

# Entropy\*

Fabian Gunzinger                      Neil Stewart  
Warwick Business School              Warwick Business School

November 26, 2021

## Contents

<b>1</b>	<b>Introduction</b>	<b>1</b>
<b>2</b>	<b>Data</b>	<b>2</b>
2.1	Preprocessing . . . . .	2
2.2	Sample selection . . . . .	2
2.3	Dataset description . . . . .	2
2.4	Dependent variable . . . . .	3
2.5	Independent variable . . . . .	4
2.6	Understanding entropy . . . . .	5

## 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.

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:

---

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

Independent variable: entropy over categories and others

Outcome variables: first-time saving, average monthly savings

## 2 Data

### 2.1 Preprocessing

Duplicate transactions

### 2.2 Sample selection

Table 1: Sample selection

	Users	Accounts	Transactions	Value (£M)
Raw sample	27,129	132,900	65,871,348	12,194.8
At least 6 months of data	23,878	125,724	65,300,510	12,102.4
At least one current account	22,547	122,158	63,247,036	11,829.8
At least 5 monthly debits totalling GBP200	14,918	79,516	46,440,907	8,681.2
Income payments in 2/3 of all observed months	10,909	61,431	35,593,720	6,736.5
Yearly incomes between 5k and 200k	5,792	30,794	18,704,936	3,278.3
No more than 10 active accounts in any year	5,359	23,151	15,954,877	2,541.3
Debits of no more than 100k in any month	5,004	20,963	14,297,483	1,765.7
Current and savings account balances available	2,716	10,028	7,531,905	956.9
Last account refresh within observed period	2,716	10,021	7,531,838	956.9
Working-age	2,357	8,607	6,792,133	847.3
Final sample	2,357	8,607	6,792,133	847.3

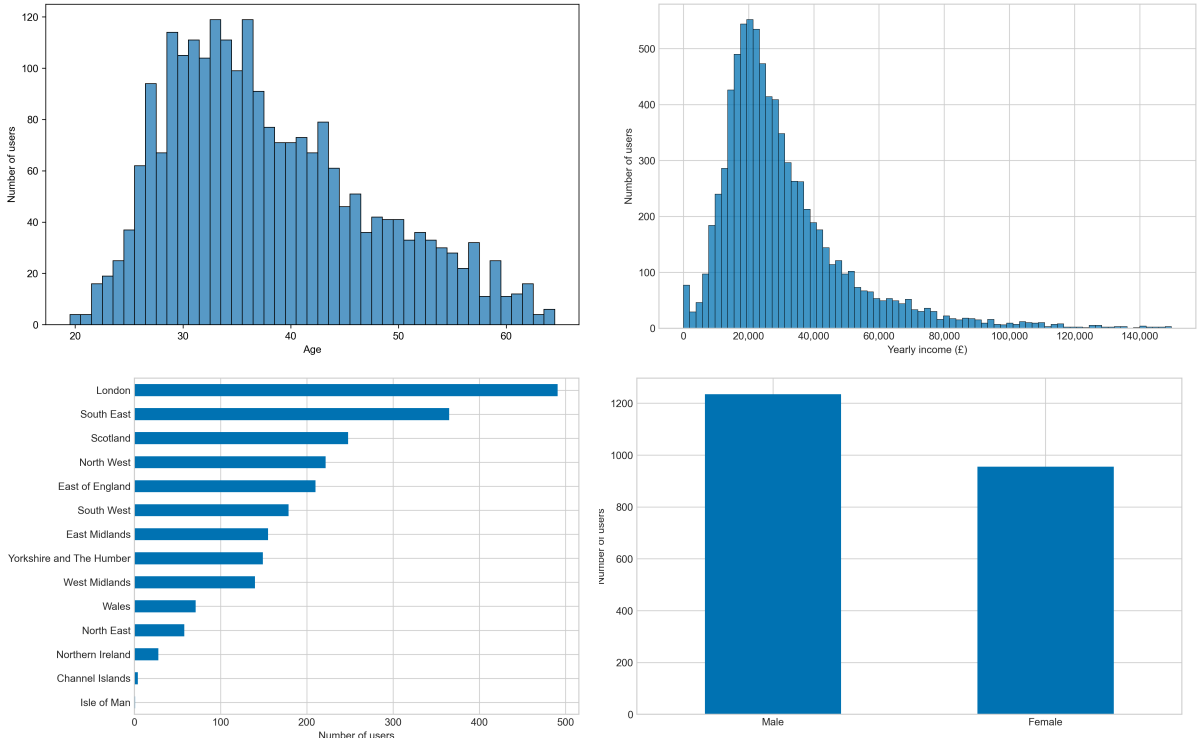
### 2.3 Dataset description

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

Figure 1: Demographic characteristics of Money Dashboard users



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 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.

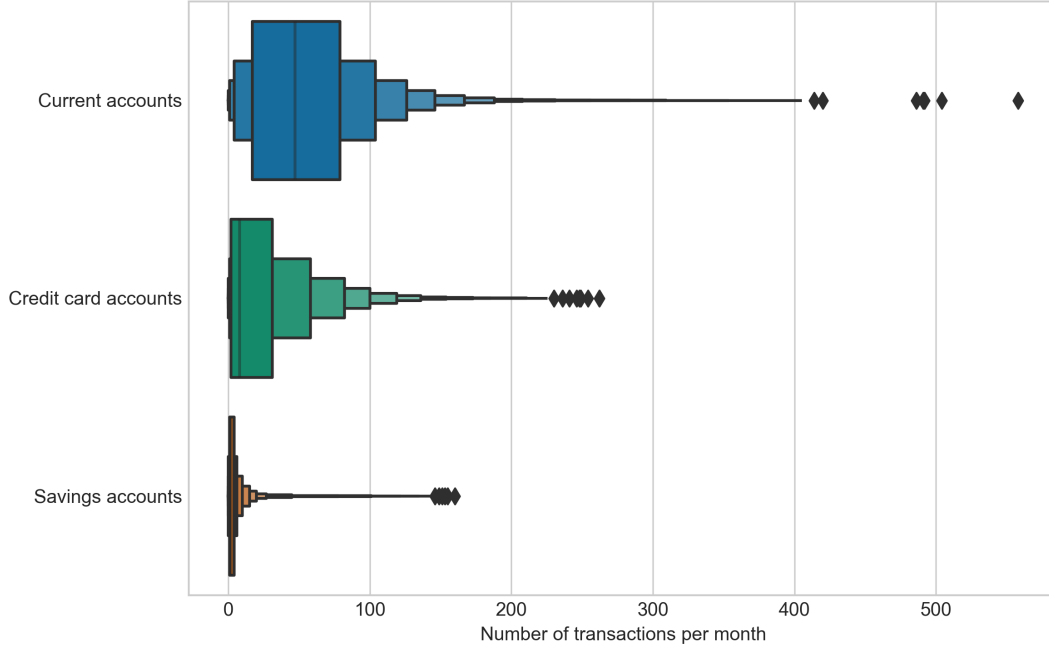
## 2.4 Dependent variable

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)

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

- Pure credit (loans only)
- Wealth held (debit - credit balance)

## 2.5 Independent variable

Spending entropy:

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

$$H = - \sum p_i \log_2(p_i), \quad (1)$$

where  $p_i$  is the probability that an individual makes a purchase in spending category  $i$ . 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.

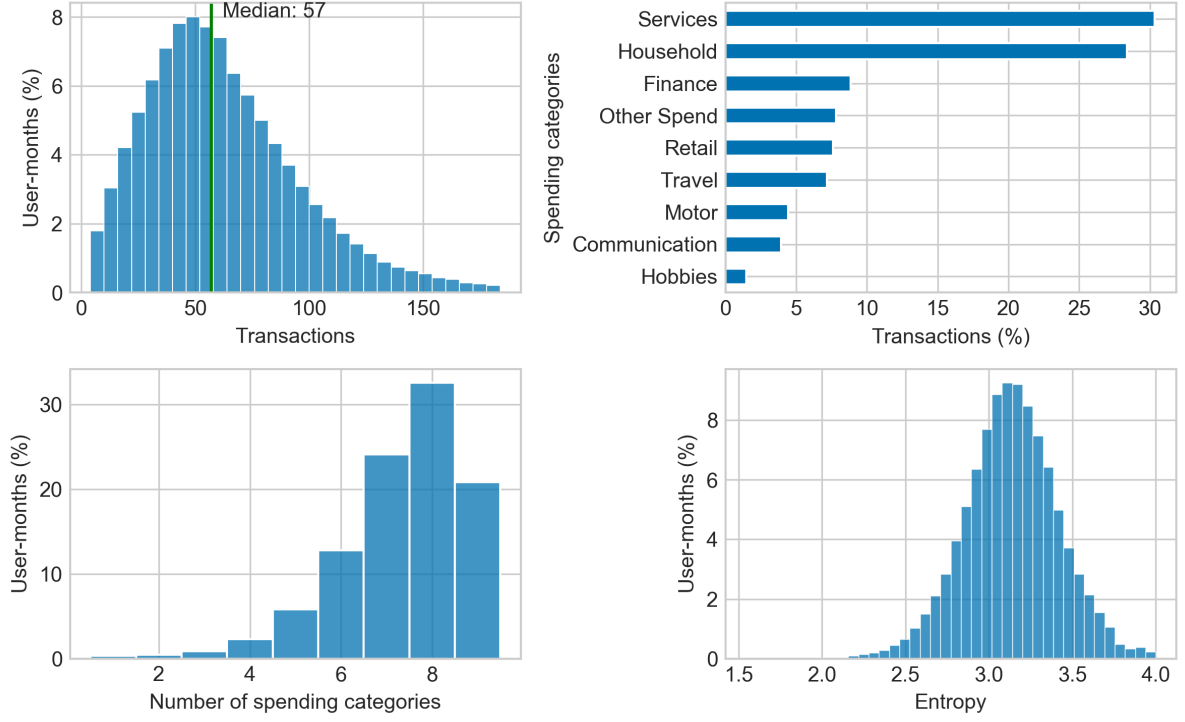
- To calculate individual entropy scores, we group spending into 9 spending categories (SC), based on the classification used by Lloyds Banking Group as discussed in Muggleton et al. (2020).
- Also following that paper, when calculating  $p_i$  we use additive smoothing and add one to the numerator and  $N_{SC}$  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

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

- The right-hand panel in Figure 3 shows the distribution of the resulting entropy scores.

Figure 3: 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.

## 2.6 Understanding entropy

- Entropy is an information theoretic concept that captures the degree of uncertainty when outcomes follow a multi-modal distribution.
- It's calculated as  $H = -\sum p_i \log_2(p_i)$ , where  $p_i$  is the probability of a given event (e.g. making a purchase from a particular category or with a particular merchant).
- Interpretation: low entropy suggests high predictability; high entropy, low predictability.
- Examples:
  - Individual only purchases one type of good ( $p_i = 1$ ):  $H = -(1)(0) = 0$
  - Individual spends uneven proportions on different types of goods ( $p_1 = 1/2, p_2 = 1/4, p_3 = p_4 = 1/8$ ):  $H = -[(1/2)(-1) + (1/4)(-2) + 2(1/8)(-3)] = 1.75$

## 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.