

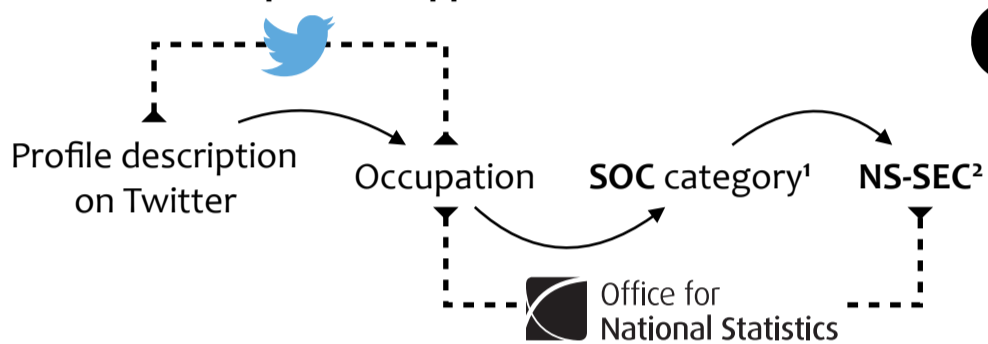
Inferring the Socioeconomic Status of Social Media Users based on Behaviour & Language

Vasileios Lamos, Nikolaos Aletras,
Jens K. Geyti, Bin Zou & Ingemar J. Cox



Summary. We present a method for determining the **socioeconomic status** of a social media (Twitter) user. Initially, we formulate a 3-way classification task, where users are classified as having an **upper, middle** or **lower** socioeconomic status. A nonlinear learning approach using a composite **Gaussian Process** kernel provides a classification accuracy of **75%**. By turning this task into a binary classification – upper vs. medium and lower class – the proposed classifier reaches an accuracy of **82%**.

How is a user profile mapped to a socioeconomic status?



1. **Standard Occupational Classification:** 369 job groupings
2. **National Statistics Socio-Economic Classification:** Map from the job groupings in SOC to a socioeconomic status, i.e. {upper, middle or lower}

Topics (word clusters) are formed by applying **spectral clustering** on daily word frequencies in T2.

Examples of topics with word samples

- Corporate:** #business, clients, development, marketing, offices
- Education:** assignments, coursework, dissertation, essay, library
- Internet Slang:** ahahaha, awwwww, hahaa, hahahaha, hmhhh
- Politics:** #labour, #politics, #tories, conservatives, democracy
- Shopping:** #shopping, asda, bargain, customers, market, retail
- Sports:** #football, #winner, ball, bench, defending, footballer

Formulating a Gaussian Process classifier

Definition:

$$f(\mathbf{x}) \sim \mathcal{GP}(m(\mathbf{x}), k(\mathbf{x}, \mathbf{x}'))$$

$$f : \mathbb{R}^d \rightarrow \mathbb{R} \quad \mathbf{x} \in \mathbb{R}^d$$

Kernel formulation:

$$k(\mathbf{x}, \mathbf{x}') = \left(\sum_{n=1}^C k_{SE}(\mathbf{c}_n, \mathbf{c}'_n) \right) + k_N(\mathbf{x}, \mathbf{x}')$$

where

$$\mathbf{x} = \{\mathbf{c}_1, \dots, \mathbf{c}_C\}, C = 5$$

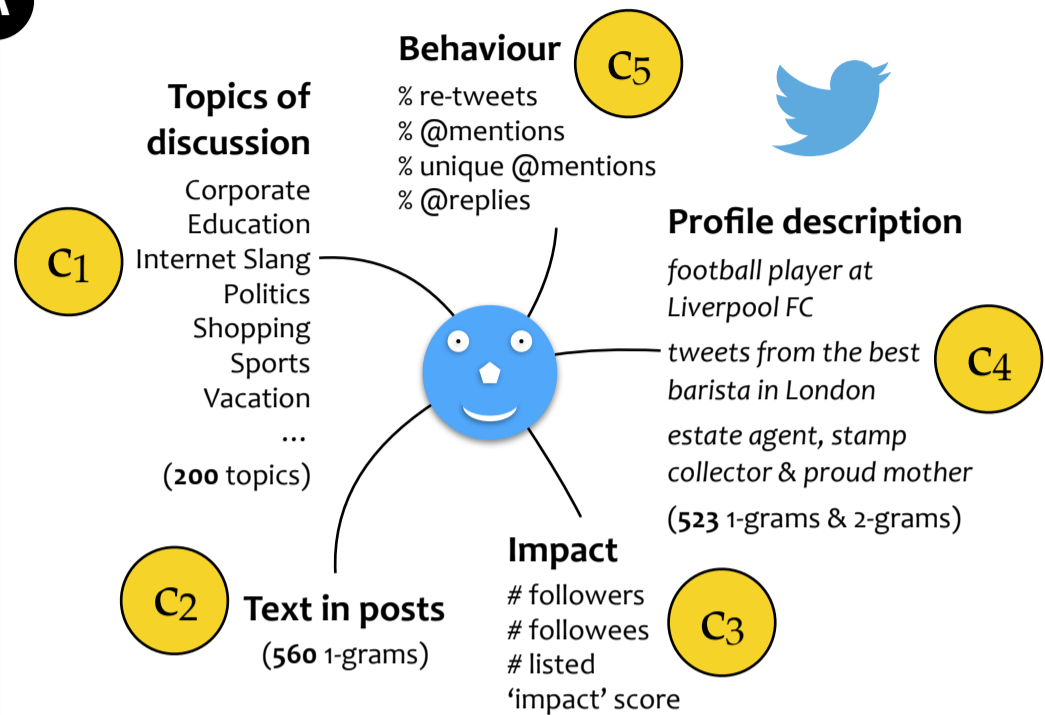
$$k_{SE}(\mathbf{x}, \mathbf{x}') = \theta^2 \exp(-\|\mathbf{x} - \mathbf{x}'\|_2^2 / (2\ell^2))$$

$$k_N(\mathbf{x}, \mathbf{x}') = \theta_N^2 \times \delta(\mathbf{x}, \mathbf{x}')$$



Download the data set

Twitter user attributes (feature categories)



Data sets

- T1:** 1,342 Twitter user profiles, 2 million tweets, from February 1, 2014 to March 21, 2015; profiles are labelled with a socioeconomic status
- T2:** 160 million tweets, sample of UK Twitter, same date range with T1, used to learn a set of 200 latent topics

Confusion matrices (aggregate)

	T1	T2	P		T1	T2	T3	P
O1	584	115	83.5%	O1	606	84	53	81.6%
O2	126	517	80.4%	O2	49	186	45	66.4%
R	82.3%	81.8%	82.0%	O3	55	48	216	67.7%
				R	85.4%	58.5%	68.8%	75.1%

O = output (inferred), T = target, P = precision, R = recall
{1, 2, 3} = {upper, middle, lower} socioeconomic status

Classification performance (10-fold CV)

Classification	Accuracy (%)	Precision (%)	Recall (%)	F1
2-way	82.05 (2.4)	82.2 (2.4)	81.97 (2.6)	.821 (.03)
3-way	75.09 (3.3)	72.04 (4.4)	70.76 (5.7)	.714 (.05)

Conclusions. (a) First approach for inferring the socioeconomic status of a social media user, (b) 75% & 82% accuracy for the 3-way and binary classification tasks respectively, and (c) future work is required to evaluate this framework more rigorously and to analyse underlying qualitative properties in detail.

Funded by the EPSRC IRC project "i-sense"



Selected References

- Lamos et al. Predicting and Characterising User Impact on Twitter. EACL, 2014.
- Preotiuc-Pietro et al. An analysis of the user occupational class through Twitter content. ACL, 2015.
- Preotiuc-Pietro et al. Studying User Income through Language, Behaviour and Affect in Social Media. PLoS ONE, 2015.
- Rasmussen and Williams. Gaussian Processes for Machine Learning. MIT Press, 2006.