

ANALYZING PODES DATA

Research Report for Faculty of IT, Monash University

Yu (Robert) Fu

Master of Data Science

Monash University

CONTENTS:

- Part I. Background and Introduction**
- Part II. Datasets Description and Pre-processing**
- Part III. Descriptive Statistics**
- Part IV. Schema Matching for PODES 2011-2014**
- Part V. Modelling and Prediction**
- Part VI. Discussion and Next Steps**

Part I. Background and Introduction

[Key Research Question]

In recent years, the Indonesian government has been carrying out major expansions of national grid (PLN) to provide rural villages with access to electricity. The pace and magnitude of each year's expansion are believed to be characterized by village economic and geographical feasibilities. It is of interest, therefore, to identify the "disadvantaged" villages (those unlikely to be connected to PLN soon) and advise with viable alternative power options.

This report details the procedures and results from the initial analysis of multi-year Indonesian Village Surveys (PODES) in understanding factors that identify these "disadvantaged" villages. While working as a winter and on-going part-time research assistant at Faculty of Information Technology at Monash University, I was luckily supervised by Dr. Christoph Bergmeir and Professor Lachlan Andrew, as long as under general guidance of Professor David Taniar and Ariel Liebman, to whom I owe my great thanks for providing me with such a learning experience.

My primary research questions are:

- a. From PODES data, find out what differs in descriptive statistics between PLN-connected villages and those unconnected.
- b. From analyzing survey responses of 2011 and 2014, predict village PLN connect status after

2014.

In order to acquire a certain level of interpretability, I was suggested to use a simple prediction model. In addition, for future simulations, it is recommended that model results be presented with probabilities.

[Main Tasks]

To answer the research questions, the first and foremost step is to record the change of village connection status in a given year. In particular, we would like to know among all disconnected villages in one survey (for instance, 2011), which would change connection status by the next survey (say, 2014). This creates a group of “**connect**” villages that differ from the rest of disconnected ones in a given year, and thus form the basis of our comparison and analysis.

There are three major tasks involved in the analysis that will be detailed in the report. The first part is running descriptive statistics between “**connect**” villages with others - namely, still-yet-to-connected villages, and already-connected ones. The second part tries to create a uniform schema for 2011 and 2014 surveys, and use that for training the predictive model. The third task is to report major findings and discuss for future steps.

[Environment and Packages]

All coding is done using R 3.31 in Jupyter Notebook for easy presentation. Packages used include:

XLConnect – for excel file input

dplyr – for data manipulation

Fselector – for calculating information gain

Glmnet – for penalized regressions

AUC – for calculating AUC and ROC

doMC – for faster parallel caculation.

H2O – for neural network setup

Part II. Datasets Description and Pre-processing

[Datasets]

Using the provided PODES survey from 2008, 2011 and 2014 as the basis of my research, the table below summarizes the number of observations (villages) and questions (variables) in each dataset. A chronological comparison of village IDs shows after each survey, around 90 percent of villages would remain till the next survey, while new villages would comprise about 10 percent in total. This is understandable due to possible administrative causes in real life, and therefore, for the interest of analysis, only villages belonging to both consecutive years (2008-2011, or 2011-2014) are used.

Dataset	# of Villages	# of Variables	# of Villages Remained from Last Survey	# of New Villages
2008	75,410	560	NA	NA
2011	77,961	580	69,873 (93%)	8,088 (10%)
2014	82,190	592	73,951 (89%)	8,239 (10%)

[Major Challenges]

The original datasets present some challenges in both language and schema, so the following treatments are applied.

- a. We used Google Translate to transform the documents from Indonesian to English in order to understand each question. But different wordings may still cause some confusion.
- b. Survey schemas are different for PODES 2011 and 2014. Mapping techniques are applied in order to create a uniform template.
- c. Some variables are missing from the datasets, or contain very different codings compared with survey documents. After confirming with John Paterson, they are excluded from the predictive models.

[Imputation and Pre-processing]

The ordinal and tiered structure of PODES survey questions introduced tremendous level of missing data for both 2011 and 2014. But they can be treated with a simple imputation method. Except for very few variables (e.g., Age), all others with missing values can be imputed with either “NA” or 0. This is because the missing value is indeed a result of “Not Applicable”, or zero in quantity.

In addition, population of villages is added to the 2014 dataset, using the provided “Podes_Tabular_2014.xlsx” by matching village IDs. For both 2011 and 2014 datasets, “**connect**” variable is created by matching village IDs between PODES data 2008 and 2011, as well as 2011 and 2014, respectively, to acquire the villages that would be connected in the later survey.

Finally, all variable types are checked before importing the dataset to ensure they have the right format.

Part III. Descriptive Statistics

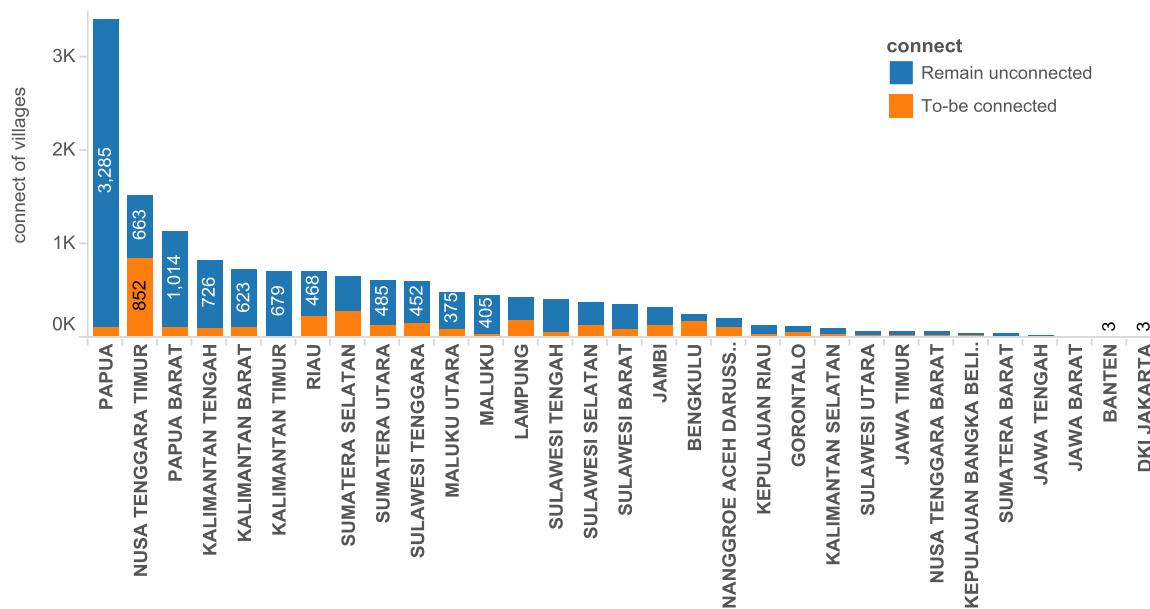
In order to gain an initial idea of how “*connect*” villages are different from the rest of the group, information gain tests are run on both 2008 and 2011 PODES to acquire the independent pair scores between “*connect*” variable with all other variables, and high-ranking pairs are selected for plotting (in Tableau). Though this selection process is rather subjective, it does provide some evidence of potential factors to the outcome variable among the five hundred plus candidates.

The top 100 variables that share highest mutual information with “*connect*” status span across several areas, including: geo-location (i.e., province, district), access to health and educational facilities, road condition and infrastructure, local markets (access to shops, for instance), security conditions (police station or village guards), population characteristics (i.e., size, ethnicity, religion), and revenue from different levels of government. The following section provides descriptive statistics in these highlighted areas between three major village groups of interest, namely:

- 1) villages that already connected to PLN
- 2) “*connect*” villages that are not yet but would connected to PLN by the time of next survey
- 3) villages that would remain unconnected to PLN

Geolocation:

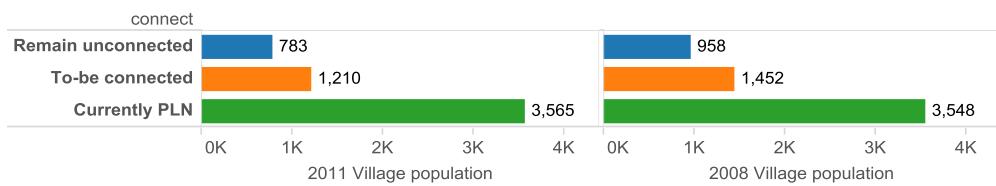
NON-PLN VILLAGES BY PROVINCE IN 2011 PODES SURVEY



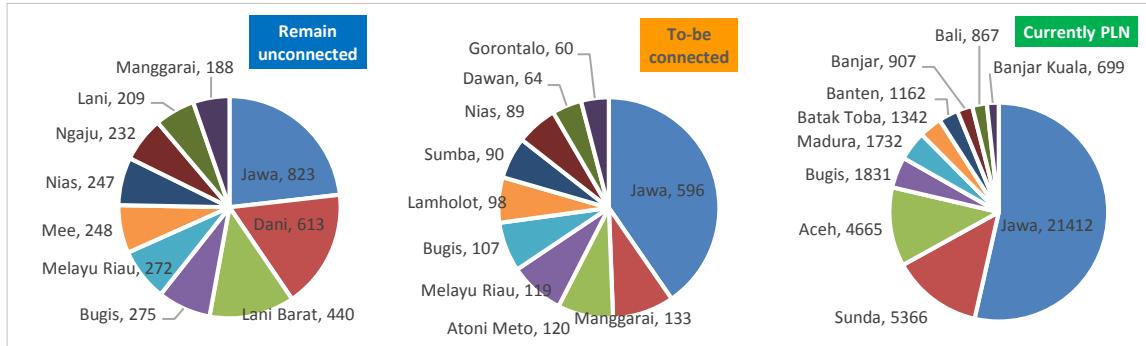
In 2011, Papua has the largest number of villages disconnected to PLN, while Nusa Tenggara Timur has the largest number of “*connect*” villages, showing there would be a great expansion in PLN network between 2011 and 2014. It is believed that province and other geological factors play an important role in PLN connection for the next phase.

Population Characteristics:

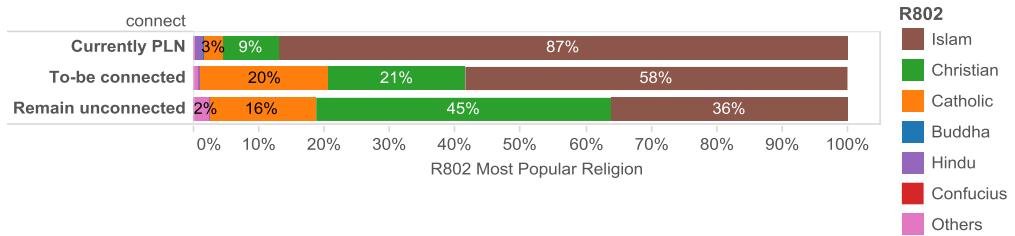
VILLAGE POPULATION SIZE



TOP 10 VILLAGE ETHNICITIES BY CONNECT STATUS IN 2011



VILLAGE PROMINENT RELIGION

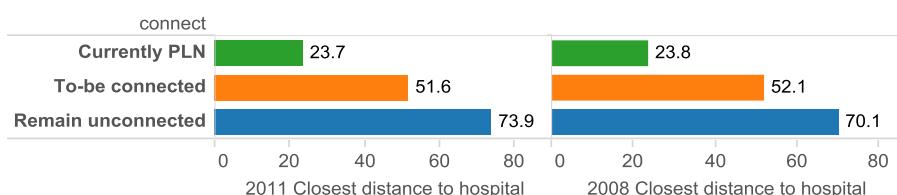


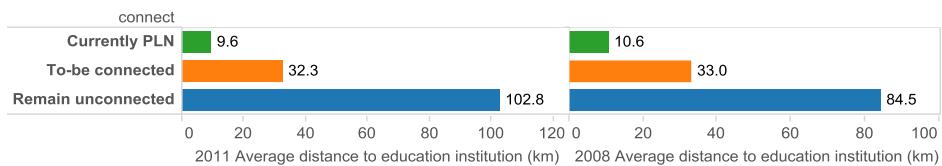
Overall, disconnected villages tend to be smaller in population size, have a greater proportion of ethnic minorities and non-Islamic religions.

Access to Health and Educational Facilities:

It is not surprising to observe the discrepancies of distance to nearby hospitals and education institutions across village types. In general, it seems village connection status are related to the proximity of these facilities.

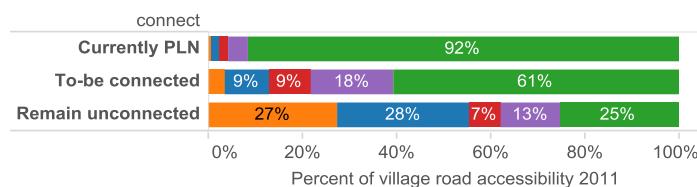
DISTANCE TO HOSPITAL AND EDUCATION INSTITUTIONS





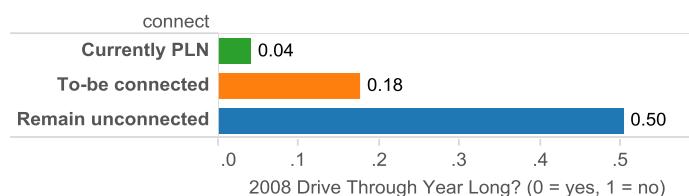
Road Condition and Infrastructure:

MAIN ROAD ACCESSIBILITY



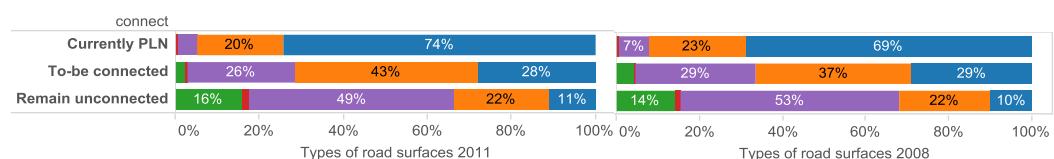
R1002C

- █ Throughout the year
- █ Throughout the year except for certain times
- █ Throughout the year except during rainy season
- █ Inaccessible to wheeled vehicles throughout the year
- █ NA



Disconnected villages have less accessible and lower level of infrastructure compared to “**connect**” and PLN-established villages. This could be an indicator that they are located in more remote areas.

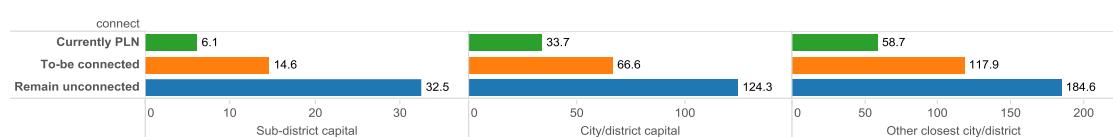
TYPE OF ROAD SURFACE

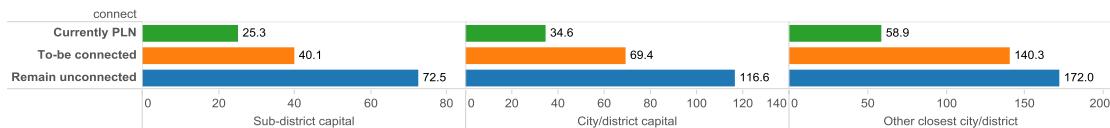


R1001B1

- █ Asphalt / Concrete
- █ Hardened(gravel, stones,etc)
- █ Soil
- █ Other
- █ NA

DISTANCE TO NEARBY CAPITALS





The distance to nearby administrative capitals suggests supporting evidence of their remote location. In particular, the distance between regional capitals and disconnected villages are twice as far as the distance to “**connect**” villages, and four times as far as PLN-connected villages.

Local Markets:

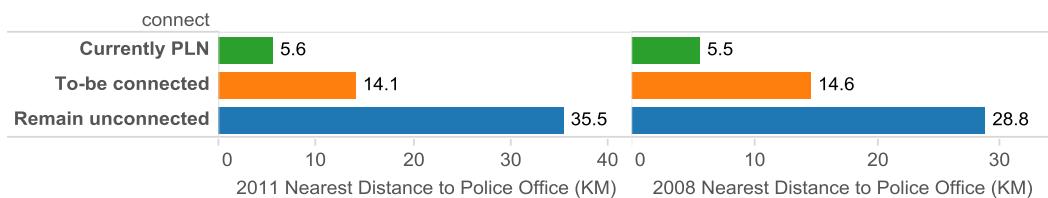
DISTANCE TO MARKETS



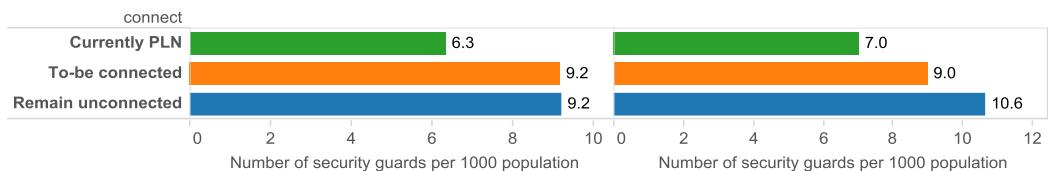
The distance and accessibility to markets and shops suggest the same thing. In general, disconnected villages are twice and four times further away compared the other types of villages.

Security:

DISTANCE TO POLICE OFFICE



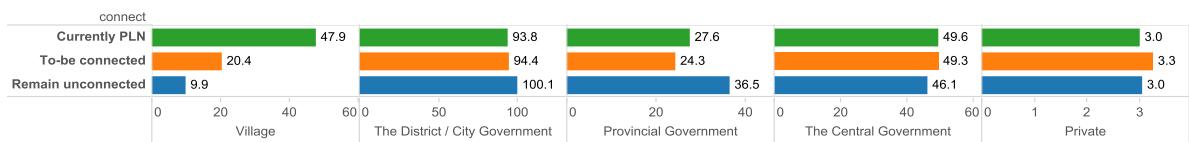
2011 (LEFT) AND 2008 (RIGHT) SECURITY GUARDS RATIO



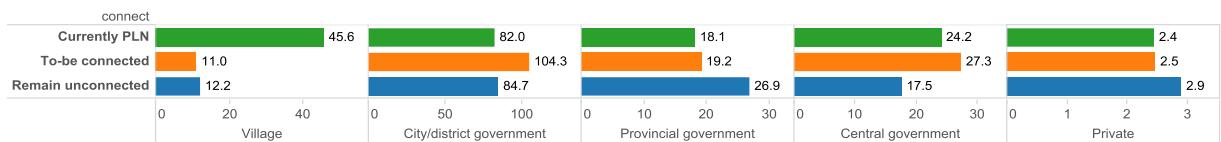
Disconnected villages are also further away from police offices. This might be related to their relatively smaller population sizes and different social context required to keep order. Instead, they have a higher ratio of security guards among population.

Revenue from Different Levels of Government:

2011 VILLAGE REVENUE BY SOURCE (MILLION RP)



2008 VILLAGE REVENUE BY SOURCE (MILLION RP)



There is a drastic difference in village-generated revenue between PLN-established villages and disconnected ones. Although the causal relationship between village connection status and revenue generated or received can be convoluted in nature (it could be very likely that connected villages tend to generate more income), this graph suggests that they might be somehow related to each other. Besides, it is also noted that disconnected villages tend to receive more financial support from the provincial government, but not at sub-district or central level. Considering their smaller population size, this benefit could translate into high number for each person.

To sum up from the descriptive statistics, there seems to be major differences in multiple social and economic aspects between villages that would be connected to PLN, those never connected, and those already connected. In particular, disconnected villages tend to be distributed unevenly across provinces, they have smaller population and greater proportion of religious or ethnical minorities, generally located further away from school, hospital, markets or regional capitals, and have lower grade of transportation and infrastructure. In the end, they seem to be poorer and receive more government support than self-generated incomes.

Part IV. Schema Matching for PODES 2011 and 2014

In order to make sensible predictions for future “*connect*” villages after 2014, a unified schema is required to train and test the mode using 2011 PODES data. While PODES 2011 and 2014 share most of common variables, there are some with different coding or types of answers. Overall, 417 variables from 585 (2011) and 594(2014) are retained in an effort to acquire as much information as possible. It is important to note that the matching process is conducted on PODES 2014 and 2011 after all PLN-connected villages are excluded. The following section reports the major changes to the original dataset by areas of variables.

1. Village Identification

<i>Global Schema</i>	<i>Correspondent 2011</i>	<i>Correspondent 2014</i>	<i>Comment</i>
<i>villageID</i>	id2011	id2013	Match
<i>provincelD</i>	KODE_PROV	R101	Recode 2014
<i>districtID</i>	KODE_KAB	R102	Match

<i>subDisID</i>	KODE_KEC	R103	Recode 2014
-----------------	----------	------	-------------

ProvinceID is used instead of Province names because it presents less inconsistency between the two datasets. The one unique province for PODES 2014 “KALIMANTAN UTARA”, which created in 2012 from province “KALIMANTAN TIMUR”, is recoded as the latter to keep consistency. There are also 17 sub-districts corresponding to 108 villages unique in PODES 2014, mostly located in Papua. Since it requires tremendous domain knowledge to effectively impute them with known sub-districts, a decision is made to delete those villages so that we can train the model properly.

2. General Information of Village

<i>Global Schema</i>	<i>Correspondent 2011</i>	<i>Correspondent 2014</i>	<i>Comment</i>
<i>villageStatus</i>	R301	R301	Recode 2014
<i>SLS</i>	R304A	R304A	Recode 2011, 2014
<i>SLS_N</i>	R304C	304BA(BCD)_K5	Recode 2014
<i>slope</i>	R305A	305B	Recode 2011
<i>atSea</i>	R305D	307A	Match
<i>atSeaUse_1</i>	R305E2A	307B1A	Match
<i>atSeaUse_2</i>	R305E2B	R307B1B	Match
<i>atSeaUse_3</i>	R305E2C	R307B1C	Match
<i>atSeaUse_4</i>	R305E2D	R307B1D	Match
<i>atSeaUse_5</i>	R305E2E	R307B1E	Match
<i>atSeaMangrove</i>	R305E3	R307B2	Match
<i>byForest</i>	R306A	R308A	Match
<i>byForestUse</i>	R306B	R308B	Match

There are minor coding differences between PODES 2011 and 2014 in village status, the number of SLS facilities, and topography. They are however straightforward and easy for recoding.

3. Population and Employment

<i>Global Schema</i>	<i>Correspondent 2011</i>	<i>Correspondent 2014</i>	<i>Comment</i>
<i>migrantM</i>	R402A	403B1	Match
<i>migrantF</i>	R402B	403B2	Match
<i>incomeSource</i>	R403A	404A	Match
<i>agriType</i>	R403B	404B1	Match

4. Housing and Environment

<i>Global Schema</i>	<i>Correspondent 2011</i>	<i>Correspondent 2014</i>	<i>Comment</i>
<i>strLight</i>	R502A	R502A	Match
<i>strLightEngy</i>	R502B	R502B	Match
<i>fuelCook</i>	R503	R503	Match

<i>toilet</i>	R504	R504	Match
<i>litter</i>	R505A	R505A	Match
<i>pubTrash</i>	R505B	R505B	Match
<i>river</i>	R506AK2	R508A_K2	Match
<i>riverUse_1</i>	R506B1K2	R508B1_K2	Match
<i>riverUse_2</i>	R506B2K2	R508B2_K2	Match
<i>riverUse_3</i>	R506B3K2	R508B3_K2	Match
<i>riverUse_4</i>	R506B4K2	R508B4_K2	Match
<i>riverUse_5</i>	R506B5K2	R508B6_K2	Match
<i>irrigation</i>	R506AK3	R508A_K3	Match
<i>lake</i>	R506AK4	R508A_K4	Match
<i>irrigationUse_1</i>	R506B1K3	R508B1_K3	Match
<i>lakeUse_1</i>	R506B1K4	R508B1_K4	Match
<i>irrigationUse_2</i>	R506B2K3	R508B2_K3	Match
<i>lakeUse_2</i>	R506B2K4	R508B2_K4	Match
<i>irrigationUse_3</i>	R506B3K3	R508B3_K3	Match
<i>lakeUse_3</i>	R506B3K4	R508B3_K4	Match
<i>irrigationUse_4</i>	R506B4K3	R508B4_K3	Match
<i>lakeUse_4</i>	R506B4K4	R508B4_K4	Match
<i>lakeUse_5</i>	R506B5K4	R508B6_K4	Match
<i>byRiver</i>	R508A	509B	Match
<i>byRiver_village</i>	R508B	R509C1	Match
<i>byRiver_house</i>	R508C	R509C2	Match
<i>byRiver_family</i>	R508D	R509C3	Match
<i>bySUTET</i>	R509A	R510A	Match
<i>bySUTET_village</i>	R509B1	R510C1	Match
<i>bySUTET_house</i>	R509B2	R510C2	Match
<i>bySUTET_family</i>	R509B3	R510C3	Match
<i>bySlum</i>	R510A	R511A	Match
<i>bySlum_village</i>	R510B1	R511B1	Match
<i>bySlum_house</i>	R510B2	R511B2	Match
<i>bySlum_family</i>	R510B3	R511B3	Match
<i>pollution_1</i>	R511AK2	R512A_K2	Match
<i>pollution_1Src</i>	R511AK3	R512A1_K3	Match
<i>pollution_1Cmplt</i>	R511AK4	R512A_K4	Match
<i>pollution_2</i>	R511BK2	R512B_K2	Match
<i>pollution_2Src</i>	R511BK3	R512B1_K3	Match
<i>pollution_2Cmplt</i>	R511BK4	R512B_K4	Match
<i>pollution_3</i>	R511CK2	R512C_K2	Match
<i>pollution_3Src</i>	R511CK3	R512C1_K3	Match
<i>pollution_3Cmplt</i>	R511CK4	R512C_K4	Match
<i>burnField</i>	R512A	R513	Match
<i>excavation</i>	R513	R514	Match

Despite some wording difference, most questions in this part of the survey are common in both PODES datasets, and are matched in coding.

5. Natural Disaster

<i>Global Schema</i>	<i>Correspondent 2011</i>	<i>Correspondent 2014</i>	<i>Comment</i>
<i>disaster_1</i>	R60101K2	R601A_K2	Match
<i>disaster_1N</i>	R60101K3	R601A_K3+K5+K7	Recode 2014
<i>disaster_1Dead</i>	R60101K4	R601A_K4+K6+K8	Recode 2014
<i>disaster_2</i>	R60102K2	R601B_K2	Match
<i>disaster_2N</i>	R60102K3	R601B_K3+K5+K7	Recode 2014
<i>disaster_2Dead</i>	R60102K4	R601B_K4+K6+K8	Recode 2014
<i>disaster_3</i>	R60103K2	R601C_K2	Match
<i>disaster_3N</i>	R60103K3	R601C_K3+K5+K7	Recode 2014
<i>disaster_3Dead</i>	R60103K4	R601C_K4+K6+K8	Recode 2014
<i>disaster_4</i>	R60104K2	R601D_K2	Match
<i>disaster_4N</i>	R60104K3	R601D_K3+K5+K7	Recode 2014
<i>disaster_4Dead</i>	R60104K4	R601D_K4+K6+K8	Recode 2014
<i>disaster_5</i>	R60105K2	R601E_K2	Match
<i>disaster_5N</i>	R60105K3	R601E_K3+K5+K7	Recode 2014
<i>disaster_5Dead</i>	R60105K4	R601E_K4+K6+K8	Recode 2014
<i>disaster_6</i>	R60106K2	R601F_K2	Match
<i>disaster_6N</i>	R60106K3	R601F_K3+K5+K7	Recode 2014
<i>disaster_6Dead</i>	R60106K4	R601F_K4+K6+K8	Recode 2014
<i>disaster_7</i>	R60107K2	R601G_K2	Match
<i>disaster_7N</i>	R60107K3	R601G_K3+K5+K7	Recode 2014
<i>disaster_7Dead</i>	R60107K4	R601G_K4+K6+K8	Recode 2014
<i>disaster_8</i>	R60108K2	R601H_K2	Match
<i>disaster_8N</i>	R60108K3	R601H_K3+K5+K7	Recode 2014
<i>disaster_8Dead</i>	R60108K4	R601H_K4+K6+K8	Recode 2014
<i>disaster_9</i>	R60109K2	R601I_K2	Match
<i>disaster_9N</i>	R60109K3	R601I_K3+K5+K7	Recode 2014
<i>disaster_9Dead</i>	R60109K4	R601I_K4+K6+K8	Recode 2014
<i>disaster_10</i>	R60110K2	R601J_K2	Match
<i>disaster_10N</i>	R60110K3	R601J_K3+K5+K7	Recode 2014
<i>disaster_10Dead</i>	R60110K4	R601J_K4+K6+K8	Recode 2014

Both PODES surveys have 10 disaster types. However, in 2014 the number of occurrences and casualties are broken down into three years (2011-2013), while in PODES 2011, they are answered in three-year's total. Therefore, we add columns in PODES 2014 to match with 2011.

6. Education and Health

<i>Global Schema</i>	<i>Correspondent 2011</i>	<i>Correspondent 2014</i>	<i>Comment</i>
<i>edu_1State</i>	R701AK2	R701A_K2	Match
<i>edu_1Private</i>	R701AK3	R701A_K3	Match
<i>edu_1Distance</i>	R701AK4	R701A_K4	Match
~	~	~	~
<i>edu_5State</i>	R701EK2	R701E_K2	Match
<i>edu_5Private</i>	R701EK3	R701E_K3	Match
<i>edu_5Distance</i>	R701EK4	R701E_K4	Match
<i>edu_6State</i>	R701FK2	R701F_K2	Match
<i>edu_6Private</i>	R701FK3	R701F_K3	Match
<i>edu_7State</i>	R701GK2	R701G_K2	Match
<i>edu_7Private</i>	R701GK3	R701G_K3	Match
<i>edu_8Private</i>	R701HK3	R701H_K3	Match
<i>edu_9Private</i>	R701IK3	R701I_K3	Match
<i>edu_10Private</i>	R701JK3	R701J_K3	Match
<i>voc_1N</i>	R702A	R703A	Match
~	~	~	~
<i>voc_7N</i>	R702G	R703G	Match
<i>otherEdu_1</i>	R703A	R702A	Match
~	~	~	~
<i>otherEdu_5</i>	R703E	R702F	Recode 2014
<i>health_1</i>	R704AK2	R704A_K2	Match
<i>health_1N</i>	R704AK3	R704A_K3	Match
<i>health_1Distance</i>	R704AK4	R704A_K4	Match
<i>health_1EasytoTravel</i>	R704AK5	R704A_K5	Match
~	~	~	Match
<i>health_3</i>	R704CK2	R704F_K2	Match
<i>health_3N</i>	R704CK3	R704F_K3	Match
<i>health_3Distance</i>	R704CK4	R704F_K4	Match
<i>health_3EasytoTravel</i>	R704CK5	R704F_K5	Match
<i>health_4</i>	R704DK2	R704C_K2 R704D_K2	Recode 2014
<i>health_4N</i>	R704DK3	R704C_K3 + R704D_K3	Recode 2014
<i>health_4Distance</i>	R704DK4	min(R704C_K4, R704D_K4)	Recode 2014
<i>health_4EasytoTravel</i>	R704DK5	round(mean(R704C_ K5, R704D_K5))	Recode 2014
<i>health_5</i>	R704EK2	R704E_K2	Match
<i>health_5N</i>	R704EK3	R704E_K3	Match
<i>health_5Distance</i>	R704EK4	R704E_K4	Match
<i>health_5EasytoTravel</i>	R704EK5	R704E_K5	Match
<i>health_6</i>	R704FK2	R704G_K2	Match
<i>health_6N</i>	R704FK3	R704G_K3	Match

<i>health_6Distance</i>	R704FK4	R704G_K4	Match
<i>health_6EasytoTravel</i>	R704FK5	R704G_K5	Match
<i>health_7</i>	R704GK2	R704H_K2	Match
<i>health_7N</i>	R704GK3	R704H_K3	Match
<i>health_7Distance</i>	R704GK4	R704H_K4	Match
<i>health_7EasytoTravel</i>	R704GK5	R704H_K5	Match
<i>health_8</i>	R704HK2	R704I_K2	Match
<i>health_8N</i>	R704HK3	R704I_K3	Match
<i>health_8Distance</i>	R704HK4	R704I_K4	Match
<i>health_8EasytoTravel</i>	R704HK5	R704I_K5	Match
<i>health_9</i>	R704IK2	R704J_K2	Match
<i>health_9N</i>	R704IK3	R704J_K3	Match
<i>health_9Distance</i>	R704IK4	R704J_K4	Match
<i>health_9EasytoTravel</i>	R704IK5	R704J_K5	Match
<i>health_10</i>	R704JK2	R704K_K2	Match
<i>health_10N</i>	R704JK3	R704K_K3	Match
<i>health_11</i>	R704KK2	R704L_K2	Match
<i>health_11N</i>	R704KK3	R704L_K3	Match
<i>health_11Distance</i>	R704KK4	R704L_K4	Match
<i>health_11EasytoTravel</i>	R704KK5	R704L_K5	Match
<i>health_12</i>	R704LK2	R704M_K2	Match
<i>health_12Distance</i>	R704LK4	R704M_K4	Match
<i>health_12EasytoTravel</i>	R704LK5	R704M_K5	Match
<i>posyandu_N</i>	R705B	R705A	Match
<i>posyandu_N2</i>	R705C	R705B	Match
<i>doctorM_N</i>	R707A1	R706A1	Match
<i>doctorF_N</i>	R707A2	R706A2	Match
<i>dentist_N</i>	R707B	R706B	Match
<i>midwife_N</i>	R707C	R706C	Match
<i>otherHealth_N</i>	R707D	R706D	Match
<i>TBA_N</i>	R707E	R708	Match
<i>disease_1</i>	R708AK2	R709A_K2	Match
<i>disease_1N</i>	R708AK3	R709A_K3	Match
<i>disease_1Died</i>	R708AK4	R709A_K4	Match
~	~	~	Match
<i>disease_6</i>	R708HK2	R709H_K2	Match
<i>disease_6N</i>	R708HK3	R709H_K3	Match
<i>disease_6Died</i>	R708HK4	R709H_K4	Match
<i>mulnutrition_N</i>	R709	R710	Match
<i>JAMKESMAS_N</i>	R711	R711A	Match
<i>SKTM_N</i>	R712	R711B	Match

Both datasets contain identical questions on ten common education institutions, seven types of vocational schools, and four types of special education schools (PODES 2014 has five and slightly

different coding scheme). They are broadly matched with each other.

For health facilities, it appears one particular type (PHC) from PODES 2011 is spited into two different types in PODES 2014 (PHC with hospitalization and those without). To retain most information, arithmetic treatment is carried to extract relevant information. All other eleven types are matched perfectly.

Six out of eight outbroken diseases in PODES 2011 are shared in PODES 2014. We retain the common ones and ignore the others.

7. Social and Cultural

<i>Global Schema</i>	<i>Correspondent 2011</i>	<i>Correspondent 2014</i>	<i>Comment</i>
<i>religion_1</i>	R80101	R80101	Match
~	~	~	Match
<i>religion_7</i>	R80107	R80107	Match
<i>mostReligion</i>	R802	R802	Recode 2011
<i>worship_1</i>	R803A	R803A	Match
~	~	~	Match
<i>worship_8</i>	R803H	R803H	Match
<i>disabled_1</i>	R805A	R805A	Match
~	~	~	Match
<i>disabled_9</i>	R805I	R805I	Match
<i>multiEthnicity</i>	R806	R804A1	Match
<i>commService</i>	R808	R807A	Match

Both datasets share identical questions in regard to seven types of religions, eight types of worship places and nine types of disabilities. The only recode is on the “most popular religion” in the village, where 3 occurrences of “Buddha” in PODES 2011 is recoded to “Other” to match with pattern in PODES 2014 (where there is no “Buddha”). It is also noted that PODES 2014 has not reported major ethnicities.

8. Entertainment

<i>Global Schema</i>	<i>Correspondent 2011</i>	<i>Correspondent 2014</i>	<i>Comment</i>
<i>cinema</i>	R901A	R902A1	Match
<i>cinemaDistance</i>	R901B	R902A2	Match
<i>pub</i>	R902A	R902B1	Match
<i>pubDistance</i>	R902B	R902B2	Match
<i>sport_1</i>	R903AK2	R903A_K2	Match
<i>sport_1Group</i>	R903AK3	R903A_K3	Match
~	~	~	Match
<i>sport_7</i>	R903GK2	R903H_K2	Match
<i>sport_7Group</i>	R903GK3	R903H_K3	Match

<i>sport_8Group</i>	R903HK3	R903F_K3	Match
<i>sport_9Group</i>	R903IK3	R903I_K3	Match
<i>sport_10Group</i>	R903JK3	R903J_K3	Match

All common variables matched.

9. Transportation, Communication and Information

<i>Global Schema</i>	<i>Correspondent 2011</i>	<i>Correspondent 2014</i>	<i>Comment</i>
<i>infrastructure</i>	R1001A	R1001A	Match
<i>roadType</i>	R1001B1	R1001B1	Match
<i>roadThrough</i>	R1001B2	R1001B2	Recode 2014
<i>distance_1</i>	R1004AK2	R1002A_K2	Match
<i>distance_1Transport</i>	R1004AK5	R1002A_K6	Recode 2014
<i>distance_2</i>	R1004BK2	R1002B_K2	Match
<i>distance_2Transport</i>	R1004BK5	R1002B_K6	Recode 2014
<i>distance_3</i>	R1004CK2	R1002D_K2	Match
<i>distance_3Transport</i>	R1004CK5	R1002D_K6	Recode 2014
<i>landPhone</i>	R1005A	R1003A	Match
<i>landPhone_N</i>	R1005B	R1003B	Match
<i>payPhone</i>	R1006	R1004A	Match
<i>BTS</i>	R1007A	R1005A	Match
<i>cellPhone</i>	R1007B	R1005B	Match
<i>telecomm</i>	R1008	R1006	Match
<i>internet</i>	R1009	R1007B	2014 recode to 1 or 2
<i>postOffice</i>	R1010A	R1008A	Match
<i>mobilePO</i>	R1011	R1008B	Match
<i>TV_1</i>	R1012A	R1009B_K2	Match
<i>TV_2</i>	R1012B	R1009C_K2	2011 recode to 1 or 2
<i>TV_3</i>	R1012C	R1009A_K2	2011 recode to 1 or 2
<i>TV_4</i>	R1012D	R1009D_K2	2011 recode to 1 or 2

The road infrastructure variables are mostly shareable between the datasets, however there are some major recordings carried out to keep PODES 2014 consistent with 2011. The different wordings from the original document (in Indonesian) for “*distance_Transport*” are tricky for decision, and therefore requires some subjective judgment. Besides, all other common variables in this section are matched with only small discrepancies.

10. Land Use

<i>Global Schema</i>	<i>Correspondent 2011</i>	<i>Correspondent 2014</i>	<i>Comments</i>
<i>landChg_1</i>	R1103AK4	R1103A_K4	Match
<i>landChg_1pct</i>	R1103AK5	R1103A_K5	Match
~	~	~	Match

<i>landChg_6</i>	R1103CK4	R1103C_K4	Match
<i>landChg_6pct</i>	R1103CK5	R1103C_K5	Match

All common variables matched.

11. Economy

<i>Global Schema</i>	<i>Correspondent 2011</i>	<i>Correspondent 2014</i>	<i>Comments</i>
<i>KUD</i>	R1201A	R1213A	Match
<i>nonKUD</i>	R1201B	R1213B	Match
<i>industry_1</i>	R1202A	R1201A	Match
<i>~</i>	<i>~</i>	<i>~</i>	Match
<i>industry_8</i>	R1202H	R1201H	Match
<i>gpShop</i>	R1203A	R1202A	Match
<i>gpShopDis</i>	R1203B	R1202B	Match
<i>kerosene</i>	R1204A	R1203A	Match
<i>LPG</i>	R1204B	R1203B	Match
<i>prmntMkt</i>	R1205A	R1204A, R1204B	Match
<i>prmntMktDis</i>	R1205B	R1204C	Match
<i>fltMkt</i>	R1206	R1205	Match
<i>miniMkt</i>	R1207	R1206	Match
<i>grocery</i>	R1208	R1207	Match
<i>tavern</i>	R1209	R1208	Match
<i>restaurant</i>	R1210	R1209	Match
<i>hotel</i>	R1211	R1210	Match
<i>lodge</i>	R1212	R1211	Match
<i>cooperative_1</i>	R1213A	R1212A	Match
<i>cooperative_2</i>	R1213B	R1212B	Match
<i>cooperative_3</i>	R1213C	R1212C	Match
<i>cooperative_4</i>	R1213D	R1212D	Match
<i>creditFcity_1</i>	R1214A	R1214A	Match
<i>creditFcity_2</i>	R1214C	R1214C	Match
<i>bank_1</i>	R1215AK2	R1215A_K2, R1215B_K2	Match
<i>bank_1N</i>	R1215AK3	R1215A_K2, R1215B_K2	Match
<i>bank_2</i>	R1215BK2	R1215C_K2	Match
<i>bank_2N</i>	R1215BK3	R1215C_K3	Match

All common variables in this section are matched, with only small change of PODES 2014 (by collapsing two variables into one) to keep consistent with 2011.

12. Security

<i>Global Schema</i>	<i>Correspondent 2011</i>	<i>Correspondent 2014</i>	<i>Comments</i>
<i>massFght</i>	R1301A	R1301A	Match
<i>fght_1N</i>	R1301B1K2	R1301B1_K2	Match

<i>fght_1Dead</i>	R1301B1K3	R1301B1_K3	Match
<i>fght_1Injured</i>	R1301B1K4	R1301B1_K4	Match
	R1301B1K5	R1301B1_K4	Not used: 2014 messy answer
~	~	~	Match
<i>fght_7N</i>	R1301B7K2	R1301B7_K2	Match
<i>fght_7Dead</i>	R1301B7K3	R1301B7_K3	Match
<i>fght_7Injured</i>	R1301B7K4	R1301B7_K4	Match
	R1302A	R1302A	Not used: 2014 unmatchable answer
	R1302B	R1302B	Not used: 2014 messy answer
<i>criminal_1</i>	R130301K2	R1303A01K3	Match
<i>criminal_1Tndcy</i>	R130301K3	R1303A01K4	Match
~	~	~	Match
<i>criminal_10</i>	R130310K2	R1303A10K3	Match
<i>criminal_10Tndcy</i>	R130310K3	R1303A10K4	Match
<i>suicide</i>	R1304	R1307	Match
<i>streetChild</i>	R1306	R1308A	Recode 2011
<i>sexWorker</i>	R1307	R1309	Match
<i>secActivities_1</i>	R1308A	R1304A	Match
<i>secActivities_2</i>	R1308B	R1304B	Match
<i>secActivities_3</i>	R1308C	R1304C	Match
<i>secActivities_4</i>	R1308D	R1304D	Match
<i>police</i>	R1309BK2	R1306A	Match
<i>policeDist</i>	R1309BK3	R1306B1	Match
<i>policeEasytoTravel</i>	R1309BK4	R1306B2	Match
<i>guards</i>	R1310	R1305	Match

All common variables are generally matched in this section, except for minor recodings. It is noted that some variables, though matched according to survey document, presents vague or messy answers in the datasets. They are therefore not used for global schema.

13. Autonomy

<i>Global Schema</i>	<i>Correspondent 2011</i>	<i>Correspondent 2014</i>	<i>Comments</i>
<i>income_1Form</i>	R1401AK2	R1501A_K2	Match
<i>income_1</i>	R1401AK3	R1501A_K3	Match
~	~	~	Match
<i>income_7Form</i>	R1401B6K2	R1501C6_K2	Match
<i>income_7</i>	R1401B6K3	R1501C6_K3	Match
<i>dev_1</i>	R1402A1K2	R1401A1_K2	Match
<i>dev_1Source</i>	R1402A1K3	R1401A1_K3	Match
~	~	~	Match

<i>dev_4</i>	R1402A4K2	R1401A4_K2	Match
<i>dev_4Source</i>	R1402A4K3	R1401A4_K3	Match
<i>fund_1</i>	R1402B1K2	R1401B1_K2	Match
<i>fund_1Source</i>	R1402B1K3	R1401B1_K3	Match
~	~	~	Match
<i>fund_3</i>	R1402B3K2	R1401B3_K2	Match
<i>fund_3Source</i>	R1402B3K3	R1401B3_K3	Match
<i>progm_1</i>	R1402C1K2	R1401C1_K2	Match
<i>progm_1Source</i>	R1402C1K3	R1401C1_K3	Match
~	~	~	Match
<i>progm_3</i>	R1402C3K2	R1401C3_K2	Match
<i>progm_3Source</i>	R1402C3K3	R1401C3_K3	Match

All common variables matched.

14. Village Administration

<i>Global Schema</i>	<i>Correspondent 2011</i>	<i>Correspondent 2014</i>	<i>Comments</i>
<i>head</i>	R1501AK2	R1601A_K2	Match
<i>headAge</i>	R1501AK3	R1601A_K3	Match
<i>headSex</i>	R1501AK4	R1601A_K4	Match
<i>headEdu</i>	R1501AK5	R1601A_K5	Recode 2014
<i>sec</i>	R1501BK2	R1601B_K2	Match
<i>secAge</i>	R1501BK3	R1601B_K3	Match
<i>secSex</i>	R1501BK4	R1601B_K4	Match
<i>secEdu</i>	R1501BK5	R1601B_K5	Recode 2014

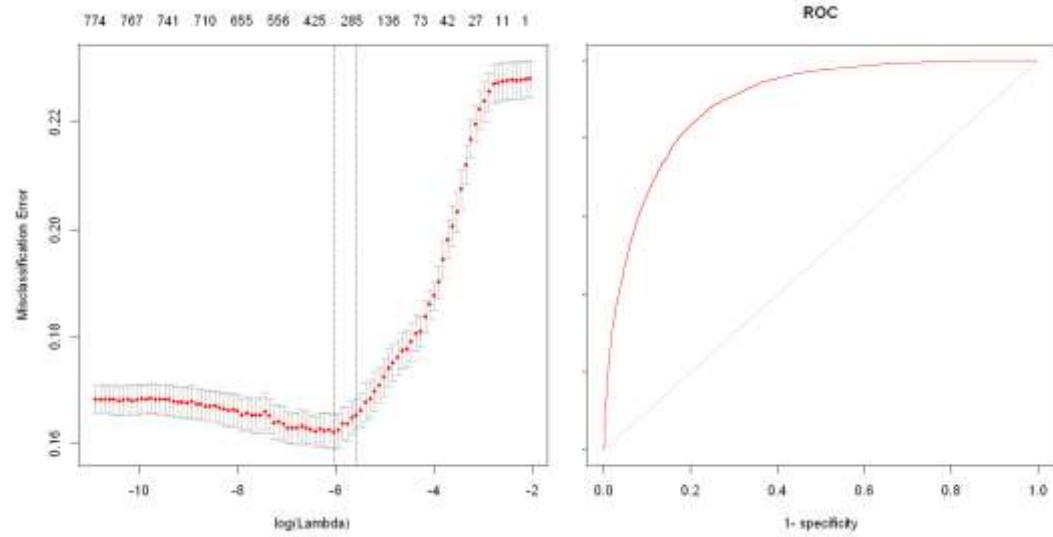
All common variables matched except for minor recoding.

Part V. Modelling and Prediction

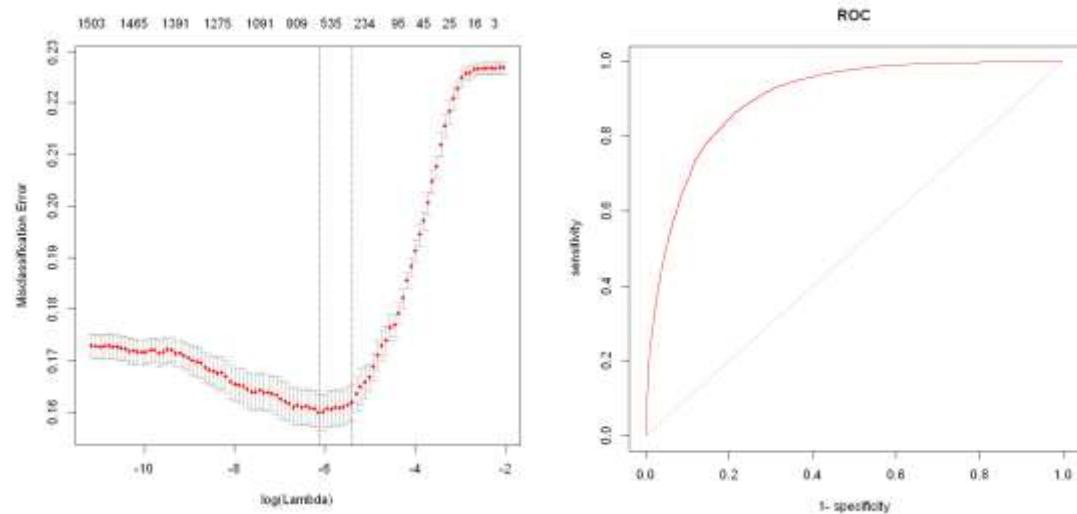
[Modelling]

For a simple and probability-based model, we use Logistic regression with L1 penalty. This is primarily because of its simplicity and the ability to choose features on its own. Given the lack of domain knowledge to build the model from the bottom up, this characteristic is extremely helpful. To control for overfitting, we use cross-validated misclassification error to find the best parameter.

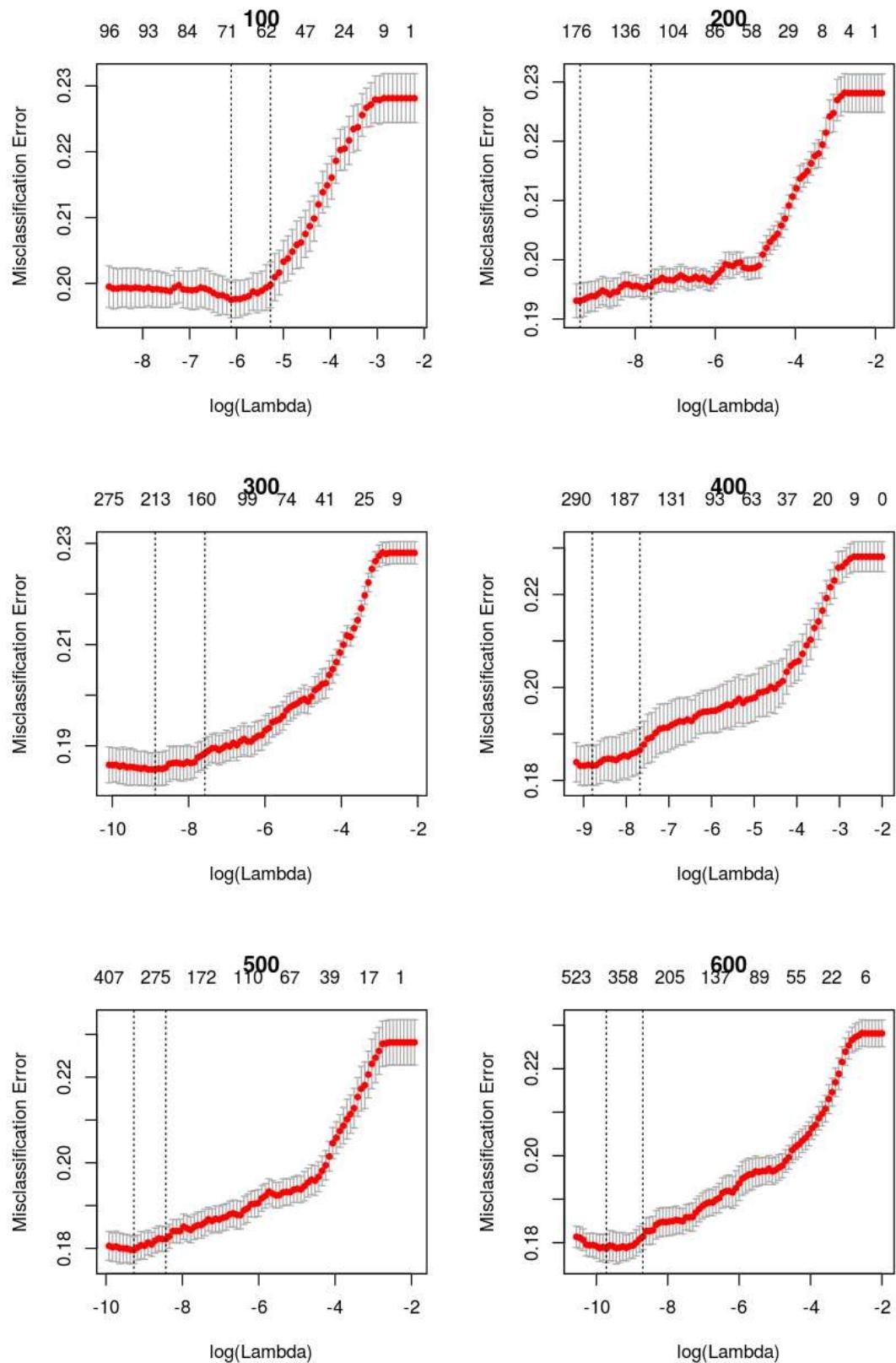
Using a global schema on the training set, the cross-validated misclassification error is 16.2%. Considering there are 3307 “**connect**” villages in the total of 14,496, a blind guess would result in 22.8% error rate. Therefore, the model improves the accuracy by 6.4%. With a one-standard-deviation lambda, 286 variables are reported to have non-zero coefficient. The AUC score is 0.896.



As a comparison, we train the model on the original PODES dataset, which contains 540 useful variables. The misclassification rate slightly decreases to 16.0% with a slight higher AUC = 0.902. Using one-standard-deviation lambda, 349 variables are reported non-zero, which reflects a more complex model.



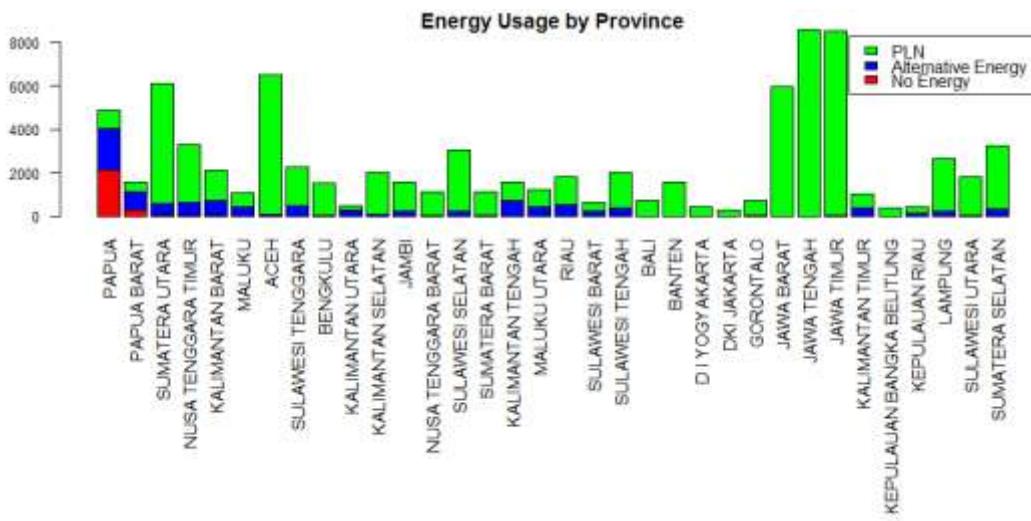
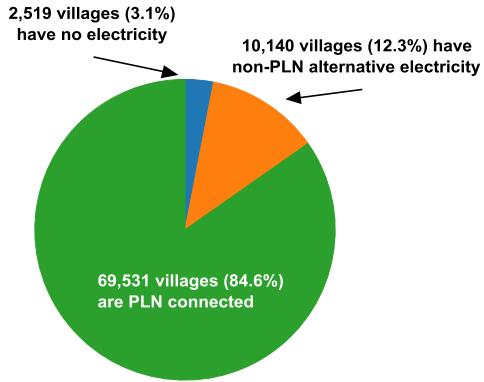
A neural network auto-encoder for self-feature generation is also tested to see if self-generating features (through an auto-encoder) can serve as better predictors than original ones in logistic regression. To save computing power, the network only contains one layer with 50 times of iterations. The following graphs show the cross-validated error rate with different number of neurons in each test. The maximum number is 600, which is less than the original model matrix generated from the dataset.



It is noted that though following a downward trend, the misclassification rate from the auto-generated features perform no better than the original predictors. Therefore we are comfortable with the original model.

[Prediction]

An overlook of 2014 PODES survey shows around 85% of all national villages are already connected to PLN network. Among 12,659 villages that are not PLN-connected, 2,519 (3.1%) have no electricity at all. We also see there is a wide unbalance between provinces in the number of disconnected villages, with New Guinea having the largest number of villages with no electricity of any kind.



Using the model trained with global schema, we calculate the predicted probability of village connect to PLN on the PODES 2014 dataset and plot them by main islands of Indonesia.

Part VI. Discussion and Next Steps

The maps below show the predicted probability of connecting to PLN by 2017 (a prediction made based on training data in 2011-2014). The darker blue area shows the region that are already PLN-connected. A red-blue color threshold indicates the value of probability that are suitable for each island. All provinces except Java (color threshold = 0.5) and New Guinea (color threshold = 0.05) has the value of 0.1, indicating the red area are less likely to be connected (probability less than 10%).

A preliminary reading of the maps indicates a wide provincial disparity of power connections. In

Java, only small pockets of area are shown in red, showing a near 50% predicted probability of future connection. On the other hand, there are vast areas in Kalimantan and New Guinea that have low probability of future connection. In particular, Papua has the largest coverage of areas that are extremely unlikely to be connected (less than 5%).

It's also interesting to note that most villages with high predicted probability (light blue in color) are closer to the villages that are already PLN-established, showing geo-location, such as province, district and sub-district play an important role in the predictive model.

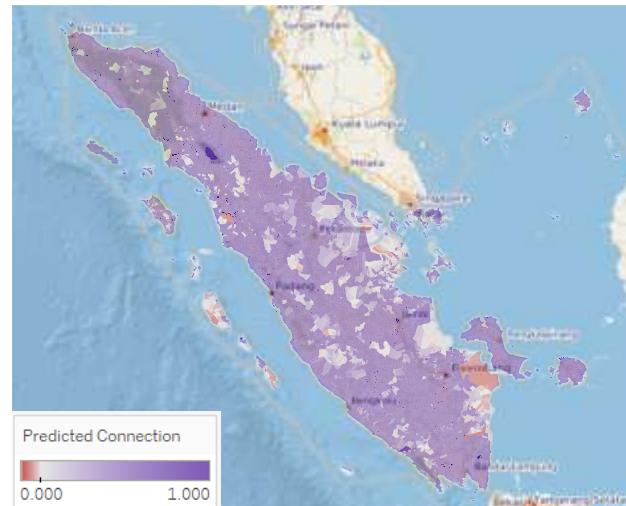
We may also observe that current PLN network seems to expand along major highways. Areas that are further away from these roads tend to be less likely to connect in the future.

Finally, smaller and separated islands that are away from the main island seem less likely to be connected in the future. It is possible due to their remote location from established network and poorer infrastructure.

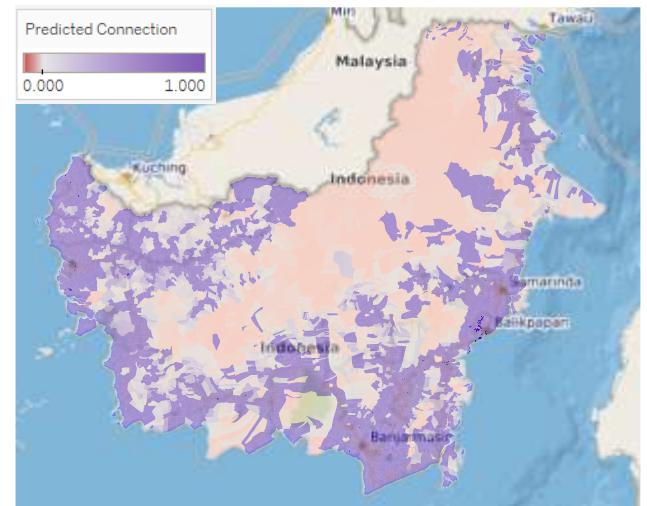
JAVA(50%)



SUMATERA



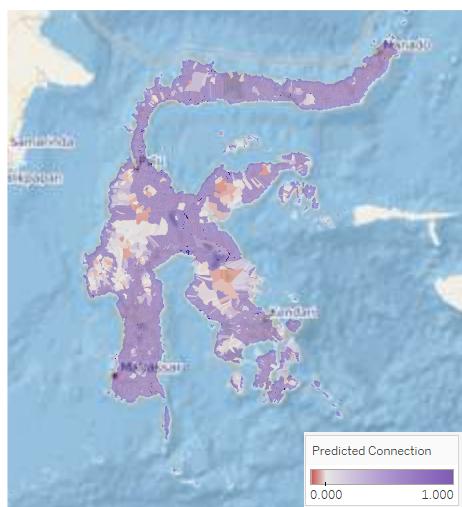
KALIMANTAN



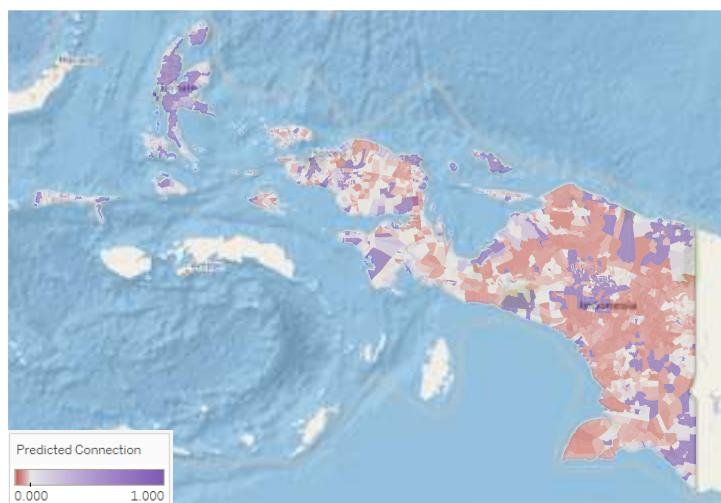
LESSER SUNDA ISLANDS



SULAWESI



NEW GUINEA (5%)



[Next Steps]

It is noteworthy that a few more steps could be taken to refine the model and prediction.

Firstly, we realize that the prediction for New Guinea could be more optimistic than reality. This is primarily because of the treatment of unique sub-districts in PODES 2014 that are not available for PODES 2011. If we can recode them to known sub-district IDs (like North Kalimantan), then we should have improved prediction for the province.

Secondly, due to time restraints and lack of domain knowledge, we have not conducted data imputation in great detail. This results in many sparse variables with large amount of unique level coded as "NA" or 0 (even they are truly a unique level). While this has limited effect on the predicted probability, they do report large coefficients and create difficulty for model interpretation. In order to gain more precise and interpretable insights from the coefficients, it would be a good idea to consider collapsing them into single variables.

Thirdly, even though neural network with auto-encoder did not provide stronger training accuracy than the original model, in the interest of further improve prediction, it is still worthwhile to run

further tests particularly with larger number of iterations. Due to computing limitation, only 50 iterations were used during my tests (maximum would be the number of training observations). It is important to allow for more parameter updates in order to verify if the technique could provide any improvements as a benchmark.

Finally, the global schema could be expanded to meet future survey templates (i.e., 2017) by not deleting unused variables in the training phase. We had deleted these on purpose for faster calculation and a cleaner schema.