

Master-Entwicklungsprojekt





"Visualisieren und Auswerten großer Datenmengen"

Anwendungsdomäne: Media Systems

Erzeugung eines Dotplots

Das Eclipse-Plugin "Dotplot" stellt mit einer aus der Genetik stammenden Methode Gemeinsamkeiten einer Menge von Zeichenketten, Wörtern, Wortsequenzen oder Sätzen – allgemein "Tokens" genannt – grafisch dar. Bild 1 zeigt anhand diverser Beispiele, wie ein Quelltext auf zwei Achsen einer Matrix verteilt wird und Matches (gleiche Tokens) als Punkt markiert werden. Spezifische Muster im Dotplot lassen den erfahrenen Betrachter schnell auf die Art und den Grad von Ähnlichkeiten schließen. Bild 2 zeigt einen Dotplot zum Text aus Bild 3. Häufig auftretende Muster sind dort synthetisch erzeugt worden. Bild 4 ist eine Selbstanwendung des Dotplot-Plugins auf seine Sourcen. Die Farbskala unter den Bildern 2 und 4 bildet die Gewichte von Matches auf Farben im Dotplot ab. Man kann so die seltenen Treffer von häufigen besser unterscheiden.

Projektziel

- Open-Source-Produkt f\u00fcr die Dotplot-Erzeugung aus beliebigen Textsorten
- ☐ Automatische Plagiaterkennung
- ☐ Refactoring von Programmen
- ☐ Autorenstile identifizieren: "Ein echter Shakespeare?"

Methode

- ☐ Grafisches Dotplot-Verfahren aus der Genforschung
- ☐ Interaktive Muster-Erkennung☐ Sequenz-Alignment

Softwaretechnik

als Eclipse-Plugin

- ☐ Open-Source-Projekt (GNU GPL)
 ☐ Java-Implementierung
- ☐ Kooperationsplattform SourceForge
- ☐ Vorgehensweise: Extreme Programming (XP)

Features

- ✓ Berechnung im Grid
- ✓ Information-Mural-Algorithmus (verlustarme Interpolation)
- ✓ Dotplot-Perspektive in Eclipse
- ✓ Export in diversen Dateiformaten (JPEG. PNG. PDF)
- ✓ PDF-Konverter und Inputfilter für Java, C++, PHP, ...

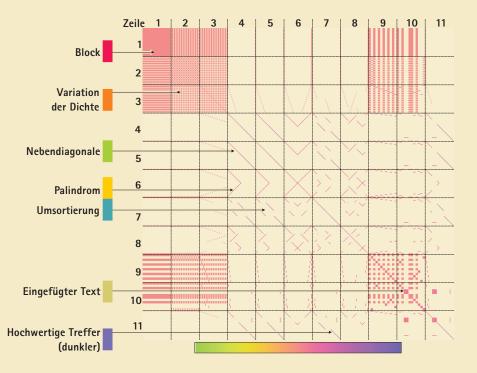


Bild 2: Dotplot (per JAI) der Tokens aus Bild 3. Farben und Zeilenangaben entsprechen den Markierungen im Quelltext. Die synthetisch erzeugten Muster sind mit den entsprechenden Bezeichnungen gekennzeichnet.

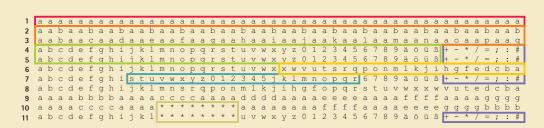
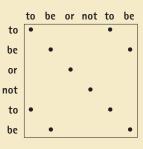
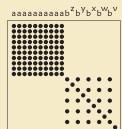
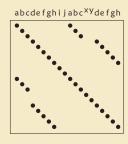
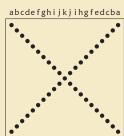


Bild 3: Datenquelle für den in Bild 2 abgebildeten Plot









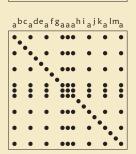


Bild 1: Beispiele zur Entstehung prinzipieller Muster im Dotplot. Von oben nach unten:

- Wörter als Tokens
- Block mit variierender Dichte
- Unterbrochene Nebendiagonalen
- Palindrom
- Dunkles Kreuz

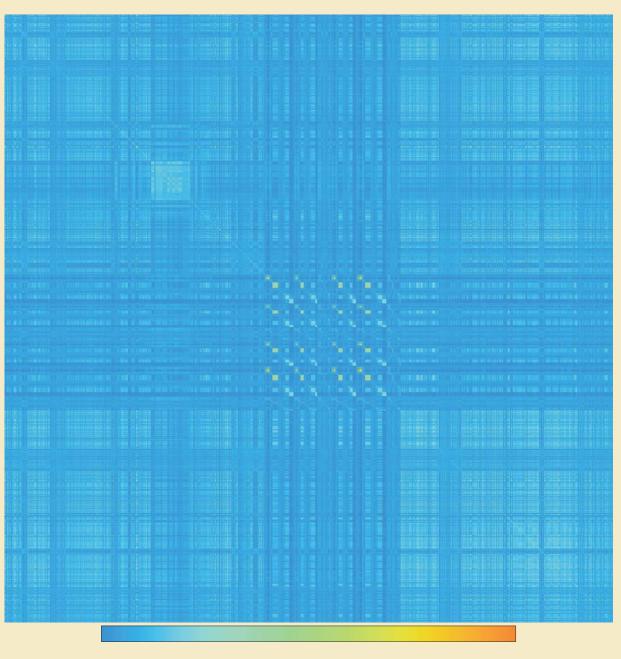


Bild 4: Dotplot (Information Mural) der Java-Sourcen des Eclipse-Plugins: 175 Klassen, 30.000 Zeilen Code

ToDos

- ☐ Automatische Wortstammreduktion und Satzende-Erkennung bei natürlichen Sprachen
- ☐ Implementierung als Rich-Client-Plattform (RCP)
- ☐ Information Mural als Navigationshilfe

Herausforderungen

- ☐ Performancesteigerung
- ☐ Webinterface für Online-Service
- ☐ Visualisierung multimedialer Daten

Projektleiter: Prof. Dr. Klaus Quibeldey-Cirkel

Tutor: cand. ing. Sascha Hemminger

SourceForge-Admin: Dipl.-Inform. (FH) Tobias Gesellchen