grams
Support Vector
Machines
Naive Bayes
Tree Learners
(J48)

Approach

Results

Assessment

Own imple-
mentation

Sources

Positive:

- Good explanation of SN-Grams
- Thorough comparison of many cases
- Clear results
- New, practical method found

Negative:

- Hard to reproduce:
 - Examples inconsistent
 - No concrete parameters given (Learners!)
 - Tool versions missing
- Small set of candidate authors (3)

Section 5

Own implementation

Own implementation

Syntactic
N-grams as
Machine
Learning
Features for
Natural
Language
Processing

Marvin
Gülzow

Basics

Syntactic
N-Grams
Support Vector
Machines
Naive Bayes
Tree Learners
(J48)

Approach

Results

Assessment

**Own imple-
mentation**

Sources

Not just yet :(

Section 6

Sources

Sources

Syntactic
N-grams as
Machine
Learning
Features for
Natural
Language
Processing

Marvin
Gülzow

Basics

Syntactic
N-Grams
Support Vector
Machines
Naive Bayes
Tree Learners
(J48)

Approach

Results

Assessment

Own imple-
mentation

Sources



G. Sidorov et. al.: *Syntactic N-grams as Machine learning Features for Natual Language Processing*, CIC Mexico, IPN Mexico, University of the Aegean (Greece)



Stanford Cousera lecture on language modeling
(<https://class.coursera.org/nlp/lecture/17>)



Efstathios Stamatatos, *A Survey of Modern Authorship Attribution Methods*, Dept. of Information and Communication Systems Eng, University of the Aegean



Michael Berthold and Iris Adae: *SVMs and Rule Learning*
Lecture held at the University of Constance, Winter term 2014/15



Laura Auria and Rouslan A. Moro: *Support Vector Machines (SVM) as a Technique for Solvency Analysis*, DIW Berlin, 2008