

# Entropy\*

Fabian Gunzinger                      Neil Stewart  
Warwick Business School          Warwick Business School

January 3, 2022

## Contents

<b>1</b>	<b>Introduction</b>	<b>1</b>
<b>2</b>	<b>Background</b>	<b>2</b>
<b>3</b>	<b>Data</b>	<b>2</b>
3.1	Preprocessing . . . . .	2
3.2	Sample selection . . . . .	3
3.3	Dataset description . . . . .	3
3.4	Dependent variable . . . . .	4
3.5	Independent variable . . . . .	5
<b>4</b>	<b>Methods</b>	<b>6</b>
4.1	Model specification . . . . .	6

## 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.

---

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

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

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 Background

Types of savings:

- Current account balances
- Savigns account balances
- ISA

Table 1: Sample selection

0	count	57158
1	mean	1060.95
2	std	3742.03
3	min	0
4	25%	50
5	50%	250
6	75%	800
7	max	237150

## 3 Data

### 3.1 Preprocessing

Duplicate transactions

Table 2: Sample selection

	Users	Accounts	Transactions	Value (£M)
Raw sample	24,159	123,597	59,635,566	11,208.3
At least 6 months of data	21,504	116,971	59,076,492	11,115.3
No missing months	18,550	98,518	51,249,027	9,605.7
Account balances available	14,715	50,837	33,581,635	7,165.4
At least 5 debits totalling £200 per month	11,461	39,175	28,811,005	5,819.7
At least one current account	11,447	39,147	28,792,062	5,818.0
Income in 2/3 of all observed months	9,315	33,100	24,669,019	5,008.5
Yearly income between £5k and £200k	6,455	21,383	16,807,675	2,990.9
No more than 10 accounts in any year	6,332	19,315	16,218,234	2,779.1
Debits of less than £100k each month	5,926	17,482	14,839,981	1,950.8
Final sample	5,926	17,482	14,839,981	1,950.8

## 3.2 Sample selection

## 3.3 Dataset description

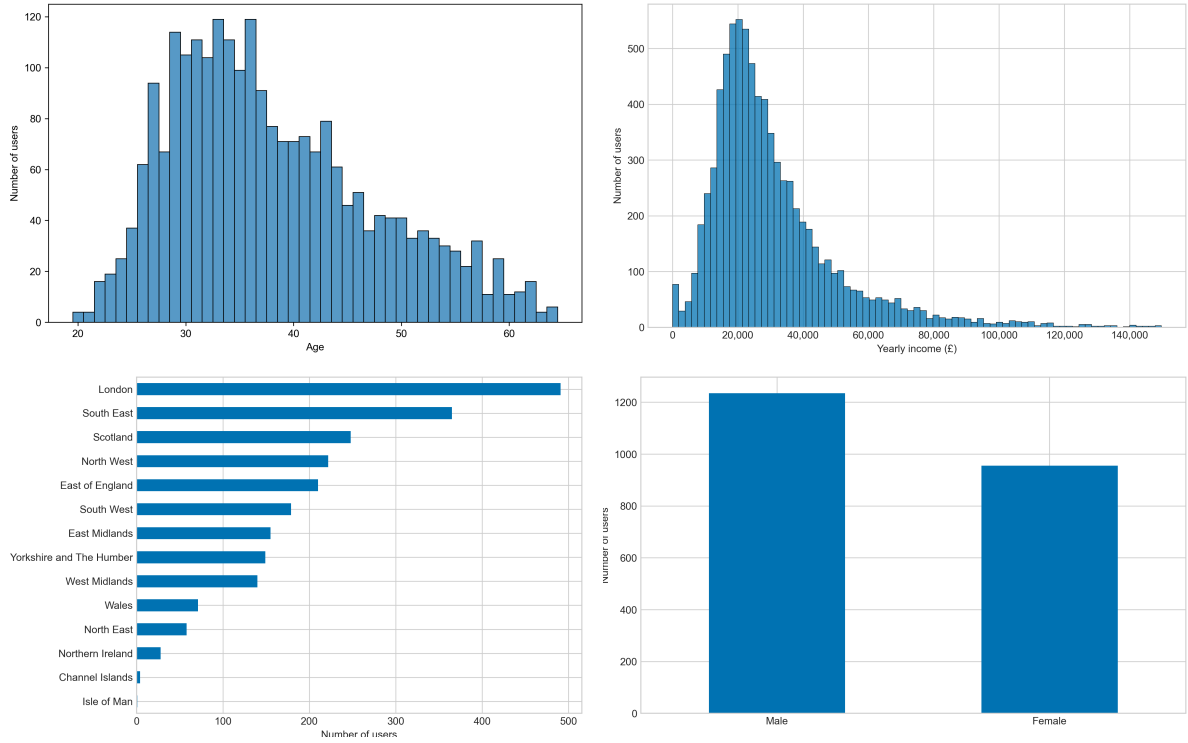
Data is provided by Money Dashboard (MDB), a UK-based financial management app that allows its users to add accounts from all their banks to obtain an integrated view of their finances. Our dataset contains information on more than 500 million transactions made between 2012 and June 2020 by more than 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).<sup>1</sup> Also, as pointed out in

<sup>1</sup>To calculate incomes, we broadly follow Hacıoglu et al. (2020) in defining total income as the sum of

Figure 1: Demographic characteristics of Money Dashboard users



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.

### 3.4 Dependent variable

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

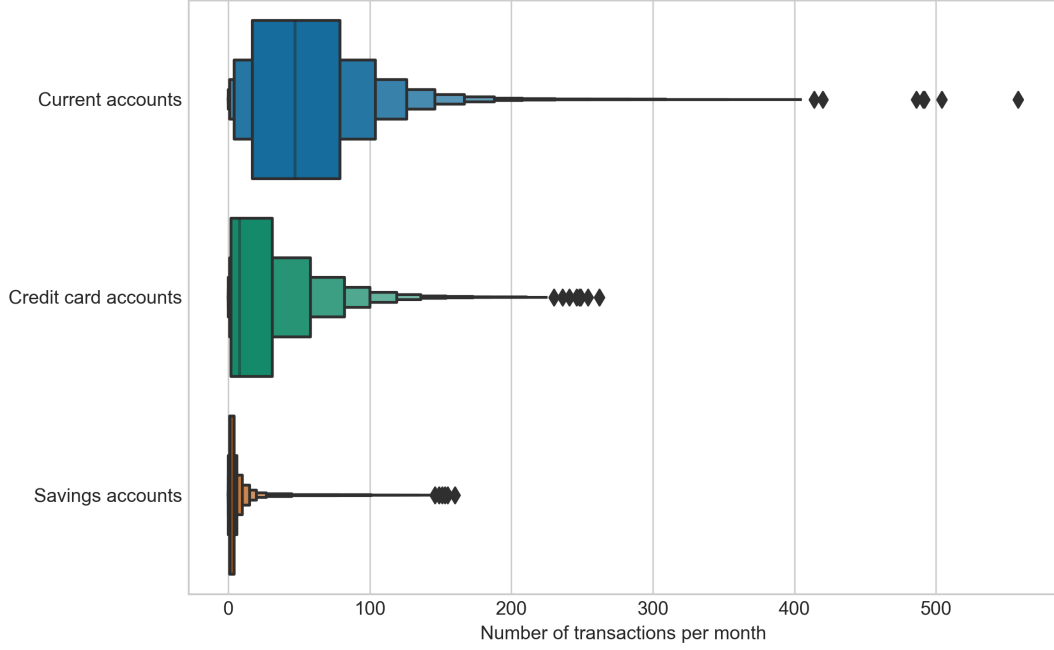
Types of balances, from Becker (2017), who treats balance at end of each month as observations:

- Current account balance
- Debit balance (savings and current account balance)
- Pure savings (savings account balance only)

---

earnings, pension income, benefits, and other income.

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.

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

### 3.5 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

Table 3: Summary statistics

	count	mean	std	min	max	25%	50%
obs	157120	81.3846	35.3007	25	174	53	77
balance_ca	156811	1125.99	2598.34	-3216.27	12724.9	-230.618	450.785
balance_sa	56253	2291.72	3342.4	-47.5005	17509.6	176.824	836.732
sa_inflows	54406	516.375	718.691	0	4000	50	225.56
sa_outflows	54376	486.503	747.158	0	3974.29	0	152.35
sa_net_inflows	51551	37.4411	563.8	-2005.5	2000	-120	25
sa_scaled_inflows	54376	24.4935	33.1865	0	177.896	2.40662	11.4112
sa_scaled_outflows	54376	23.2617	35.1391	0	183.834	0	7.74404
sa_scaled_net_inflows	51514	1.76159	25.985	-91.7285	90.1245	-6.03842	1.15717
entropy_sptac	156976	2.58346	0.180323	2.15702	2.91818	2.45466	2.59619
total_monthly_spend	156819	1488.82	926.403	121.57	4390.58	784.98	1277.33
tag_spend_household	156824	571.313	497.255	5.3	2161.4	180.165	414.56
tag_spend_other_spend	156860	159.475	168.906	-105.94	740	30	103.07
tag_spend_services	160645	239.838	196.456	0	877.72	87.72	187.75
tag_spend_travel	164641	42.7647	68.7196	0	340.58	0	12.25
tag_spend_hobbies	165173	10.3662	17.5902	0	88.99	0	0
tag_spend_retail	156819	68.0719	83.4032	-11.23	369.61	0	35.05
tag_spend_finance	157007	163.048	218.66	0	1023.5	11.315	69.12
tag_spend_communication	164755	48.3293	38.2553	0	161.76	17.115	42.45
tag_spend_motor	165354	40.6831	53.3651	0	211.7	0	10.025

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

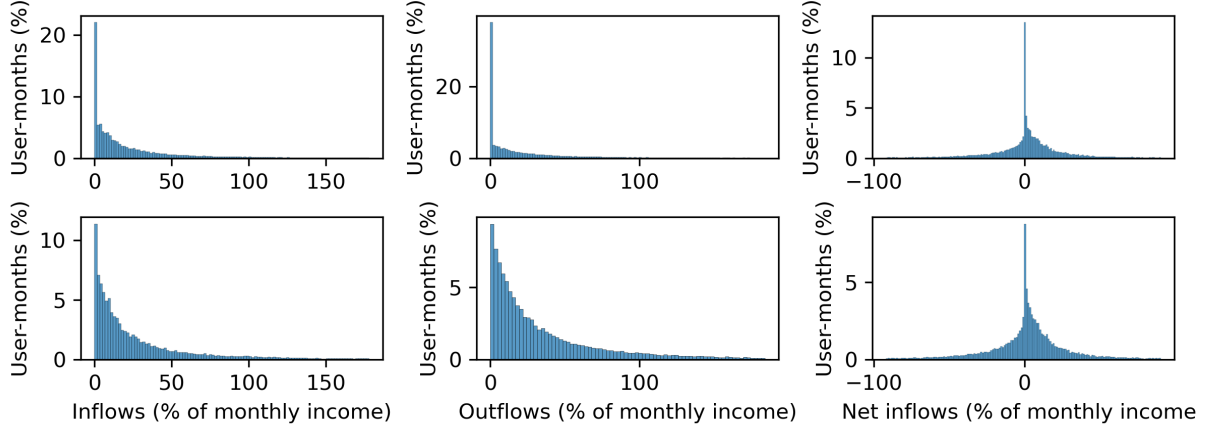
$$p_i = \frac{\text{Count of purchases in } SC_i + 1}{\text{Count of all purchases} + 9} \quad (2)$$

## 4 Methods

### 4.1 Model specification

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

Figure 3: Monthly flows in and out of savings accounts

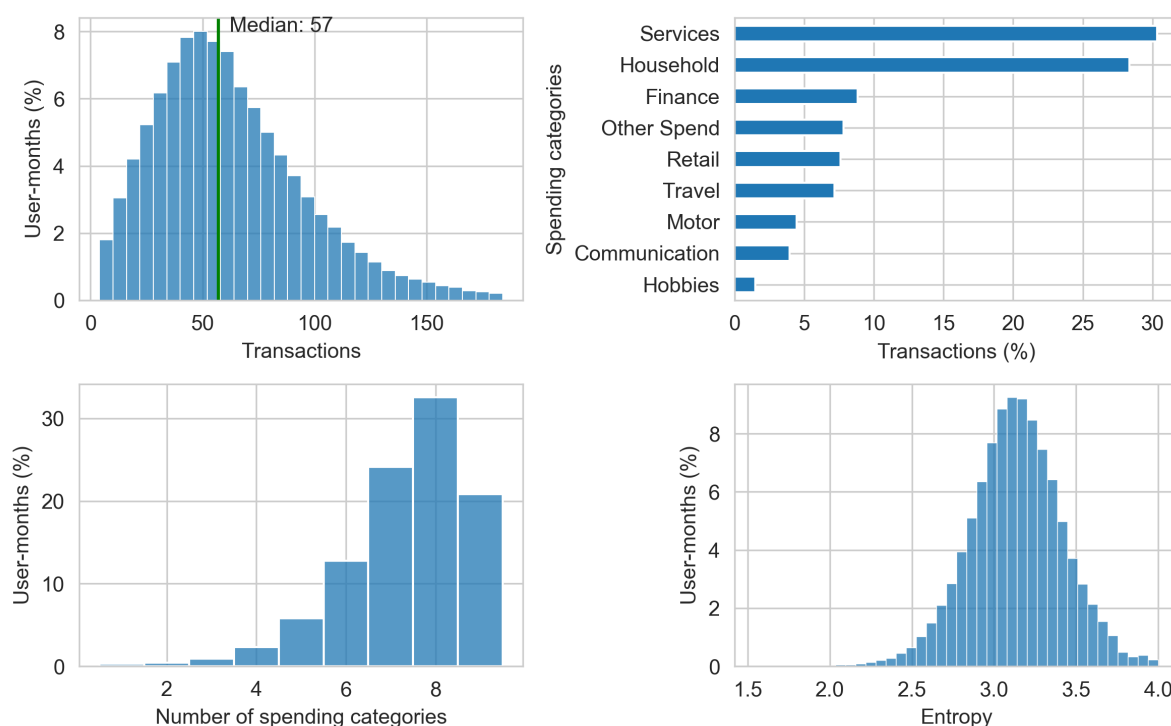


Notes: Flows are calculated for each user-month as the total inflows, outflows, and flows (calculated as inflows - outflows) into all of a user's savings accounts. Zero net flows represent months where inflows are either perfectly balanced by outflows or where there were no flows at all.

$s_{i,t}$  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.

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.

## 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](https://psyarxiv.com/qabgm).
- Shannon, Claude Elwood (1948). “A mathematical theory of communication”. In: *The Bell system technical journal* 27.3, pp. 379–423.