# Sentiment classification of chess annotations

Master-Thesis von Florian Beck
1. Gutachten: Prof. Dr. Johannes Fürnkranz
2. Gutachten:

TECHNISCHE
UNIVERSITÄT
DARMSTADT

Fachbereich Informatik
Knowledge Engineering Group

Sentiment classification of chess annotations

Vorgelegte Master-Thesis von Florian Beck

1. Gutachten: Prof. Dr. Johannes Fürnkranz
2. Gutachten:

Tag der Einreichung:

# Erklärung zur Master-Thesis
# gemäß § 22 Abs. 7 und § 23 Abs. 7 APB TU Darmstadt

Hiermit versichere ich, Florian Beck, die vorliegende Master-Thesis gemäß § 22 Abs. 7 APB der TU Darmstadt ohne Hilfe Dritter und nur mit den angegebenen Quellen und Hilfsmitteln angefertigt zu haben. Alle Stellen, die Quellen entnommen wurden, sind als solche kenntlich gemacht worden. Diese Arbeit hat in gleicher oder ähnlicher Form noch keiner Prüfungsbehörde vorgelegen.

Mir ist bekannt, dass im Falle eines Plagiats (§38 Abs.2 APB) ein Täuschungsversuch vorliegt, der dazu führt, dass die Arbeit mit 5,0 bewertet und damit ein Prüfungsversuch verbraucht wird. Abschlussarbeiten dürfen nur einmal wiederholt werden.

Bei der abgegebenen Thesis stimmen die schriftliche und die zur Archivierung eingereichte elektronische Fassung gemäß § 23 Abs. 7 APB überein.


Darmstadt, den March 18, 2019

_____

(Florian Beck)

**Abstract**

Dies ist der Abstract der Arbeit. Er gibt wertungsfrei, kurz und prägnant den Inhalt der wissenschaftlichen Arbeit wieder.

# Contents

**List of Figures**

# 1  Introduction

## 1.1  Motivation

## 1.2  Problem description

## 1.3  Goal of the thesis

- is it possible to "convert" a chess annotation comment to the appropiate symbol by using a classifier

## 1.4  Structure of the thesis

## 2 Basics

Sentiment Analysis and classification, Ordinal Classification, Word Embeddings, TF-IDF mindestens 5 Seiten

| Symbol | Meaning | Example |
|--------|---------|---------|
| x | capture | |
| + | check | |
| # | checkmate | |
| 0-0 | castling kingside | |
| 0-0-0 | castling queenside | |
| = | promotion | |

**Table 1:** Basic chess notations

## 3 Concept

- general problem: structured data easier to evaluate than unstructured data, possibility to build statistics (how good is a product rated?, political attitude?, player/game statistics in chess -> average count of good/bad moves in a chess game, average count of good/bad moves by a specific player (-> talent scouting?)) –> EXTRACT KNOWLEDGE OUT OF UNSTRUCTURED INFORMATION - general [in this case]: comments [chess annotations] should be converted to symbols (=classes) [chess symbols] - step 1: define input and output - possible input types (symbol: detection of handwritten letters, word/sentence: detection of language, page/file: detection of author - possible output types (letter, language, author) - "precision": language = language family|language|dialect -> chess: good move vs. brilliant/good/slightly good move - step 2: find a database (or similar source) with sufficient information to extract data from - step 3: define conditions that filtered data sets has to fulfill - language restriction - minimal comment length - supervised/unsupervised learning - in case of

  - step 5: tokenize the text - handling of punctuation - step 6: token preprocessing - remove stopwords - lowercase - stemming
    - how to handle order of classes?
    - how to handle differences in class counts?
    - how to handle too many attributes? -> attribute selection

## 4 Approach

### 4.1 Data set extraction

Nur Weka-Classification oder auch ersten Ansatz der NLTK-Classification mit reinnehmen?

#### 4.1.1 PGN-format

Structure of PGN-file, comments in PGN-file

#### 4.1.2 Chessbase-DB

Transformation to PGN, language detection using polyglot

#### 4.1.3 NAGs

Symbols and corresponding NAGs

#### 4.1.4 Python NLTK

RegExp parsing, tokenization, extraction of comment and class

### 4.2 Prepare data for classification

#### 4.2.1 Feature extraction

simple features: count(word), advanced features: tf-idf, bigrams, trigrams

#### 4.2.2 Generating training instances
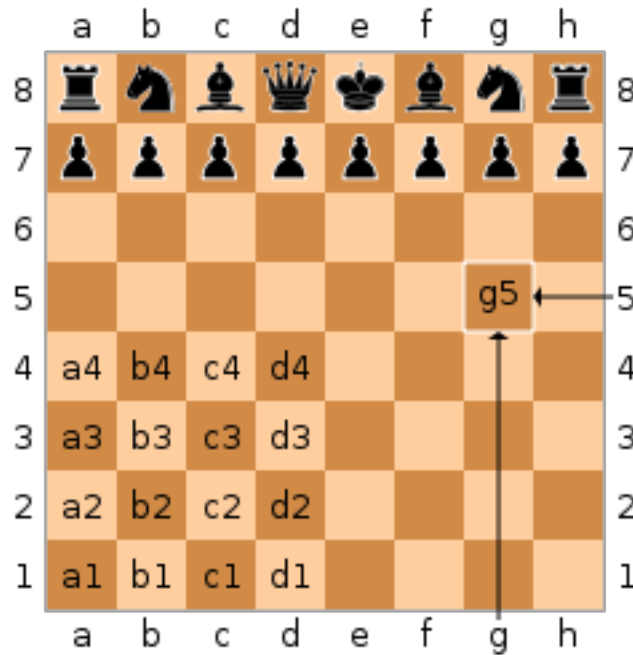
structure of arff-file

**Figure 1:** Square names in algebraic notation

```
[Event "Deutschland "]
[Site "?"]
[Date "1995.??.??"]
[Round "?"]
[White "Lutz, Ch"]
[Black "Kramnik, V."]
[Result "0−1"]
[ECO "B33"]
[PlyCount "70"]
[EventDate "1995.??.??"]
```

1. e4 {B33: Sicilian: Pelikan and Sveshnikov Variations} 1... c5 2. Nf3 Nc6 3.
d4 cxd4 4. Nxd4 Nf6 5. Nc3 e5 6. Ndb5 d6 7. Bg5 a6 8. Na3 b5 9. Nd5 Be7 10.
Bxf6 Bxf6 11. c3 O–O 12. Nc2 Bg5 13. a4 bxa4 14. Rxa4 a5 15. Bc4 Rb8 16. b3 Kh8
17. O–O g6 18. Qe2 Bd7 19. Rfa1 19... Bh6 {last book move} 20. g3 {
Consolidates f4} (20. Nde3 20... Be6 $14) 20... f5 $11 21. exf5 gxf5 22. b4
22... e4 {Black wins space.} 23. bxa5 Ne5 24. Rb4 Rxb4 25. cxb4 f4 26. Nd4 e3
27. fxe3 (27. Nxf4 $2 {doesn't work because of} 27... exf2+ 28. Qxf2 28... Bxf4
$19) 27... f3 {He broke from his leash} (27... fxg3 28. hxg3 Qg5 29. Kh2 Nxc4
30. Nf4 $19) 28. Qa2 f2+ 29. Kg2 Qe8 30. Be2 30... Ng4 {
The pressure on the isolated pawn grows} 31. Bf3 $4 (31. Qd2 Qh5 32. Bxg4 Qxg4
33. Nf4 Bxf4 34. exf4 Qh3+ 35. Kxf2 Qxh2+ 36. Ke1 Qxg3+ 37. Kd1 Qg1+ 38. Ke2
Bg4+ 39. Kd3 Qxa1 40. f5 $19) 31... Nxe3+ $19 32. Nxe3 Qxe3 33. Qxf2 $4 {
sad, but how else could White save the game?.} (33. Rd1 Bg7 34. Qb3 Bxd4 35.
Qxe3 Bxe3 36. Be2 $19) 33... Bh3+ $1 {the final blow} 34. Kg1 {
Black now must not overlook the idea Re1} (34. Kxh3 {A deflection} 34... Qxf2)
34... Qc3 35. Re1 Bd2 (35... Bd2 36. Ne2 36... Qxf3 $19 (36... Bxe1 $6 {
is clearly weaker} 37. Nxc3 Bxf2+ 38. Kxf2 $19) (36... Rxf3 $2 37. Nxc3 Rxf2
38. Kxf2 Bxc3 39. Re7 $18)) 0−1

**Figure 2:** Sample PGN game

## 5 Experiment setup

### 5.1 Selecting classes

distribution how many instances per class, splitting into several problems: 2 classes (!,?), 6 classes (!!,!,!?,?!,?,??), 2 classes (+,-), 7 classes –> introduction of dictionary difference if even or odd number of classes ("neutral class")

### 5.2 Tokenizer tuning

punctuation, special chess notations (#ce etc.)

### 5.3 NLTK Parameters

stopwords, stemming, threshold(hapax), lowercase, bigram, trigram

### 5.4 Classifiers

which classifier to use? –> MCC (x3), OCC, RF, NBM

## 6  Evaluation of results

tables with number of attributes, tables with accuracies, comparison of confusion matrix
each for simple approach, tf-idf, word embedding

# 7 Conclusion

## 7.1 Summary

## 7.2 Outlook

**References**