



MarMot und
Morfessor

Christoph
Teichmann

Problem

Ansätze

MarMot

Morfessor

Übungen

MarMot und Morfessor

Christoph Teichmann

Saarland University

MarMot und
Morfessor

Christoph
Teichmann

Problem

Ansätze

MarMot

Morfessor

Übungen

Morbiditätsrisikoausgleichszahlungsgesetze

- Was ist Part-Of-Speech?
- Was ist Themenbereich?
- Was ist Bedeutung?
- ...

MarMot und
Morfessor

Christoph
Teichmann

Problem

Ansätze

MarMot

Morfessor

Übungen

Unbekannte Worte Kernproblem in CL – mehr Morphologie →
mehr unbekannte Worte

Sentiment: "Der Film war w_{932} "

Unbekannte Worte Kernproblem in CL – mehr Morphologie →
mehr unbekannte Worte

Sentiment: "Der Film war Attahöhlenschlecht"

- Regelbasierte Systeme können Regeln basierend auf Morphologie haben
- Machine Learning wird einfacher wenn es wenige, oft wiederholte Muster gibt

MarMot und
Morfessor

Christoph
Teichmann

Problem

Ansätze

MarMot

Morfessor

Übungen

- Segmentierung:
Attahöhlenschlecht / Atta-höhle-n-schlecht (heute)

MarMot und
Morfessor

Christoph
Teichmann

Problem

Ansätze

MarMot

Morfessor

Übungen

- Segmentierung:
Attahöhlenschlecht / Atta-höhle-n-schlecht (heute)
- Tagging: trug / VT|3pers|sing|prät (heute)



MarMot und
Morfessor

Christoph
Teichmann

Problem

Ansätze

MarMot

Morfessor

Übungen

- Segmentierung:
Attahöhlenschlecht / Atta-höhle-n-schlecht (heute)
- Tagging: trug / VT|3pers|sing|prät (heute)
- Lemmatisierung: trug / tragen (in NLTK)

MarMot und
Morfessor

Christoph
Teichmann

Problem

Ansätze

MarMot

Morfessor

Übungen

- Segmentierung:
Attahöhlenschlecht / Atta-höhle-n-schlecht (heute)
- Tagging: trug / VT|3pers|sing|prät (heute)
- Lemmatisierung: trug / tragen (in NLTK)
- Flektion: tragen+VT|3pers|sing|prät / trug

MarMot und
Morfessor

Christoph
Teichmann

Problem

Ansätze

MarMot

Morfessor

Übungen

- Segmentierung:
Attahöhlenschlecht / Atta-höhle-n-schlecht (heute)
- Tagging: trug / VT|3pers|sing|prät (heute)
- Lemmatisierung: trug / tragen (in NLTK)
- Flektion: tragen+VT|3pers|sing|prät / trug
- ...



MarMot und
Morfessor

Christoph
Teichmann

Problem

Ansätze

MarMot

Morfessor

Übungen

- Morphologischer Tagger von LMU CL



MarMot und
Morfessor

Christoph
Teichmann

Problem

Ansätze

MarMot

Morfessor

Übungen

- Morphologischer Tagger von LMU CL
- Kommandozeilenbasiert



MarMot und
Morfessor

Christoph
Teichmann

Problem

Ansätze

MarMot

Morfessor

Übungen

- Morphologischer Tagger von LMU CL
- Kommandozeilenbasiert
- Wie Part-Of-Speech Tagger aber optimiert für viele Tags/strukturierte Tags



MarMot und
Morfessor

Christoph
Teichmann

Problem

Ansätze

MarMot

Morfessor

Übungen

- Morphologischer Tagger von LMU CL
- Kommandozeilenbasiert
- Wie Part-Of-Speech Tagger aber optimiert für viele Tags/strukturierte Tags
- Machine Learning basiert



MarMot und Morfessor

Christoph
Teichmann

Problem

Ansätze

MarMot

Morfessor

Übungen

- Morphologischer Tagger von LMU CL
- Kommandozeilenbasiert
- Wie Part-Of-Speech Tagger aber optimiert für viele Tags/strukturierte Tags
- Machine Learning basiert
- Vortrainierte Modelle / einfach neue Modelle zu trainieren



MarMot und
Morfessor

Christoph
Teichmann

Problem

Ansätze

MarMot

Morfessor

Übungen

Ausprobieren:

MarMot download

`http://cistern.cis.lmu.de/marmot/bin/CURRENT/ and
model for German`

`http://cistern.cis.lmu.de/marmot/models/CURRENT/`



Ausprobieren:

Mit existierendem Modell taggen

```
java -cp marmot-2015-10-22.jar marmot.morph.cmd.Annotator  
-model-file <wo-model> -test-file  
form-index=<index-wort>,<Daten> -pred-file <Resultat>
```




Ausprobieren:

Neues Modell trainieren

```
java -Xmx2G -cp marmot-2015-10-22.jar  
marmot.morph.cmd.Trainer -train-file  
form-index=<index-wort>,tag-index=<index-pos-tag>,morph-  
index=<index-morph-tag>,<Daten> -tag-morph true  
-model-file <speicherort>
```



MarMot und
Morfessor

Christoph
Teichmann

Problem

Ansätze

MarMot

Morfessor

Übungen

- Tool für morphologische Segmentierung

MarMot und
Morfessor

Christoph
Teichmann

Problem

Ansätze

MarMot

Morfessor

Übungen

- Tool für morphologische Segmentierung
- Unüberwacht → trainiert aus Daten ohne Annotation

MarMot und
Morfessor

Christoph
Teichmann

Problem

Ansätze

MarMot

Morfessor

Übungen

- Tool für morphologische Segmentierung
- Unüberwacht → trainiert aus Daten ohne Annotation
- Nicht notwendigerweise linguistisch korrekt

MarMot und
Morfessor

Christoph
Teichmann

Problem

Ansätze

MarMot

Morfessor

Übungen

- Tool für morphologische Segmentierung
- Unüberwacht → trainiert aus Daten ohne Annotation
- Nicht notwendigerweise linguistisch korrekt
- Versucht widerkehrende Elemente zu finden



MarMot und
Morfessor

Christoph
Teichmann

Problem

Ansätze

MarMot

Morfessor

Übungen

- Tool für morphologische Segmentierung
- Unüberwacht → trainiert aus Daten ohne Annotation
- Nicht notwendigerweise linguistisch korrekt
- Versucht widerkehrende Elemente zu finden
- Minimum Description Length basiert



MarMot und
Morfessor

Christoph
Teichmann

Problem

Ansätze

MarMot

Morfessor

Übungen

- Tool für morphologische Segmentierung
- Unüberwacht → trainiert aus Daten ohne Annotation
- Nicht notwendigerweise linguistisch korrekt
- Versucht widerkehrende Elemente zu finden
- Minimum Description Length basiert
- Nützlich für downstream Machine Learning



MarMot und
Morfessor

Christoph
Teichmann

Problem

Ansätze

MarMot

Morfessor

Übungen

Ausprobieren:

Clonen mit Git von:

`https://github.com/aalto-speech/morfessor`

MarMot und
Morfessor

Christoph
Teichmann

Problem

Ansätze

MarMot

Morfessor

Übungen

Ausprobieren:

Neues Modell trainieren

```
morfessor-train - -traindata-list -s <speicherort> <data>
```

MarMot und
Morfessor

Christoph
Teichmann

Problem

Ansätze

MarMot

Morfessor

Übungen

Ausprobieren:

Mit existierendem Modell segmentieren

```
morfessor-segment -l <model> <data>
```



MarMot und
Morfessor

Christoph
Teichmann

Problem

Ansätze

MarMot

Morfessor

Übungen

- Annotieren sie einen kurzen Text mit MarMot
- Schreiben sie ein kleines Trainingsfile für MarMot und trainieren sie ein neues Modell
- Schreiben sie ein kleines Trainingsfile für Mofessor (viele Morphemwiederholungen) und trainieren sie ein Modell
- Segmentieren sie ihr File mit Morfessor
- Versuchen sie ihr Trainingsfile mit hilfreichen Beispielen zu erweitern