

Entropy*

Fabian Gunzinger Neil Stewart
Warwick Business School Warwick Business School

January 6, 2022

Contents

1	Introduction	1
2	Data	2
2.1	Preprocessing	2
2.2	Sample selection	2
2.3	Dataset description	3
2.4	Dependent variable	4
2.5	Independent variable	6
3	Methods	7
3.1	Model specification	7

1 Introduction

Nomenclature:

- user : Individual - ‘tag’ : Spending categories

Literature:

Muggleton et al. (2020) find that consumption entropy over categories correlates with financial distress.

Davenport et al. (2020) study the impact of COVID-19 on the spending and savings behaviour of MDB users.

Baker and Kueng (2021) summarises literature that uses mass financial transaction data to study household financial behaviour.

*This research was supported by Economic and Social Research Council grant number ES/V004867/1. WBS ethics code: E-414-01-20.

Becker (2017) finds that access to a fintech money management app increases first-time savings and savings account balances among 65,000 customers of a large European bank but that update is negatively correlated with financial sophistication.

Colby and Chapman (2013) has useful literature review on short-term savings and suggests that subgoals can increase willingness to forego short-amounts in the present because they move the reference point in a prospect-theory framework.

Paper:

Independent variable: entropy over categories and others

Outcome variables: first-time saving, average monthly savings

2 Data

2.1 Preprocessing

The MDB data is noisy. We perform a number of preprocessing steps to deal with that.

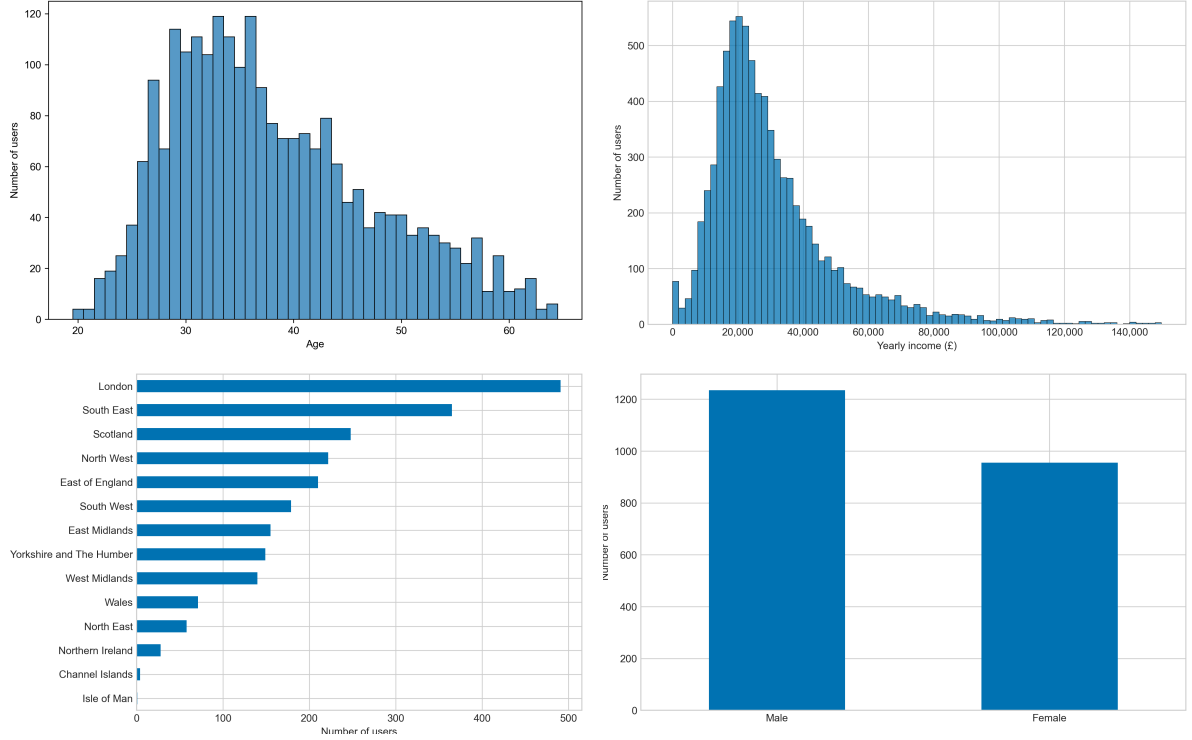
- Duplicates handling.
- We trim all variables at the 1-percent level on the upper end of the distribution for variables that take non-negative values only and on both ends of the distribution for all other variables. We trim (replace outliers with missing values) rather than winsorise (replace outliers with the cutoff percentile value) because we believe that outliers result from errors in the data rather than represent genuine information.

2.2 Sample selection

Table 1: Sample selection

	Users	Accounts	Transactions	Value (£M)
Raw sample	23,785	122,154	58,189,203	10,946.7
At least 6 months of data	20,949	114,941	57,629,233	10,857.3
No missing months	18,212	97,573	50,350,143	9,444.9
Account balances available	14,469	84,889	43,467,849	8,488.5
At least 10 debits totalling £200 per month	11,702	67,688	37,011,881	7,148.0
At least one current account	11,559	67,088	36,688,783	7,102.6
Income in 2/3 of all observed months	9,601	57,832	32,080,562	6,257.4
Yearly income between £5k and £200k	6,612	36,123	21,016,001	3,590.5
No more than 10 accounts in any year	6,080	26,613	17,874,988	2,795.9
Debits of less than £100k each month	5,677	24,140	16,074,957	1,953.5
Final sample	5,677	24,140	16,074,957	1,953.5

Figure 1: Demographic characteristics of Money Dashboard users



2.3 Dataset description

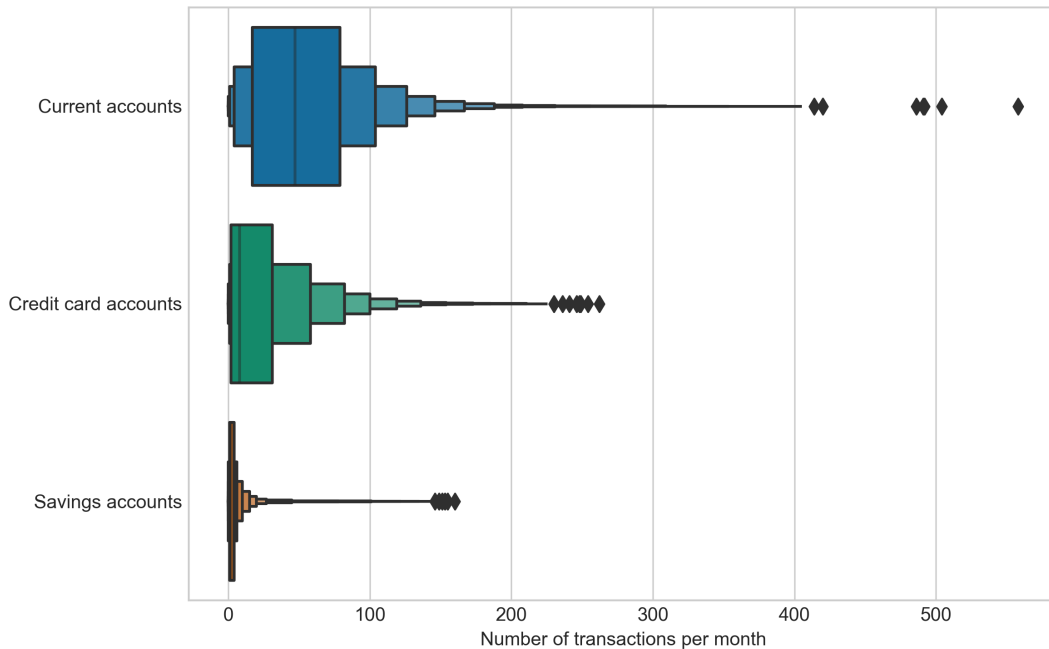
Data is provided by Money Dashboard (MDB), a UK-based financial management app that allows its users to add accounts from all their banks to obtain an integrated view of their finances. Our dataset contains information on more than 500 million transactions made between 2012 and June 2020 by more than 250,000 users. For each transaction, we can see the amount, date, and description of the transaction, as well as transaction *tags*, classifications added by MDB that indicate the type of the transaction (e.g. ‘groceries’, ‘insurance’). We also have basic information on each user (e.g. year of birth, postcode sector) as well information about each bank account (e.g. type of account, date added).

The main advantages of the data for the study of consumer financial behaviour are its high (transaction-level) frequency, that it is automatically collected and updated and thus less prone to errors and unaffected by biases that bedevil survey measures, and that it offers a view of consumers’ entire financial life across all their accounts, rather than just a view of their accounts held at a single bank (provided they added all their accounts to MDB).

The main limitation is the non-representativeness of the sample relative to the population as a whole. Financial management apps are known to be used disproportionately by men, younger people, and people of higher socioeconomic status (Carlin et al. 2019, MAS 2014), a fact that is reflected in our data. The four panels in Figure 1 provide an overview

of demographic characteristics of our sample. It makes clear that Money Dashboard users are not a representative sample of the UK population: they are predominantly males in their thirties who live in London or the South East and are relatively well off (the income distribution is shifted to the right relative to the UK as a whole).¹ Also, as pointed out in Gelman et al. (2014), a willingness to share financial information with a third party might not only select on demographic characteristics, but also for an increased need for financial management, or, one could argue, for a higher degree of financial sophistication. However, while non-representativeness could partially be addressed by re-weighting the sample, as was done in Bourquin et al. (2020), it is not of much consequence for our purpose here, since our ability to infer behaviour traits from transaction data is not dependent on having a representative sample of people.

Figure 2: Monthly transactions by account type



Notes: The two innermost boxes in the [letter-value plots](#) are identical to those in a boxplot, with the center line corresponding to the median and the left and right edges to the first and third quartiles, respectively – or half of the remaining data on either side of the median. Additional boxes on either side extend that principle by corresponding to half of the remaining data on that side. For instance, the second box to the right of the median in the current accounts plot indicates that half of all account-month observations to the right of the third quartile have fewer than about 105 transactions. Boxes of the same height correspond to the same level, individually drawn observations are outliers.

2.4 Dependent variable

- Add notes on excluding non-standing orders from ipynb.

¹To calculate incomes, we broadly follow Hacıoglu et al. (2020) in defining total income as the sum of earnings, pension income, benefits, and other income.

Table 2: Summary statistics

	count	mean	std	min	max	25%	75%
obs	163915	98.0689	51.2574	11	656	63	83
balance_ca	159116	1130.2	4834.21	-12257.1	34628.6	-884.008	8344.8
balance_sa	64935	2625.16	5496.73	-1644.62	42168.6	2.21001	8344.8
sa_inflows	60583	780.844	1507.55	0	13800	60	8344.8
sa_outflows	60583	749.975	1452.1	0	12075.5	0	8344.8
sa_net_inflows	61195	75.7461	3330.15	-69750	120000	-180	8344.8
sa_scaled_inflows	59977	0.335841	0.548266	0	4.13782	0.0319006	0.1
sa_scaled_outflows	59977	0.332819	0.57624	0	4.03946	0	0.1
sa_scaled_net_inflows	59971	0.00713739	0.614697	-4.03481	4.18434	-0.0845304	0.008
entropy_sptac	160635	2.56947	0.215097	1.90212	2.99952	2.43239	2.99952
total_monthly_spend	160635	1811.98	1416.9	-1145.77	9431.91	869.835	8344.8
tag_spend_household	160635	695.566	688.575	-184.88	3896.09	194.725	8344.8
tag_spend_hobbies	161944	18.2805	36.1727	0	289.99	0	8344.8
tag_spend_retail	160635	97.0552	164.035	-990.36	946.76	7.99	8344.8
tag_spend_services	160635	336.221	317.696	-219.87	2051.06	121.92	8344.8
tag_spend_other_spend	160639	171.52	289.513	-1599.22	1600	29	8344.8
tag_spend_finance	160635	221.244	354.74	-70.5	2540.71	10.4	8344.8
tag_spend_travel	161412	94.7894	184.964	0	1558.76	0	8344.8
tag_spend_communication	161558	56.7774	49.098	0	283.04	20.36	8344.8
tag_spend_motor	162073	63.666	79.9709	0	428.82	0	8344.8

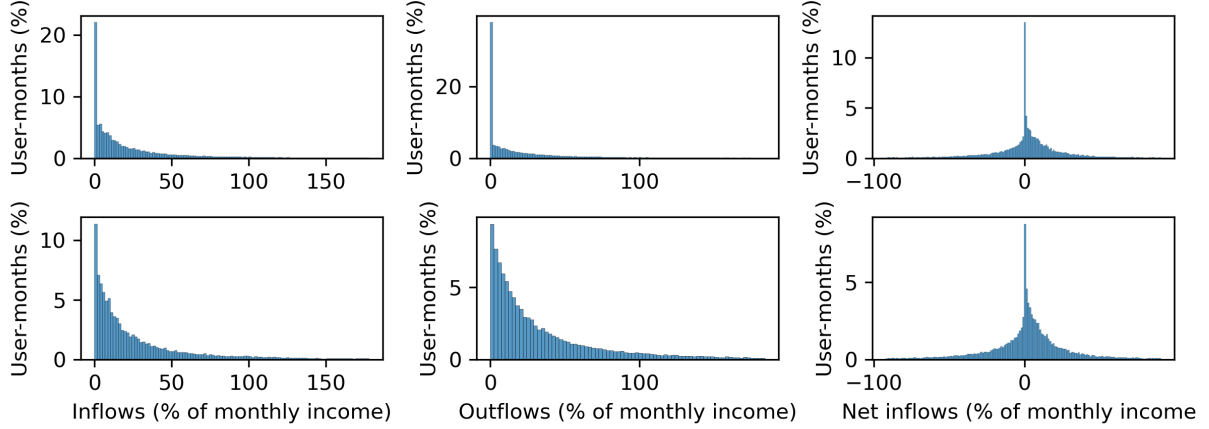
Types of savings:

- Current account balances
- Savings account balances
- ISA

Types of balances, from Becker (2017), who treats balance at end of each month as observations:

- Current account balance
- Debit balance (savings and current account balance)
- Pure savings (savings account balance only)
- Credit balance (loans and negative current account)
- Pure credit (loans only)
- Wealth held (debit - credit balance)

Figure 3: Monthly flows in and out of savings accounts



Notes: Flows are calculated for each user-month as the total inflows, outflows, and flows (calculated as inflows - outflows) into all of a users's savings accounts. Zero net flows represent months where inflows are either perfectly balanced by outflows or where there were no flows at all.

2.5 Independent variable

Spending entropy:

- We calculate spending entropy using the Shannon entropy H (Shannon 1948), defined as

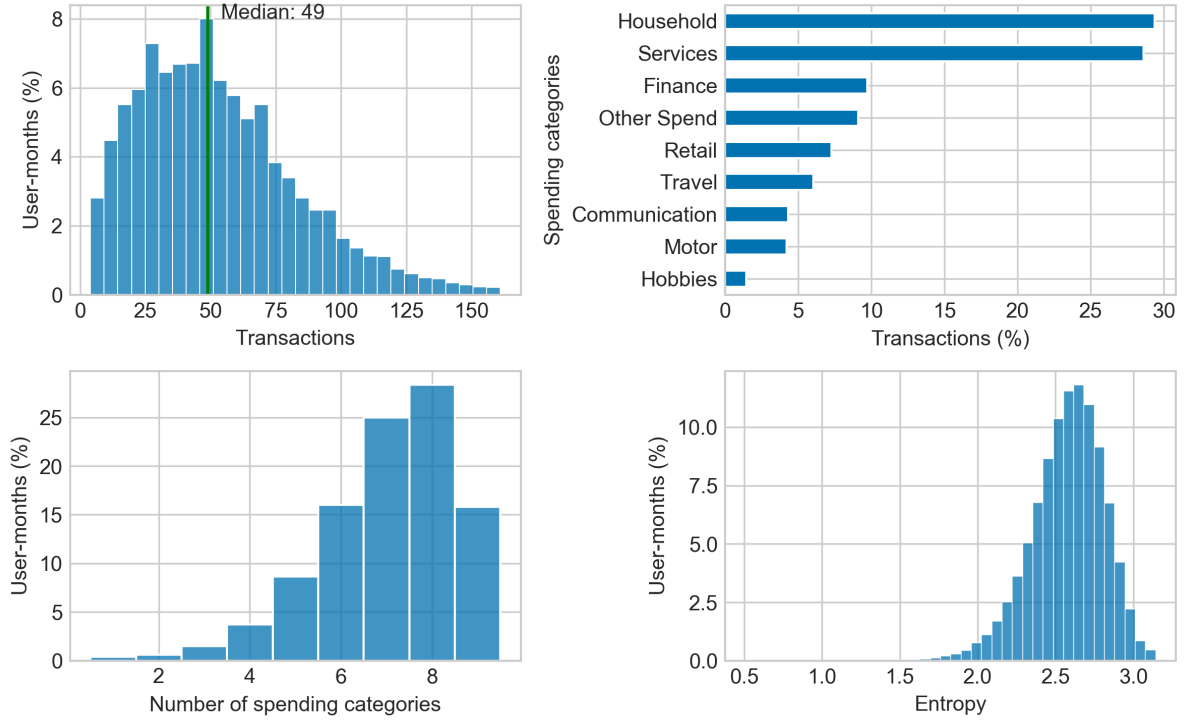
$$H = - \sum p_i \log(p_i), \quad (1)$$

where p_i is the probability that an individual makes a purchase in spending category i , and \log is the base 2 logarithm. The measure can broadly be interpreted as the degree to which an individual's spending pattern is predictable, with a higher score indicating less predictability.

- To calculate individual entropy scores, we group spending into 9 spending categories (SC), based on the classification used by Lloyds Banking Group as discussed in Muggleton et al. (2020).
- Also following that paper, when calculating p_i we use additive smoothing and add one to the numerator and N_{SC} to the denominator to avoid taking logs of zero counts in cases where an individual makes no purchases in a given spending category. p_i is thus calculated as

$$p_i = \frac{\text{Count of purchases in } SC_i + 1}{\text{Count of all purchases} + 9} \quad (2)$$

Figure 4: Transactions distributions



Notes: From top-left to bottom-right: distribution of spending transactions per user-month, breakdown of spending transactions into spending categories; breakdown of number of spending categories spent on in user-month; distribution of user-month entropy scores.

3 Methods

3.1 Model specification

$$s_{i,t} = \alpha_i + \lambda_t + \beta H_{i,t} + X'_{i,t} \delta + \epsilon_{i,t} \quad (3)$$

$s_{i,t}$ is individual i 's savings rate in month t , calculated as the total inflow of funds in month t into all savings accounts held by i , divided by i 's estimated monthly income.

The vector of control variables, $X_{i,t}$, contains the monthly spend for each spending category, total monthly spend across all categories, and annual income.

References

- Baker, Scott R and Lorenz Kueng (2021). “Household Financial Transaction Data”. Tech. rep. National Bureau of Economic Research.
- Becker, G (2017). “Does fintech affect household saving behavior? findings from a natural field experiment”. Tech. rep. mimeo.
- Bourquin, Pascale, Isaac Delestre, Robert Joyce, Imran Rasul, and Tom Walters (2020). “The effects of coronavirus on household finances and financial distress”. In: *IFS Briefing Note BN298*.
- Carlin, Bruce, Arna Olafsson, and Michaela Pagel (2019). “Generational Differences in Managing Personal Finances”. In: *AEA Papers and Proceedings*. Vol. 109, pp. 54–59.
- Colby, Helen and Gretchen B Chapman (2013). “Savings, subgoals, and reference points”. In:
- Davenport, Alex, Robert Joyce, Imran Rasul, and Tom Waters (2020). “Spending and saving during the COVID-19 crisis: evidence from bank account data”. In: *Institute for Fiscal Studies, Briefing Note 308*.
- Gelman, Michael, Shachar Kariv, Matthew D Shapiro, Dan Silverman, and Steven Tadelis (2014). “Harnessing naturally occurring data to measure the response of spending to income”. In: *Science* 345.6193, pp. 212–215.
- Hacioglu, Sinem, Diego Känzig, and Paolo Surico (2020). “The Distributional Impact of the Pandemic”. In:
- MAS, Money Advice Service (2014). *Money Lives: the financial behaviour of the UK*. URL: <https://www.moneyadviceservice.org.uk/en/corporate/money-lives> (visited on 04/07/2020).
- Muggleton, Naomi K, Edika G Quispe-Torreblanca, David Leake, John Gathergood, and Neil Stewart (2020). “Evidence from mass-transactional data that chaotic spending behaviour precedes consumer financial distress”. Tech. rep. DOI: [10.31234/osf.io/qabgm](https://doi.org/10.31234/osf.io/qabgm). URL: psyarxiv.com/qabgm.
- Shannon, Claude Elwood (1948). “A mathematical theory of communication”. In: *The Bell system technical journal* 27.3, pp. 379–423.