# 計量政治学における統計的テキスト分析

持橋大地　　数理・推論研究系 准教授　daichi@ism.ac.jp

本研究は, 佐々木智也氏・江島舟星氏 (東京大学 法学政治学研究科 D3・D1) との共同研究です。



発表論文:

Eshima, Shusei (The University of Tokyo) and Daichi Mochihashi (The Institute of Statistical Mathematics), "Tree-Structured Topic Model: a nonparametric Bayesian approach to model texts in a continuous space", Asian Polmeth V (Poster), Seoul National University, 2018.

Sasaki, Tomoya (The University of Tokyo) and Daichi Mochihashi (Institute of Statistical Mathematics), "Detecting Topic Changes among Texts", Asian Polmeth V (Poster), Seoul National University, 2018.

大学共同利用機関法人 情報・システム研究機構
統計数理研究所

The Institute of Statistical Mathematics