

Google Trends can improve surveillance of Type 2 diabetes

Nataliya Tkachenko^{1,*}, Sarunkorn Chotvijit¹, Neha Gupta¹, Emma Bradley³, Charlotte Gilks³, Weisi Guo¹, Henry Crosby¹, Eliot Shore¹, Malkiat Thiarai¹, Rob Procter^{1,2}, and Stephen Jarvis^{1,2}

¹Warwick Institute for the Science of Cities, University of Warwick, Coventry, CV4 7AL, UK

²Department of Computer Science, University of Warwick, Coventry, CV4 7AL, UK

³Experian, The Sir John Peace Building, Experian Way, NG2 Business Park, Nottingham, NG80 1ZZ, UK

*Corresponding Author: Nataliya.Tkachenko@warwick.ac.uk

SUPPLEMENTARY MATERIAL:

1. Database components constituting Experian Mosaic Public Sector variables (Fig 1)
2. Pearson correlation matrix between diabetes risk variables, selected from EMPS (Fig 2 and Fig 3)
3. Resulting metrics from Multilinear Regression (MLR) and Stepwise Multilinear Regression (SMR) (Table 1)
4. List of risk- and disease related keywords, used to extract weekly search volumes from the Google Trends

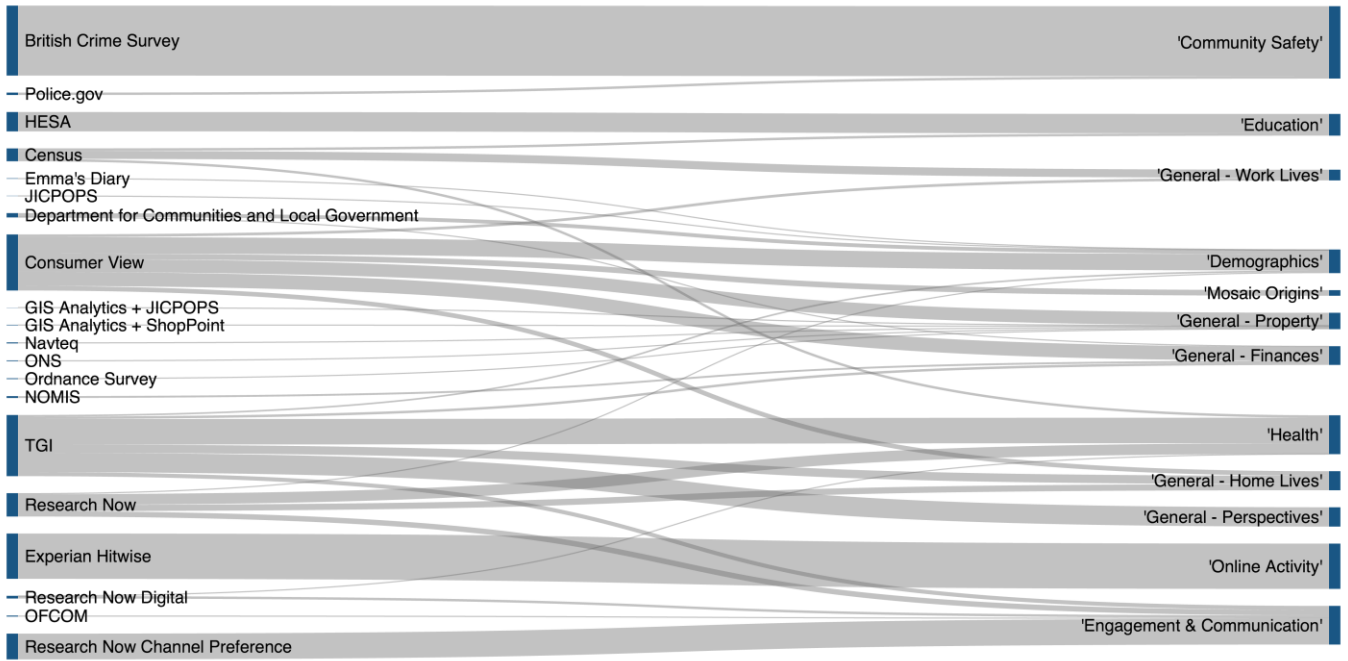


Fig 1. Database components constituting Experian Mosaic Public Sector variables

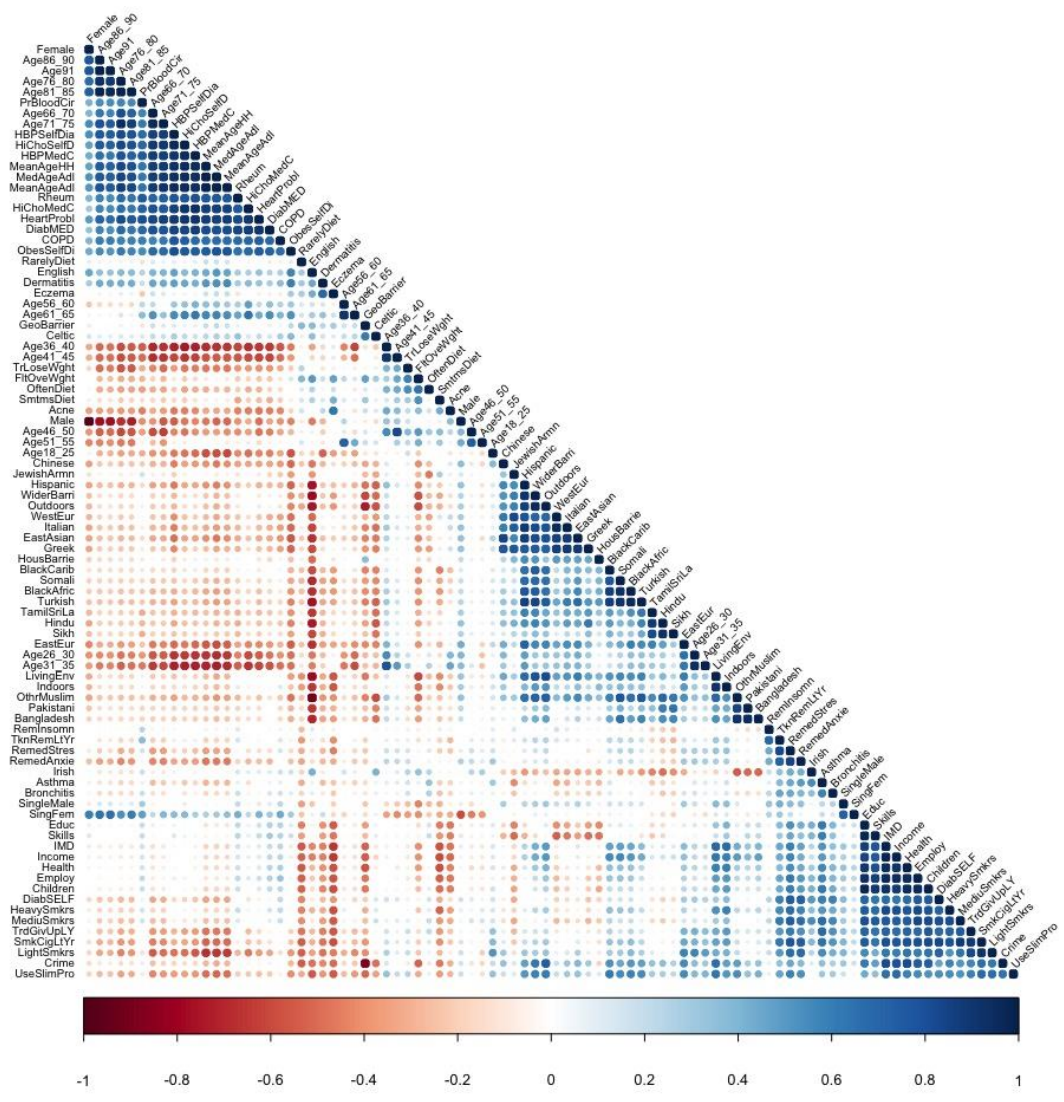


Fig 2. Pearson correlation matrix between diabetes risk variables, selected from EMPS. This illustration was produced in R corrgram v1.9.0

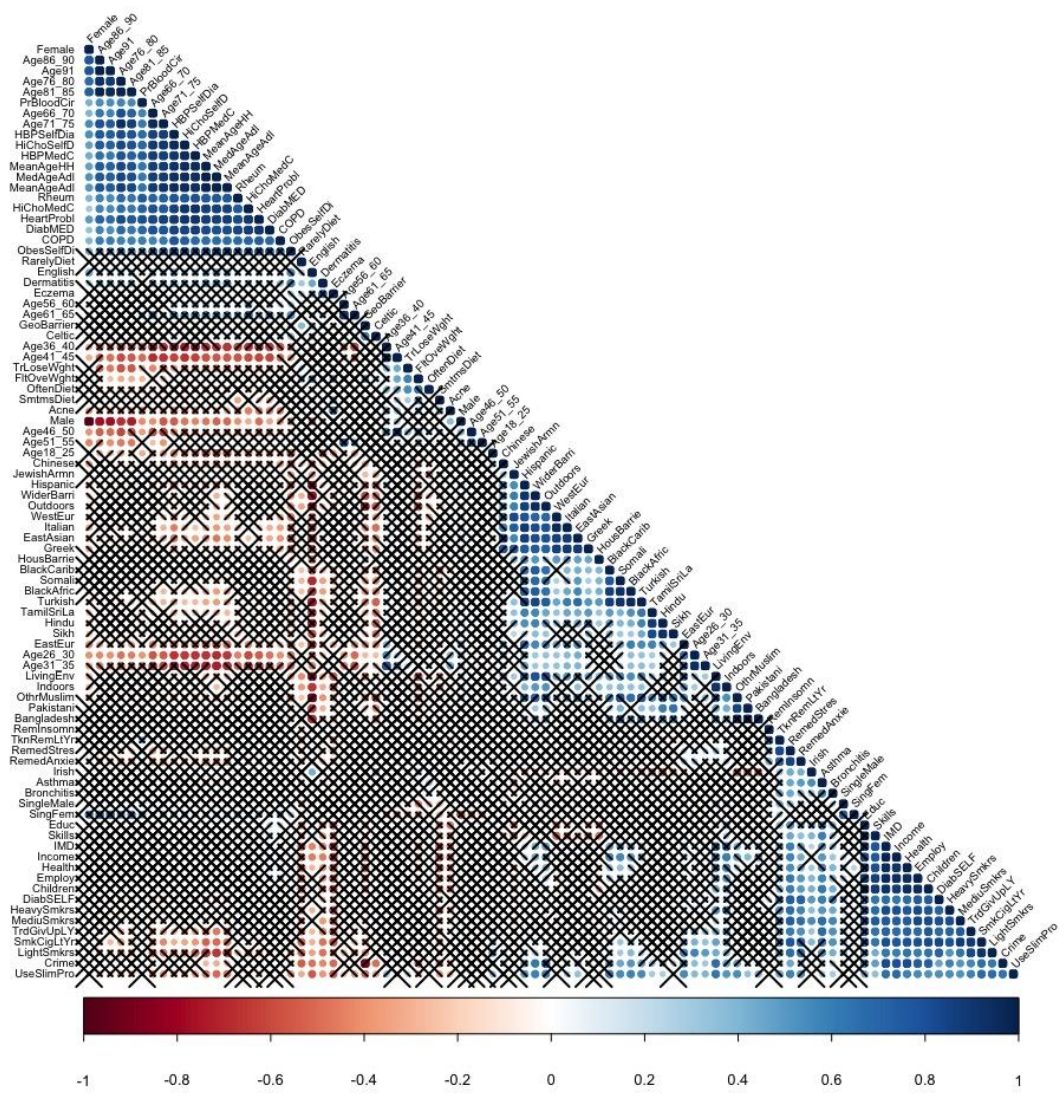


Fig 3. Pearson correlation matrix between diabetes risk variables, selected from EMPS: correlation significance at $p < .05$. This illustration was produced in R corrgram v1.9.0

	Without ISB			With ISB		
	F*	RSS	AIC**	F*	RSS	AIC**
Multilinear	9.83	4.05		26.25	1.51	
Backward AIC-Stepwise	31.65	2.56	53.15	45.52	1.09	16.88
Forward AIC-Stepwise	29.16	2.96	71.17	41.18	1.15	26.62

Table 1. Resulting metrics from regression scenarios demonstrate that the models with self-assessed diabetes variable (With ISB) perform better than traditional diabetes model variables (Without ISB). In this case, the Backward AIC-Stepwise model outperformed both Forward AIC-Stepwise and classic multilinear models, with the metrics being consistent across both scenarios.

(*p<.0001; **Akaike Information Criterion)

	Without ISB		With ISB	
	Group	Keywords	Group	Keywords
Disease	<i>'diabetes', 'Type 2 diabetes', 'diabetes mellitus'</i>			
Risk factors	Use of Corticosteroids	<i>'anxiety', 'eczema', 'acne'</i>	Use of Corticosteroids	<i>'anxiety', 'acne', 'insomnia'</i>
	Smoking	<i>'smoking', 'how to give up smoking', 'give up smoking'</i>	NA	
	Deprivation	<i>'deprivation', 'housing', 'cheap houses'</i>	Deprivation	<i>'deprivation', 'housing', 'cheap houses'</i>
	Ethnicity	<i>'pakistani', 'irish', 'celtic', 'black caribbean'</i>	Ethnicity	<i>'irish', 'somalil', 'sikh', 'celtic', 'black caribbean', 'eastern european'</i>
	Treated Hypertension	<i>'rheumatism'</i>	Treated Hypertension	<i>'high cholesterol'</i>
	BMI	<i>'obesity'</i>	BMI	<i>'slimming products', 'how to lose weight', 'diet'</i>

Table 2. List of risk- and disease related keywords, used to extract weekly search volumes from the Google Trends. Data was extracted on the 7th of April 2017, for each risk group only the data from the search query with the highest propensity have been used in the analysis.