

Fighting Poverty with Data.

Demonstrates that mobile phone metadata can be used to accurately predict the socioeconomic status of individuals and map the geographic wealth in a developing country like Rwanda.

A phone survey with 856 subscribers was combined with a massive dataset of call records.

Label space (y): The wealth index or household consumption derived from survey results of the 856 individuals, representing how poor or wealthy each person was.

Input features (X): Derived from mobile phone usage patterns of those same 856 individuals.

The features included

Call Volume, airtime, mobility, contacts e.t.c

CDRs \rightarrow Call detail records.

The training process included two main key steps:-

1. Feature Engineering: The raw phone log was transformed into several thousand quantitative metrics. These captured:-

- Communication patterns.
- Social network structure
- Mobility and migration

2. Model Training: They trained the model only using the 856 labeled examples and their data and then applied the trained model to millions of other mobile users with their features known (x) but not their y (wealth).

To avoid overfitting an "elastic net" regularization was used. This technique automatically eliminates irrelevant phone metrics and selects a simple and generalizable model that best predicts the wealth index from the survey.

Results and Validation

Individual level \rightarrow model predicted individual wealth with a cross-validated correlation of $r = 0.68$. Also had the ability to predict specific asset ownership like fridges, electricity e.t.c, with high accuracy (AUC up to 0.88)

Geographic level \rightarrow Model was applied to 1.5 million non-surveyed users. Their predicted wealth, ~~and~~ combined ~~income~~ with geographic data from cell tower was used to create high spatial resolution maps of wealth distribution, down to micro-regions of just a few households.

Validation: The results were compared with a 5 year old traditional, more expensive survey, Demographic and Health survey (DHS). It was nearly the same as DHS with a high correlation of 0.92. (0.916 to be exact).

This validated the model's accuracy.

The implications and applications.

The model had the ability to create detailed maps where no other data exists, with profound results.

It was 10x faster and a 1000x cheaper, costing only \$12000 dollars and 4 weeks compared to DHS that took millions of dollars and over a year to complete.

What do we learn from this?

A tiny set of high quality labels + a huge set of informative features can produce population level results.

→ teaches us to value labeled data strategically rather than needing millions of labels.

→ faster and detailed results allowing cost-effective methods to be of use, helping NGOs and governments to provide aid where needed.

Challenge faced:

Mobile phone data is highly personal and sensitive → privacy concerns.

856 surveys small dataset for training, on overfitting.

③ The model trained ~~on~~ data from Rwanda might not apply to other countries.

④ In survey, people might lie or not actually own a phone. Poor people in rural areas might not even own phones.

⑤ Linking survey responses accurately with their metadata is difficult, any mismatch could have corrupted the data labels.