

データの分析と知識発見

Introduction to Data Analysis

今回の構成

0101010101010101010101010101010101

テキストマイニングや
そのツールについて知る

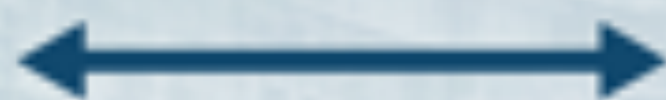
RMeCabを用いて
テキストマイニングを行う

講義全体の振り返りを行う

自然言語処理

01

自然言語



プログラミング言語

形態素解析

文をその構成要素である語に分解する

(例) 本 / を / 読ん / だ

構文解析

文をその構成要素である語と語の関係を分析

(例)

```
graph TD; A[本 / を / 読ん / だ] --- B[本 / を]; A --- C[読ん / だ]; B --- D[本]; B --- E[を]; C --- F[読ん]; C --- G[だ]
```

本 / を / 読ん / だ

形態素解析

010

形態素解析 文を語(形態素)に分解(分かち書き)

(例) 「横浜に行った」

横浜	名詞	固有名詞
に	助詞	格助詞
行っ	動詞	五段活用「行く」の連用形(促音便)
た	助動詞	

絶対値, 印刷教材, 超伝導 / 超電導などのような例も...

青空文庫

010101010101010101010101010101010101

青空文庫

■ 著作権の保護期間が切れたものをボランティアの手で校正・公開されたもの

■ 右の7つの短編を題材にする

芥川龍之介

トロッコ

芥川龍之介

鼻

芥川龍之介

羅生門

有島武郎

一房の葡萄

梶井基次郎

檸檬

小泉八雲

耳なし芳一

新美南吉

ごんぎつね

RMeCabを用いた形態素解析

```
> c1 <- t(a4)
```

```
> c2 <- dist(c1)
```

```
> c3 <- cmdscale(c2,eig=T); c4 <- c3$points
```

```
> c5 <- kmeans(c4,2); c6 <- c5$cluster
```

```
> plot(c4,xlim=c(-25,30),ylim=c(-12,16))
```

```
> text(c4,row.names(c1),col=c6,pos=3)
```

