

Initial report on analysis of the Anganwadi 2010 dataset

Avishek Sen Gupta
ThoughtWorks

October 31, 2011

Contents

| | | |
|----------|---|-----------|
| 1 | Abstract | 4 |
| 2 | Methodology | 5 |
| 2.1 | CRISP-DM | 5 |
| 2.2 | Relevance of CRISP-DM to this report | 5 |
| 3 | Data Preparation | 6 |
| 3.1 | Nature of the source data | 7 |
| 3.2 | Data representation | 8 |
| 3.2.1 | Data store | 8 |
| 3.2.2 | Schema | 9 |
| 3.3 | Data migration: Identifying invalid data | 11 |
| 4 | Bias | 12 |
| 4.1 | Sampling bias | 13 |
| 5 | Shape of the Data | 14 |
| 5.1 | Univariate distributions | 15 |
| 5.1.1 | Observations/Notes | 15 |
| 5.2 | Outlier analysis | 18 |
| 5.3 | Bivariate distribution | 19 |
| 5.3.1 | Observations/Notes | 19 |
| 5.4 | Summary plots | 21 |
| 5.4.1 | Observations/Notes | 21 |
| 5.5 | Rank Order Charts: Geoclusters | 22 |
| 5.5.1 | Observations/Notes | 23 |
| 6 | Data Analysis and Exploration | 24 |
| 6.1 | Parallel Coordinates | 25 |
| 6.2 | Covariance plot | 26 |
| 6.2.1 | Observations/Notes | 28 |
| 6.3 | Geographical distribution | 29 |
| 7 | Models and Statistical Tests | 30 |
| 7.1 | Tests for conformance to distributions | 31 |
| 7.1.1 | Jarque-Bera test | 32 |
| 7.1.2 | Quantile-Quantile plots | 32 |
| 7.2 | Answer distribution | 33 |
| 7.3 | Effectiveness of intervention | 37 |
| 7.3.1 | Intervention effect on individual responses: McNemar's Test | 37 |
| 7.4 | Modeling responses as Bernoulli trials | 38 |
| 7.4.1 | The Binomial Distribution | 39 |
| 7.4.2 | Typical questions | 40 |
| 7.5 | Test for variable independence | 40 |
| 7.5.1 | Chi-square test | 41 |

| | | |
|-----------|--|-----------|
| 8 | Prediction and Classification | 42 |
| 8.1 | Decision Trees | 43 |
| 8.2 | Bayes classifier | 44 |
| 8.3 | Density estimators | 45 |
| 8.3.1 | Naive Bayes density | 45 |
| 8.3.2 | Kernel density estimation | 45 |
| 8.3.3 | Results | 45 |
| 9 | Dimension reduction/Factor analysis | 46 |
| 9.1 | Principal Component Analysis | 48 |
| 10 | Technical notes | 49 |

List of Figures

| | | |
|----|---|----|
| 1 | Probability distributions of pre-, post-intervention, and improvement | 16 |
| 2 | Bivariate probability distribution of pre- vs. post-intervention scores | 20 |
| 3 | Box plots of pre- and post-intervention scores, broken down by language | 21 |
| 4 | Parallel Coordinates showing language, gender, school, pre- and post-intervention scores | 25 |
| 5 | Covariance plot of pre-intervention responses | 26 |
| 6 | Covariance plot of post-intervention responses | 27 |
| 7 | Geocoded populations by cluster | 29 |
| 8 | Quantile-Quantile plot of score improvement vs. theoretical Gaussian | 32 |
| 9 | Curve-fitted theoretical exponential for reflected post-intervention probability distribution | 34 |
| 10 | Quantile-Quantile plot of post-intervention score (reflected) vs. theoretical exponential | 35 |
| 11 | Answer distribution vs. question number | 36 |
| 12 | Bayes posterior distribution of language from score improvement | 45 |
| 13 | Bayes posterior distribution of gender from score improvement | 46 |
| 14 | Bayes posterior distribution of geocluster from score improvement | 47 |

1 Abstract

This report summarises the results of exploration of the Anganwadi dataset provided by the Akshara Foundation. The analysis aims to characterise the structure of the data, and reveal trends (which would otherwise be obscured by the format of the source data) which may inform strategy through subsequent prediction and/or classification procedures.

2 Methodology

2.1 CRISP-DM

CRISP-DM is a process model distilled from the most common approaches used in data mining procedures. It stands for Cross Industry Standard Process for Data Mining. Not so much a prescription as a collection of 'good practices' followed by data mining professionals, **CRISP-DM** has the following characteristics.

- Domain-neutral
- Tool-neutral
- Provides a structural approach to the data mining process

CRISP-DM segregates data mining endeavours into the following phases.

- Business Understanding
- Data Understanding
- Data Preparation
- Modeling
- Evaluation
- Deployment

2.2 Relevance of CRISP-DM to this report

As far as this report is concerned, the relevant or most significant phases we focus on are:

- Data Understanding
- Data Preparation
- Modeling

Work on the Evaluation step is still preliminary, and will probably be the subject of another report. In a full-fledged project, the rest of the activities upstream and downstream to the above list will assume more importance, and require corresponding investment.

3 Data Preparation

3.1 Nature of the source data

The dataset comes from the education domain. The source data is a file, with each line corresponding to a single student evaluation record. Roughly, there are 29000 records, prior to any data sanitisation. Each line is pipe(|)-delimited into multiple fields. The fields salient to this analysis are listed below:

- Location of the student's school
- Language of the student
- Student's score before intervention
- Student's score after intervention

The score is not a single number, it is a set of 56 responses marked as 0/1. Generally, a 1 may be treated as a favourable answer, therefore, adding them up to get a single aggregate score has natural ordering: a sense of who did better. We reproduce two such records below, with the original formatting.

[illegible]

Looking at the second row, we see that the location of the Anganwadi is BADAMAKAAN I, the student is female and speaks Urdu. The first contiguous set of 0s and 1s is the pre-intervention score, and the next set is the post-intervention one.

3.2 Data representation

3.2.1 Data store

Before any sort of sanitisation or analysis may be performed, it is important to ensure that the source data is stored in a format/datastore which makes querying and modifying the data relatively painless. This decision is largely driven by technological considerations, like:

- Scale of data (centralised/distributed store?)
- Sophistication of queries (OLAP/OLTP?)
- Structure of data, or lack thereof (SQL/NoSQL?)

We were dealing with only about 29000 records, and most of the analysis would probably be performed outside the database. Thus, we opted to use MySQL as our datastore.

3.2.2 Schema

The decisions when creating the database schema affect the ease of querying for relevant information. Apart from the attributes of interest, we wanted to store the individual binary responses as well. One way is to create one column for each response, giving us a total of 112 columns for storing these responses (56 for pre-intervention, 56 for post-intervention). The other way, and that is the one that we chose was to store this information as a 64-bit integer (bigint for MySQL). When required, we could unpack the individual response bits from this number.

A `desc responses;` command on the table reveals the schema we ended up with.

| Field | Type | Null | Key | Default | Extra |
|------------------|------------|------|-----|---------|----------------|
| student_id | int(11) | YES | | NULL | |
| area | char(50) | YES | | NULL | |
| pre_performance | bigint(20) | YES | | NULL | |
| post_performance | bigint(20) | YES | | NULL | |
| language | char(50) | YES | | NULL | |
| gender | char(20) | YES | | NULL | |
| pre_total | int(11) | YES | | NULL | |
| post_total | int(11) | YES | | NULL | |
| id | int(11) | NO | PRI | NULL | auto_increment |
| school_id | int(11) | YES | | NULL | |
| year | int(11) | YES | | NULL | |

The most important use of the reference data is to locate the schools geographically. Given that geocoding the school from its name, we used the cluster to locate schools in 2D space. We deal with geographical analysis in a later section. The schema of the master data mostly mirrors the CSV master data file format, with the addition of latitude and longitude, like so:

| Field | Type | Null | Key | Default | Extra |
|-------------|----------------|------|-----|---------|----------------|
| district | char(50) | YES | | NULL | |
| block | char(50) | YES | | NULL | |
| cluster | char(50) | YES | | NULL | |
| school_id | int(11) | YES | | NULL | |
| school_code | char(20) | YES | | NULL | |
| school_name | char(50) | YES | | NULL | |
| id | int(11) | NO | PRI | NULL | auto_increment |
| latitude | decimal(20,10) | YES | | NULL | |
| longitude | decimal(20,10) | YES | | NULL | |

+-----+-----+-----+-----+-----+-----+

To identify latitude and longitude, we used Google's Map API to geocode the cluster information. It is to be noted that there may be some clusters which weren't located by the Map API, and some more work is needed to cross-validate the coordinate information taken from the Map API.

3.3 Data migration: Identifying invalid data

It is natural to expect missing or corrupted data. The most crucial attribute are the score data, as any misinterpretation of that data may adversely bias the quality of our analysis. Thus, specific checks were put in place to ensure that none of the binary responses was null or some string other than 0 or 1.

Using this check, we found 1067 responses which violated it. All of them had either empty pre- or post-intervention scores. We did not migrate these response records, though it may be possible to do Monte Carlo simulations to predict the missing data.

As a result, out of a total of 28535 records in the original source, 27468 were migrated to the database.

We also found a large fraction of records which did not have a LANGUAGE attribute, i.e., that field was empty. Nevertheless, they were included in the migration.

4 Bias

Analysis is most susceptible to bias in the data collection stage. Sampling is one such activity. If, for a statistical study, participating individuals are not equally likely to have been selected, it may be difficult to distinguish between the actual phenomenon and this biased sampling. This sort of bias is called sampling bias.

4.1 Sampling bias

To find evidence for bias, we looked at a few parameters. Here is the breakdown of the population by language, with the biggest language bucket highlighted.

Unspecified=869
URDU=3564
KANNADA=18685
TELUGU=1688
TAMIL=2051
MARATHI=91
OTHER=239
HINDI=243
KONKANI=18
GUJARATHI=12
NOT KNOWN=3
ORIYA=2
MULTI LNG=1
BENGALI=1
NEPALI=1

There is an overwhelming proportion of students who speak Kannada as their mother tongue (leading by an order of magnitude), a fact that is very likely to bias any sort of analysis where language is involved. We must remain cognizant of such biases, and interpret the results accordingly.

Here is the breakdown of the population by gender.

Girl=14822
Boy=12646

There is not a huge disparity between the two sexes, which indicates that any analysis/prediction based on gender may be less biased.

5 Shape of the Data

Before embarking on any deep analysis of data, it behooves us to look at the shape of the raw data. There are a few reasons why we want to do this.

- **Evident trends/outliers:** Visualisation of the raw data set is always a quick way to spot trends without doing too much analysis. Of course, visualisation is best suited for 1,2 and 3-dimensional data: data of higher dimensionality must usually be either sliced prior to visualisation, or have its dimensionality reduced, before projecting it onto the 2D plane.

Having said that, there are other ways of visualising the data without sacrificing any dimensions at all, such as Parallel Coordinates, though it is more suited to data exploration.

- **Evidence of conformance to well-known distributions:** There exist many probability distributions, some of whose properties are well-studied and well-known, like the Normal distribution. If the data approximates one of these distributions, there are several mature statistical methods which may be applied to test different hypotheses and properties of the data.

Indeed, many of the classical statistical analyses make the assumption that the underlying data is (approximately) normally distributed.

In the following sections, we shall explore the shape of the Anganwadi 2010 dataset, and record our observations on it.

5.1 Univariate distributions

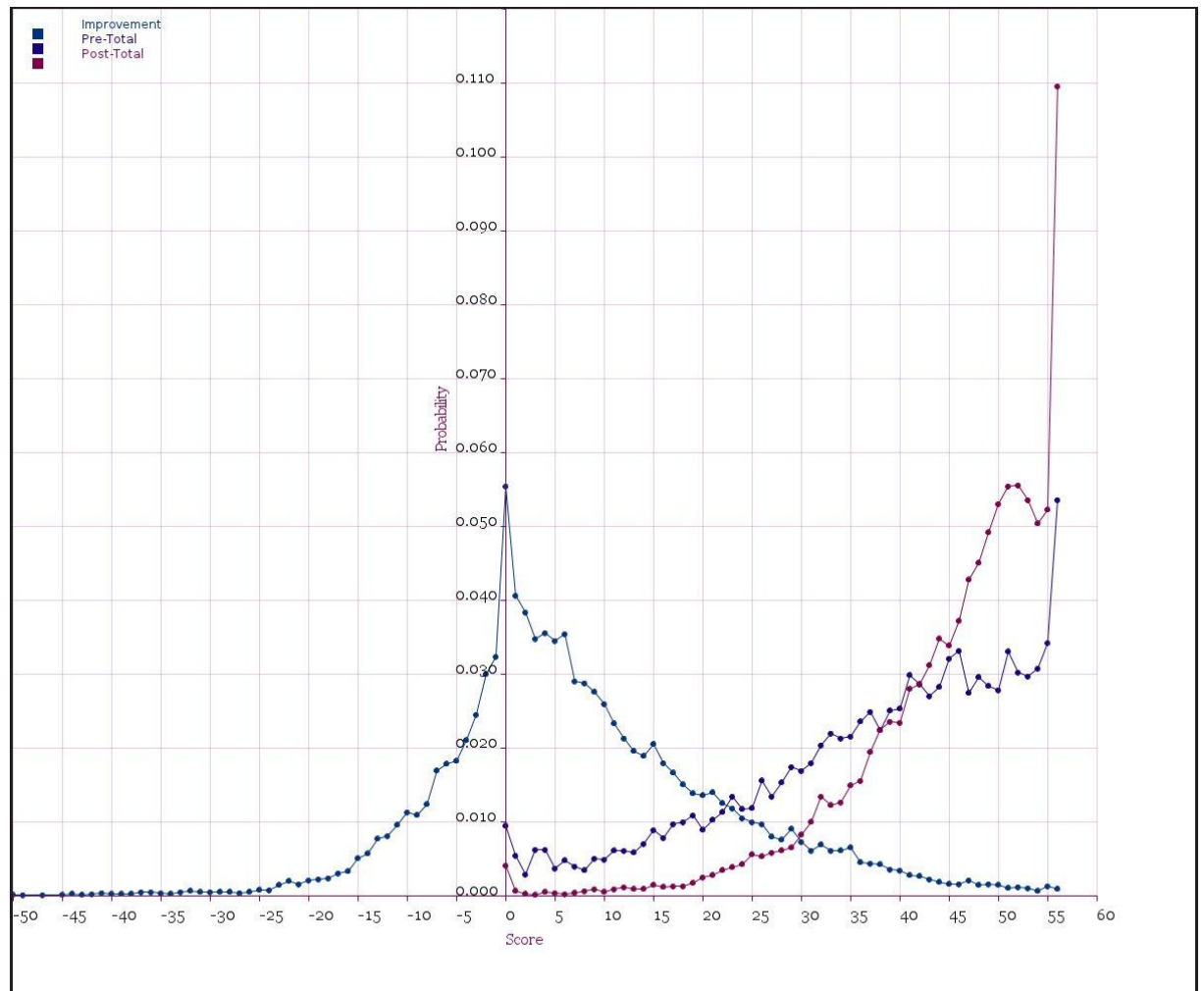
Figure 1 shows 3 distributions: the pre-intervention scores, the post-intervention scores, and the improvements.

5.1.1 Observations/Notes

- The pre-intervention score with the highest number of students is 56, which is the highest score possible. This implies that, even prior to intervention, a sizeable fraction of the students have scored very high on the test.
- The post-intervention score also follows the same trend, albeit with a steeper curve, which implies that many students have scored better in the post-test than in the pre-test.
- The improvement distribution is peaked, the peak being near zero. This makes sense, because if a large fraction of students answered all 56 questions as 1 in the pre-intervention score, there really is no way for them to improve. This is assuming that their performance did not worsen in the post-test. In fact, the calculated mean is around 7.
- There is a significant fraction of students whose performance has worsened in the post-test. This number is 7081. Out of those, we noticed that the worsening was dramatic for a small set. We have listed down the records which had a regression of 40 or more, below.

| student_id | area | pre_total | post_total | language |
|------------|-------------------|-----------|------------|----------|
| 1506619 | VABASANDRA | 52 | 0 | KANNADA |
| 1444355 | KITTAGANA COLONY | 56 | 0 | KANNADA |
| 1426910 | PRIYA DARSHINI | 53 | 0 | KANNADA |
| 1445387 | KYALASANAHALLI | 51 | 1 | TELUGU |
| 1445382 | KOTHANUR | 44 | 0 | TAMIL |
| 1445383 | KOTHANUR | 41 | 0 | KANNADA |
| 1442911 | BETTANA PALYA | 55 | 12 | KANNADA |
| 1445160 | KODIGEHALI | 52 | 0 | KANNADA |
| 1457095 | KAVERI NAGARA | 41 | 0 | KANNADA |
| 1457090 | KAVERI NAGARA | 44 | 0 | TELUGU |
| 1457092 | KAVERI NAGARA | 41 | 0 | TELUGU |
| 1507686 | REHMATH NAGAR | 44 | 0 | URDU |
| 1448385 | MALSANDRA | 50 | 0 | KANNADA |
| 1455798 | KANTEERAVA COLONY | 55 | 15 | KANNADA |
| 1444466 | PRIYA DARSHINI | 52 | 0 | KANNADA |
| 1444467 | PRIYA DARSHINI | 44 | 0 | KANNADA |
| 1445337 | RACHENAHALLI | 52 | 0 | KANNADA |

Figure 1: Probability distributions of pre-, post-intervention, and improvement



| | | | | |
|---------|-----------------------|----|----|---------|
| 1425269 | VINYAKNAGAR | 41 | 0 | KANNADA |
| 552534 | KYALASANAHALLI | 52 | 1 | TELUGU |
| 1448976 | KRISHNA SAGARA COLONY | 54 | 14 | KANNADA |
| 1507157 | BELTHURU | 52 | 7 | KANNADA |
| 1444897 | MUNESHWARA NAGAR | 45 | 0 | TAMIL |
| 1444890 | MUNESHWARA NAGAR | 40 | 0 | KANNADA |
| 1542861 | VERABADRA NAGAR 1 | 56 | 1 | KANNADA |
| 1358415 | KOTHANUR | 48 | 8 | |
| 1445437 | THRIVENINAGARA | 40 | 0 | TELUGU |
| 1442907 | BETTANA PALYA | 55 | 8 | KANNADA |
| 1457106 | KAVERI NAGARA | 40 | 0 | TELUGU |
| 1445444 | THRIVENINAGARA | 41 | 0 | TELUGU |
| 1444474 | PRIYA DARSHINI | 54 | 0 | KANNADA |
| 1445159 | KODIGEHALLI | 41 | 0 | KANNADA |
| 511998 | KODIGEHALLI | 55 | 0 | KANNADA |
| 1356274 | KAVERI NAGARA | 56 | 0 | |
| 1358429 | KOTHANUR | 42 | 0 | |
| 1497510 | JALAHALLI 1 | 42 | 0 | KANNADA |
| 1443914 | BIDARAHALLI | 56 | 12 | KANNADA |
| 1366955 | PRIYA DARSHINI | 53 | 0 | KANNADA |
| 1366956 | PRIYA DARSHINI | 53 | 0 | KANNADA |
| 1447019 | MADAPPANA HALLI | 42 | 0 | TELUGU |
| 1355112 | BYRATHI BANDE | 55 | 14 | |
| 1504259 | MAYASANDRA A | 56 | 0 | KANNADA |
| 1457120 | KAVERI NAGARA | 43 | 0 | TAMIL |
| 1457089 | KAVERI NAGARA | 43 | 0 | TAMIL |
| 1445171 | KODIGEHALLI | 44 | 0 | KANNADA |
| 1457203 | SARAIPALYA | 51 | 0 | URDU |
| 1457179 | BYRATHI BANDE | 55 | 14 | KANNADA |
| 1504538 | KEMPEGOWDA NAGAR | 51 | 7 | KANNADA |
| 1444486 | PRIYA DARSHINI | 49 | 0 | KANNADA |
| 1444488 | PRIYA DARSHINI | 42 | 0 | KANNADA |
| 1451831 | AMBED NAGAR | 50 | 9 | TAMIL |

- The post-intervention score distribution seems to follow a power law. We shall consider modeling this attribute further on.

5.2 Outlier analysis

5.3 Bivariate distribution

So far, we've been looking at single variables in isolation. Figure 2 shows a bivariate histogram of pre- vs. post-intervention scores. The lighter a cell, the more the number of records in that 'bucket'.

$$y = 0.224.x + 37.134$$

5.3.1 Observations/Notes

- Many of the scores seem to be clustered near the top right. To further highlight this, we have draw a linear regression line as a rough indicator of a trend (To model the trend in more detail, we could use LOESS). What is somewhat puzzling is that there are not a few students whose performance has dropped after the intervention. This is evident even without doing a linear regression.
- The immediate outliers which are visible are the ones on the extreme left (pre 0, post 56) and at the origin (pre=0, post=0). The latter outlier(s) may be an artifact of corrupted data collection; we cannot say.

Figure 2: Bivariate probability distribution of pre- vs. post-intervention scores

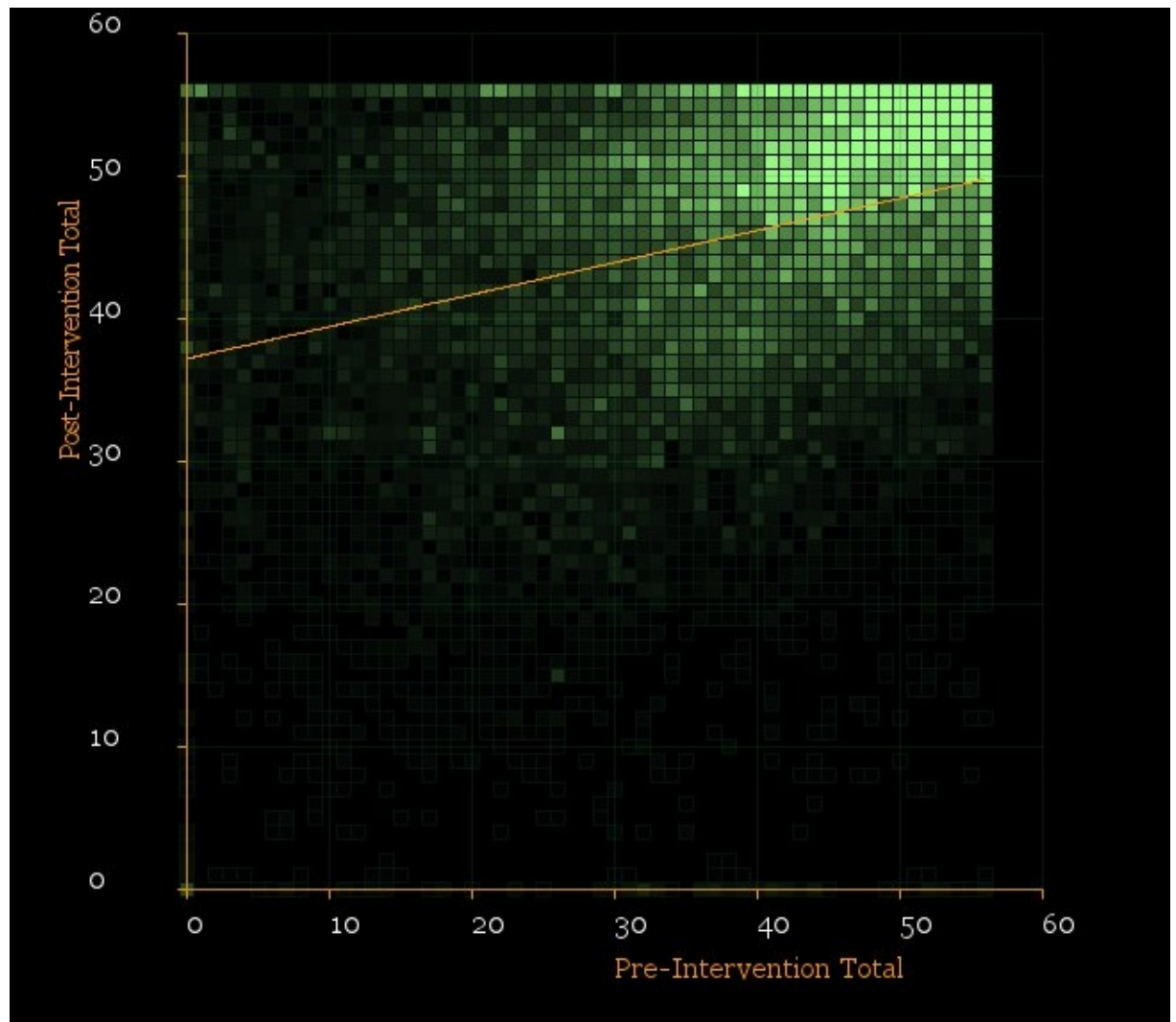
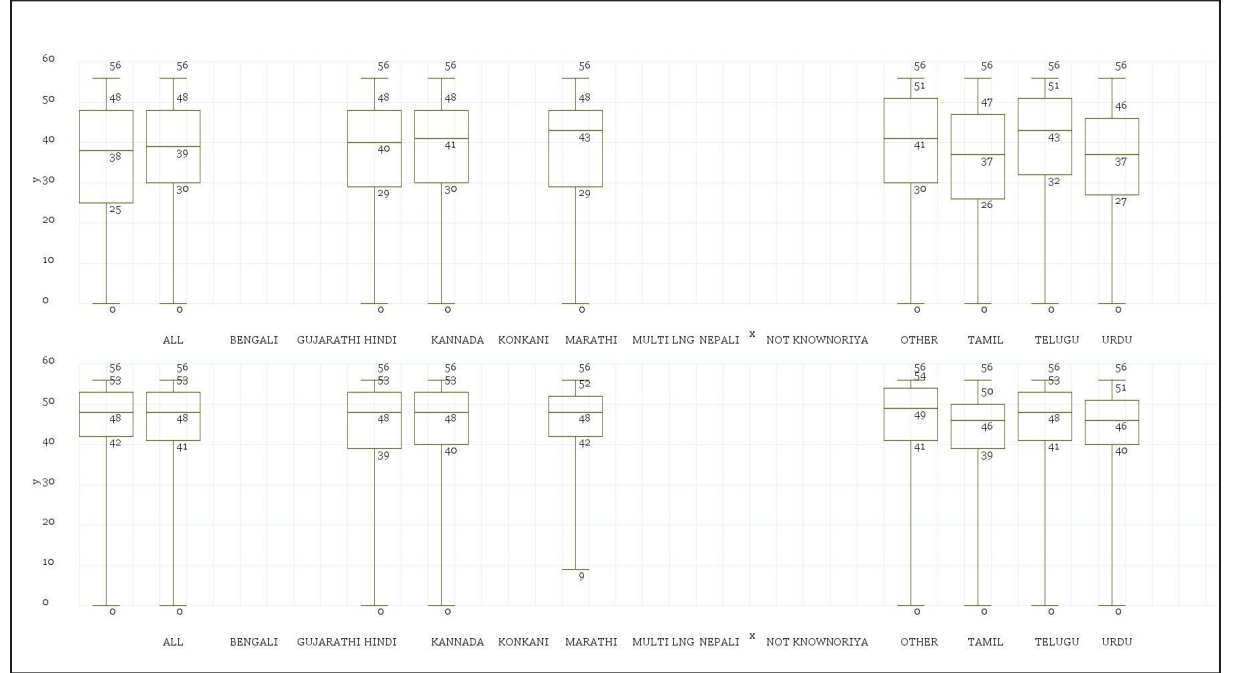


Figure 3: Box plots of pre- and post-intervention scores, broken down by language



5.4 Summary plots

Summary plots are so called for their ability to summarise up a data set as a set of numbers, which can be easily interpreted. We used Box Plots to summarise the data, broken down by language. Figure 3 shows that breakdown. The top row represents the box plots for the pre-intervention assessment, the bottom one for the post-intervention assessment.

5.4.1 Observations/Notes

- Out of all the languages, we were not able to create box plots because the corresponding samples were too few in number to differentiate between the different quartiles. These languages are Bengali, MultiLng, Not Known and Oriya.
- There are no dramatic differences between the plots in each set. The medians, 1st and the 2nd quartiles are rather close to each other.

5.5 Rank Order Charts: Geoclusters

One of the types of data that is pretty common is categorical nominal data. Categorical data is that which is non-numeric. Nominal data is that which does not lend itself to any intrinsic ordering. Names of places (if we ignore their coordinate system representation) are categorical data that is nominal.

The best way to summarise categorical data, especially if it is the Independent Variable, is to order records based on an ordinal Dependent Variable. This sort of a representation is called a Rank Order chart. We have rank ordered the geographical clusters based on the student population in each cluster.

| Index | Cluster | Probability | Cumulative Distribution |
|-------|---------------------------------|----------------------|-------------------------|
| 1. | Sir M Vishweshwariah 1st Circle | 0.030690257754477937 | 0.030690257754477937 |
| 2. | Sultanpalya | 0.030289791757681667 | 0.0609800495121596 |
| 3. | konnana kunte | 0.02668559778651522 | 0.08766564729867482 |
| 4. | kadugodi | 0.02584825979321392 | 0.11351390709188874 |
| 5. | Laggere 4th circle | 0.025156545798747633 | 0.13867045289063637 |
| 6. | K.R.Puram | 0.024756079801951363 | 0.16342653269258772 |
| 7. | Murphy Town | 0.024064365807485073 | 0.1874908985000728 |
| 8. | Netaji 2nd Circle | 0.02340905781272754 | 0.21089995631280034 |
| 9. | Abbigere | 0.02326343381389253 | 0.23416339012669288 |
| 10. | Begure | 0.023190621814475027 | 0.2573540119411679 |
| 11. | Thavare kere | 0.022462501820299987 | 0.2798165137614679 |
| 12. | Kodihalli | 0.022244065822047472 | 0.30206057958351534 |
| 13. | Dodda Banaswadi | 0.02202562982379496 | 0.3240862094073103 |
| 14. | Jigani | 0.021952817824377458 | 0.3460390272316878 |
| 15. | Koramangala | 0.021370321829037427 | 0.3674093490607252 |
| 16. | Byrathi | 0.020387359836901122 | 0.3877967088976263 |
| 17. | Rajanakunte | 0.020059705839522355 | 0.4078564147371487 |
| 18. | Varthur | 0.01944080384447357 | 0.42729721858162223 |
| 19. | Hebbagodi | 0.01907674384738605 | 0.4463739624290083 |
| 20. | Hreohalli | 0.018639871850881024 | 0.4650138342798893 |
| 21. | Yalahanka | 0.01802096985583224 | 0.48303480413572153 |
| 22. | Banashankari | 0.017438473860492208 | 0.5004732779962138 |
| 23. | K Gollahalli | 0.0172928498616572 | 0.517766127857871 |
| 24. | Sarjapura | 0.017183631862530944 | 0.5349497597204019 |
| 25. | Mallasandra | 0.01711081986311344 | 0.5520605795835154 |
| 26. | Kodigenahalli | 0.01689238386486093 | 0.5689529634483763 |
| 27. | kasaba | 0.01667394786660842 | 0.5856269113149848 |
| 28. | Haragadde | 0.016564729867482163 | 0.6021916411824669 |
| 29. | Makali | 0.016455511868355904 | 0.6186471530508229 |
| 30. | Doddakannelli | 0.016419105868647154 | 0.63506625891947 |
| 31. | Mahantalingapura | 0.016200669870394643 | 0.6512669287898646 |
| 32. | Baglur | 0.015836609873307123 | 0.6671035386631717 |
| 33. | chandapura | 0.015836609873307123 | 0.6829401485364789 |

| | | | |
|-----|-----------------------------------|------------------------|--------------------|
| 34. | Bettahalasuru | 0.015690985874472114 | 0.698631134410951 |
| 35. | Hesaragatta | 0.015690985874472114 | 0.7143221202854232 |
| 36. | Kaggalipura | 0.0154725498762196 | 0.7297946701616428 |
| 37. | Sondekoppa | 0.015217707878258336 | 0.7450123780399012 |
| 38. | Wilson Garden | 0.01510848987913208 | 0.7601208679190332 |
| 39. | Dommasanddra | 0.01510848987913208 | 0.7752293577981653 |
| 40. | Kithuru Rani Channamma 3rd Circle | 0.01474442988204456 | 0.7899737876802099 |
| 41. | Triveni 5th Circle | 0.014489587884083296 | 0.8044633755642931 |
| 42. | Palace Guttalli Circle | 0.014453181884374545 | 0.8189165574486676 |
| 43. | Kengeri circle | 0.014343963885248289 | 0.8332605213339159 |
| 44. | Yashwanthpura | 0.014052715887578273 | 0.8473132372214942 |
| 45. | Kuvempunagar 4th Circle | 0.013543031891655745 | 0.8608562691131499 |
| 46. | Chikkajala | 0.012960535896315713 | 0.8738168050094657 |
| 47. | Kadugondanahalli | 0.012596475899228193 | 0.8864132809086939 |
| 48. | Attibele | 0.012523663899810689 | 0.8989369448085045 |
| 49. | Indlavadi | 0.012341633901266929 | 0.9112785787097715 |
| 50. | Sonnenahali | 0.012305227901558177 | 0.9235838066113297 |
| 51. | Uttrahalli | 0.012232415902140673 | 0.9358162225134703 |
| 52. | Machohalli | 0.012159603902723169 | 0.9479758264161935 |
| 53. | Gavipura Guttahalli | 0.009210717926314256 | 0.9571865443425077 |
| 54. | Marasur | 0.009137905926896752 | 0.9663244502694045 |
| 55. | Netravathi | 0.00906509392747925 | 0.9753895441968837 |
| 56. | Gopalapura | 0.008955875928352992 | 0.9843454201252367 |
| 57. | Hennagara | 0.00873743993010048 | 0.9930828600553372 |
| 58. | PANTHARAPALYA | 0.002621231979030144 | 0.9957040920343674 |
| 59. | YARUB NAGAR | 0.001783893985728848 | 0.9974879860200963 |
| 60. | VERABADRA NAGAR | 0.001310615989515072 | 0.9987986020096113 |
| 61. | Yalahanka upnagar | 0.00043687199650502403 | 0.9992354740061163 |
| 62. | DODDA BANASAWADI | 0.000400465996796272 | 0.9996359400029126 |
| 63. | ASHRAYA NAGAR | 0.00036405999708752004 | 1.0 |

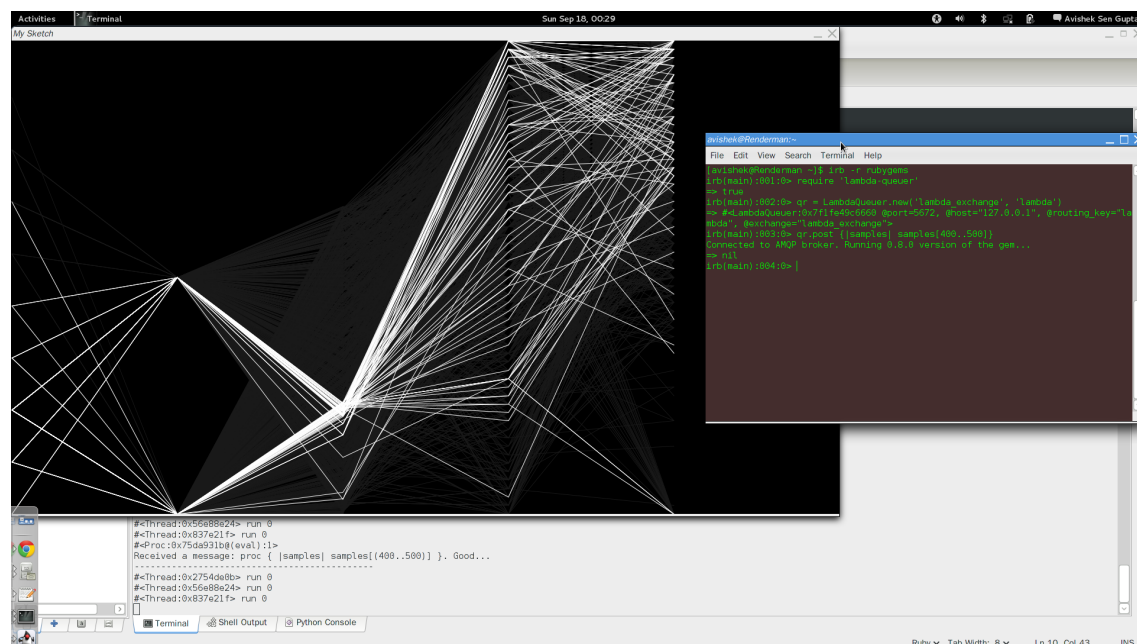
5.5.1 Observations/Notes

- Out of 63 clusters, the top 38 clusters account for 75% of the population.

6 Data Analysis and Exploration

There are multiple ways of exploring datasets when simple visual inspection is tedious and unintuitive. Many of them are standard, but some of them may reveal more about the nature of the data. Often, they are motivated by specific business drivers and questions, and some exploration may be custom to the dataset under analysis. Data exploration is also commonly done through queries to an OLAP database.

Figure 4: Parallel Coordinates showing language, gender, school, pre- and post-intervention scores



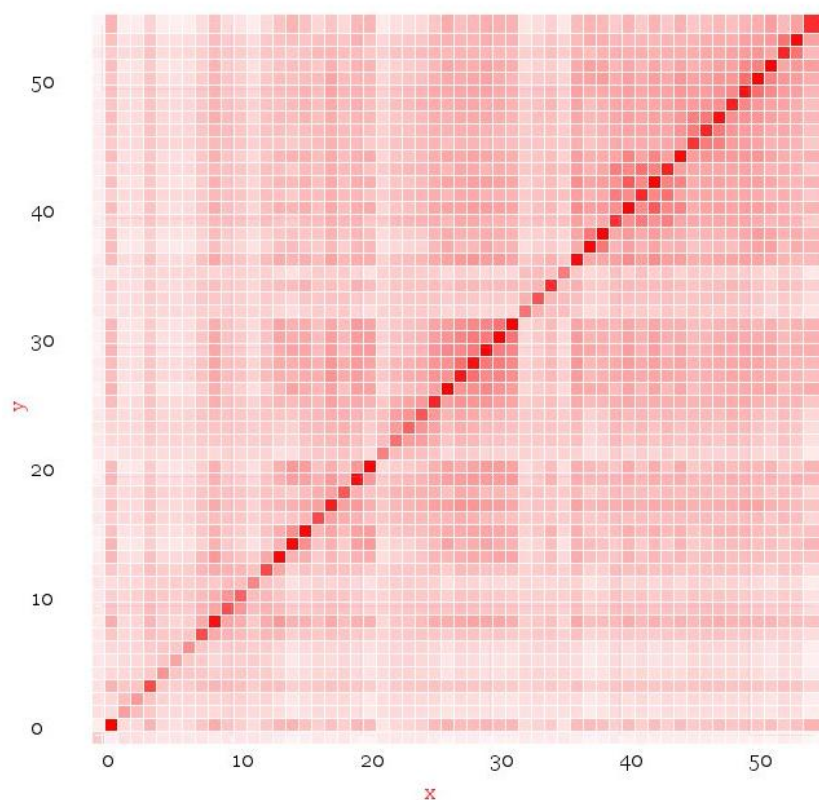
6.1 Parallel Coordinates

Parallel coordinates are an interesting way to explore high-dimensional data without discarding any detail. The only downside is that the elegance of presentation is somewhat sacrificed.

Essentially, instead of aligning axes in orthogonal directions (we can visualise only upto 3), the axes are stacked horizontally, giving the impression of n parallel lines. A single data point in this n -dimensional representation is a set of broken lines spanning the widths between the axes.

The way to explore this visualisation is to highlight the samples of interest, based on some criteria. The highlighted samples can then be inspected visually to discover trends, if any. Figure 4 shows an example of exploration using parallel coordinates on the Anganwadi data.

Figure 5: Covariance plot of pre-intervention responses

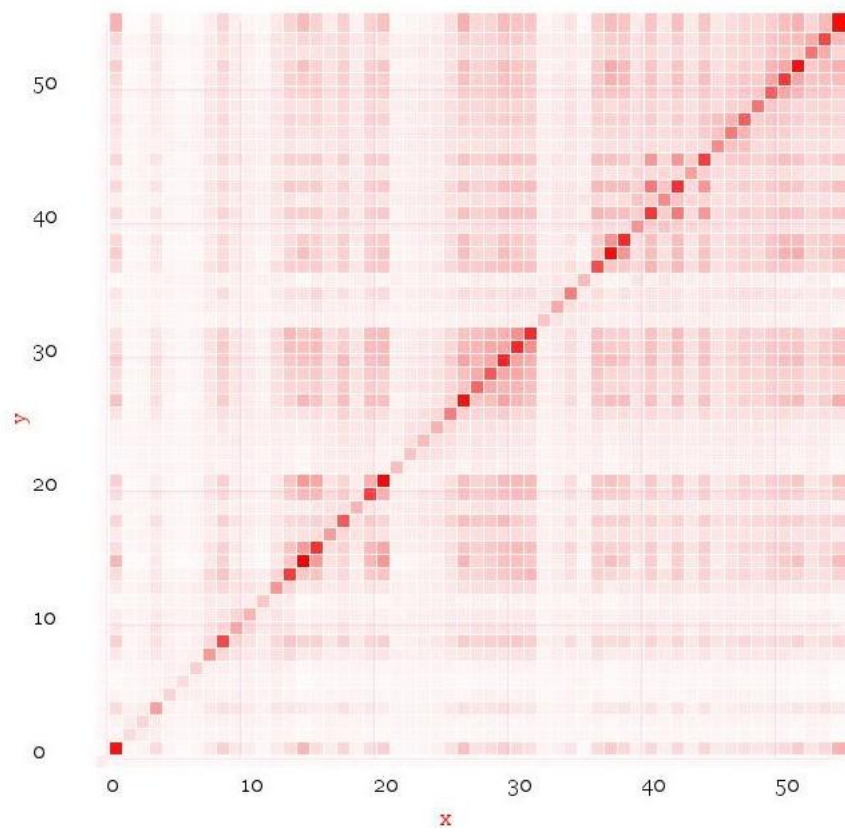


6.2 Covariance plot

The covariance plot is a first step towards looking at the correlations between responses to individual questions. Basically, the kind of question such analysis can answer is, for example, how dependent is a response to question 40 on the response to question 3? Do they change together, i.e., is there a strong dependency between answer to question X and the answer to question Y?

Figure 5 shows the covariance plot for the pre-intervention responses. Figure 6 shows the covariance plot for the post-intervention responses.

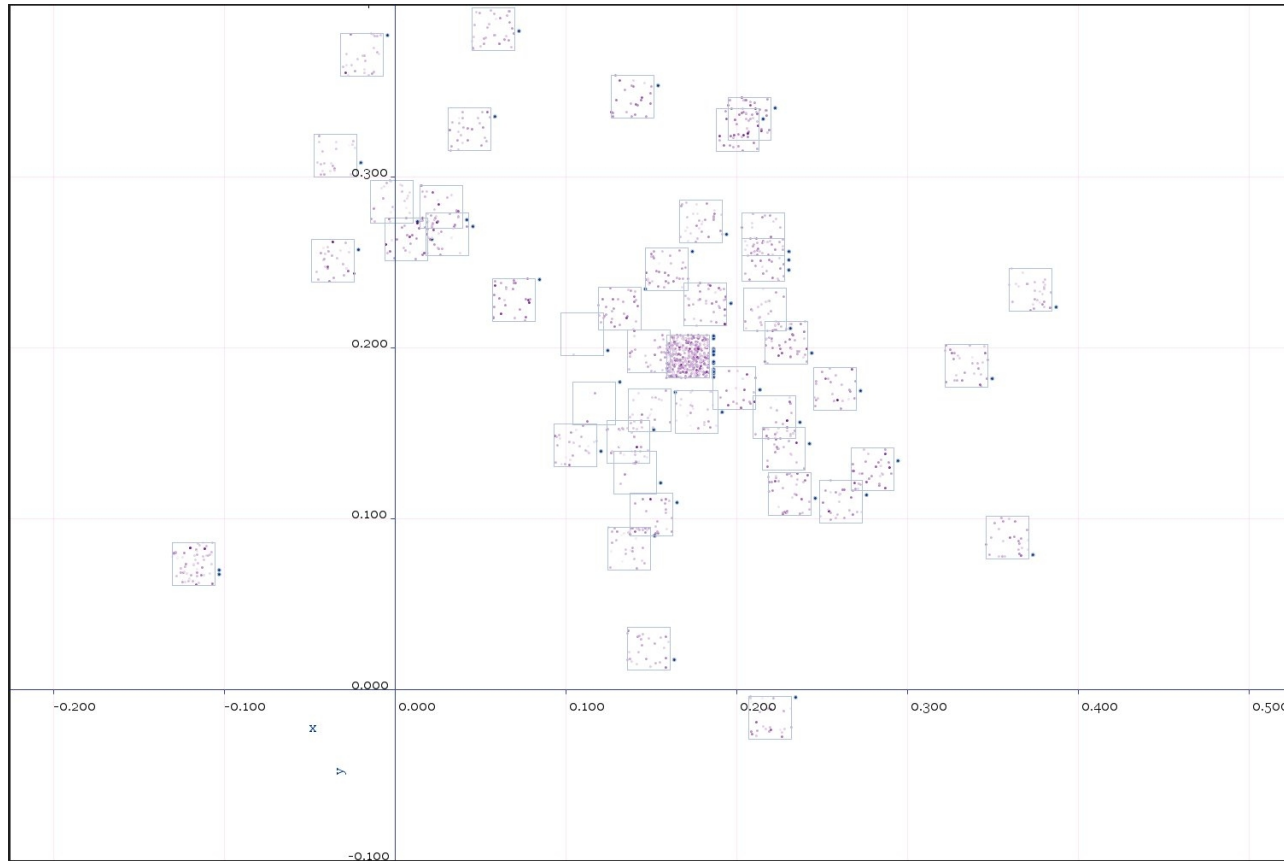
Figure 6: Covariance plot of post-intervention responses



6.2.1 Observations/Notes

- The plots are color coded such that the darker the hue, higher the covariance. The plot is symmetric about the diagonal, and the diagonal is much darker because it is essentially the covariance of a single response with itself. The data has not been normalised yet, hence the diagonal hues vary.
- The covariance between the responses appears weakened in the post-intervention scores, indicating that the degree of dependence between the responses has decreased.
- There appear to be clusters of relatively high covariance, for example, between the questions in the range of (40-43), (26-31), etc. Further investigation may be needed to answer specific questions in this area.
- Correlation does not imply causation. Inferred links between responses may be result of a deeper phenomenon; thus a high covariance doesn't necessarily imply a direct causative link between two responses.

Figure 7: Geocoded populations by cluster

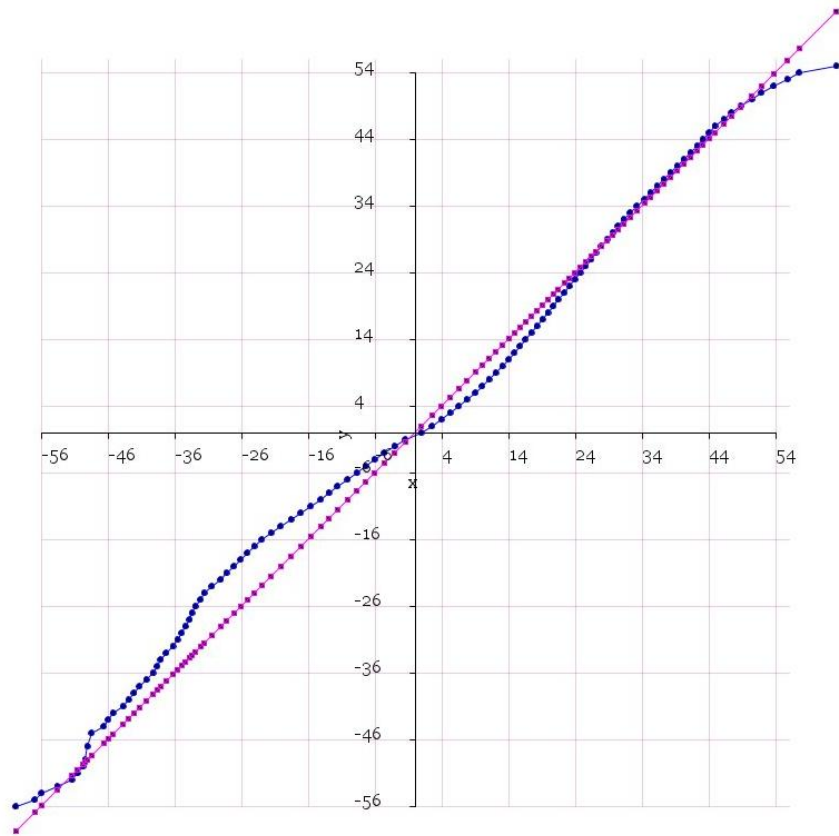


6.3 Geographical distribution

7 Models and Statistical Tests

7.1 Tests for conformance to distributions

Figure 8: Quantile-Quantile plot of score improvement vs. theoretical Gaussian



7.1.1 Jarque-Bera test

7.1.2 Quantile-Quantile plots

7.2 Answer distribution

lolol

Figure 9: Curve-fitted theoretical exponential for reflected post-intervention probability distribution

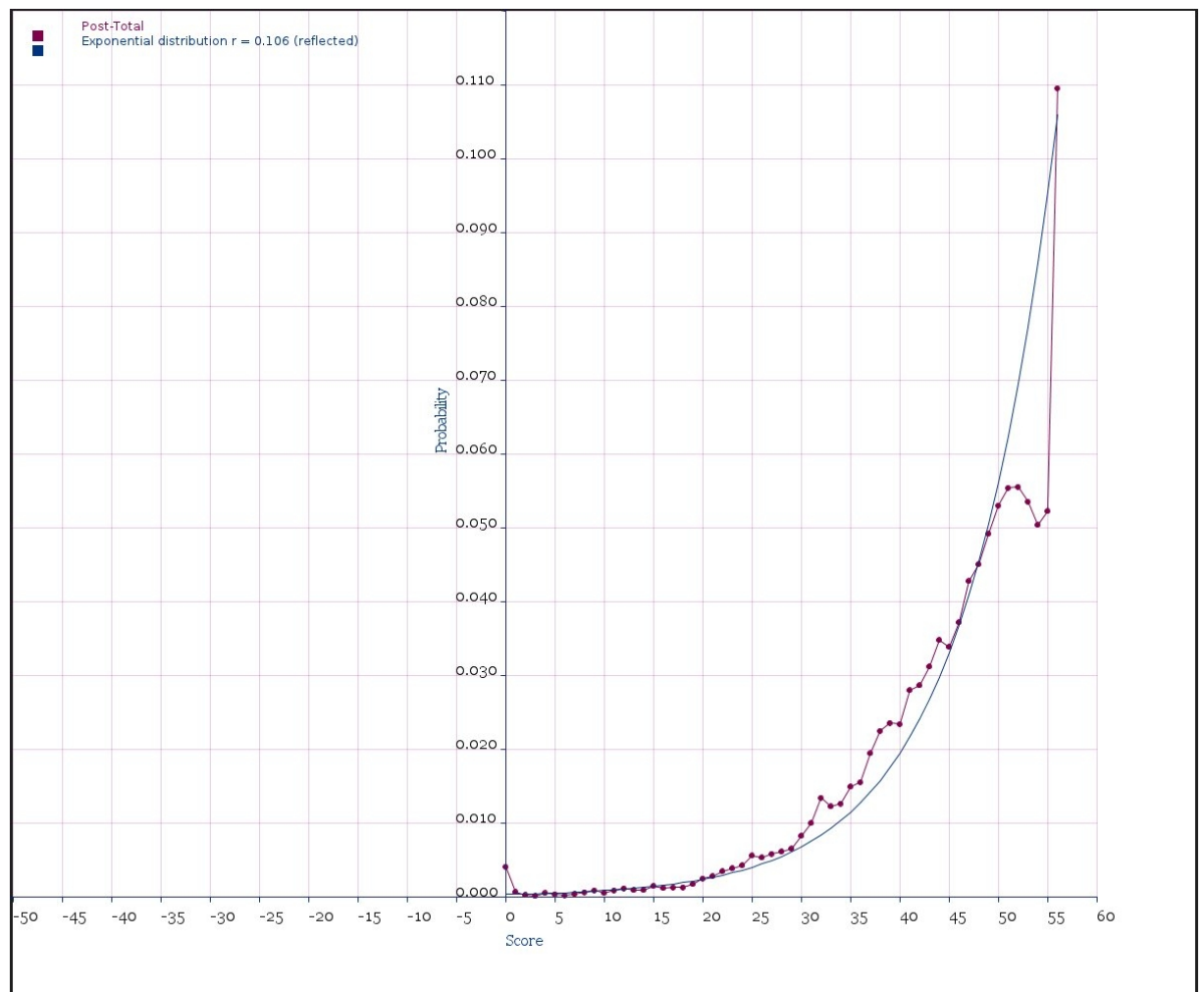


Figure 10: Quantile-Quantile plot of post-intervention score (reflected) vs. theoretical exponential

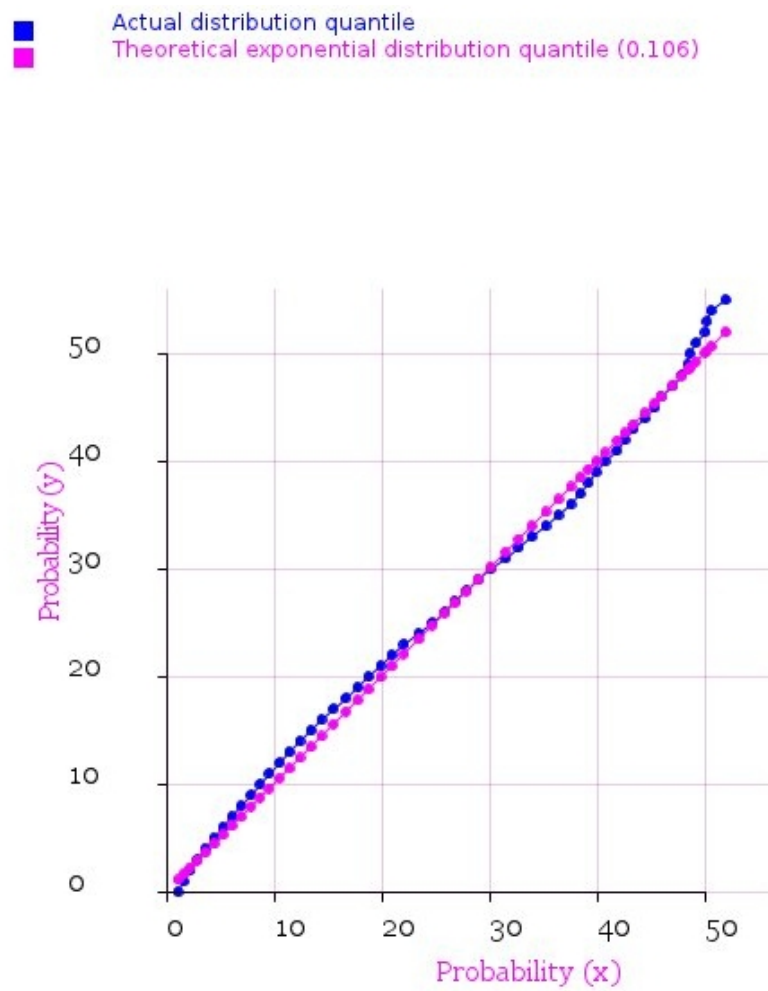
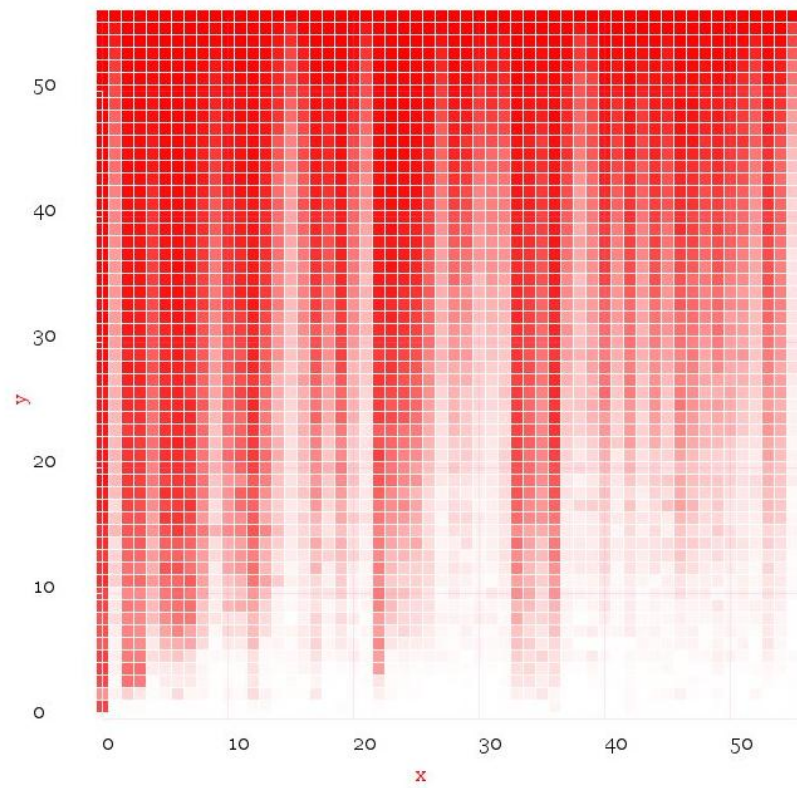


Figure 11: Answer distribution vs. question number



7.3 Effectiveness of intervention

7.3.1 Intervention effect on individual responses: McNemar's Test

7.4 Modeling responses as Bernoulli trials

7.4.1 The Binomial Distribution

7.4.2 Typical questions

7.5 Test for variable independence

7.5.1 Chi-square test

Null hypothesis: Area and Improvement are NOT related.

For area vs. improvement

Chi-Square statistic = 56499.4692602837

$X^2 = 9652.9739$

Degrees of freedom = 9426

Null hypothesis rejected

^c

Null hypothesis: Area and Pre-Score are NOT related.

For area vs. pre-score

Chi-Square statistic = 58665.7089390644

$X^2 = 8062.2959$

Degrees of freedom = 7855

Null hypothesis rejected

Null hypothesis: Area and Post-Score are NOT related.

For area vs. post-score

Chi-Square statistic = 38567.0016158761

$X^2 = 8062.2959$

Degrees of freedom = 7855

Null hypothesis rejected

Null hypothesis: Language and Post-Score are NOT related.

For language vs. post-score

Chi-Square statistic = 280.234448946825

$X^2 = 96.2166$

Degrees of freedom = 75

Null hypothesis rejected

Null hypothesis: Language and Improvement are NOT related.

For language vs. improvement

Chi-Square statistic = 232.464548410971

$X^2 = 113.1452$

Degrees of freedom = 90

Null hypothesis rejected

Null hypothesis: Language and Pre-Score are NOT related.

For language vs. pre-score

Chi-Square statistic = 277.85501653079

$X^2 = 96.2166$

Degrees of freedom = 75

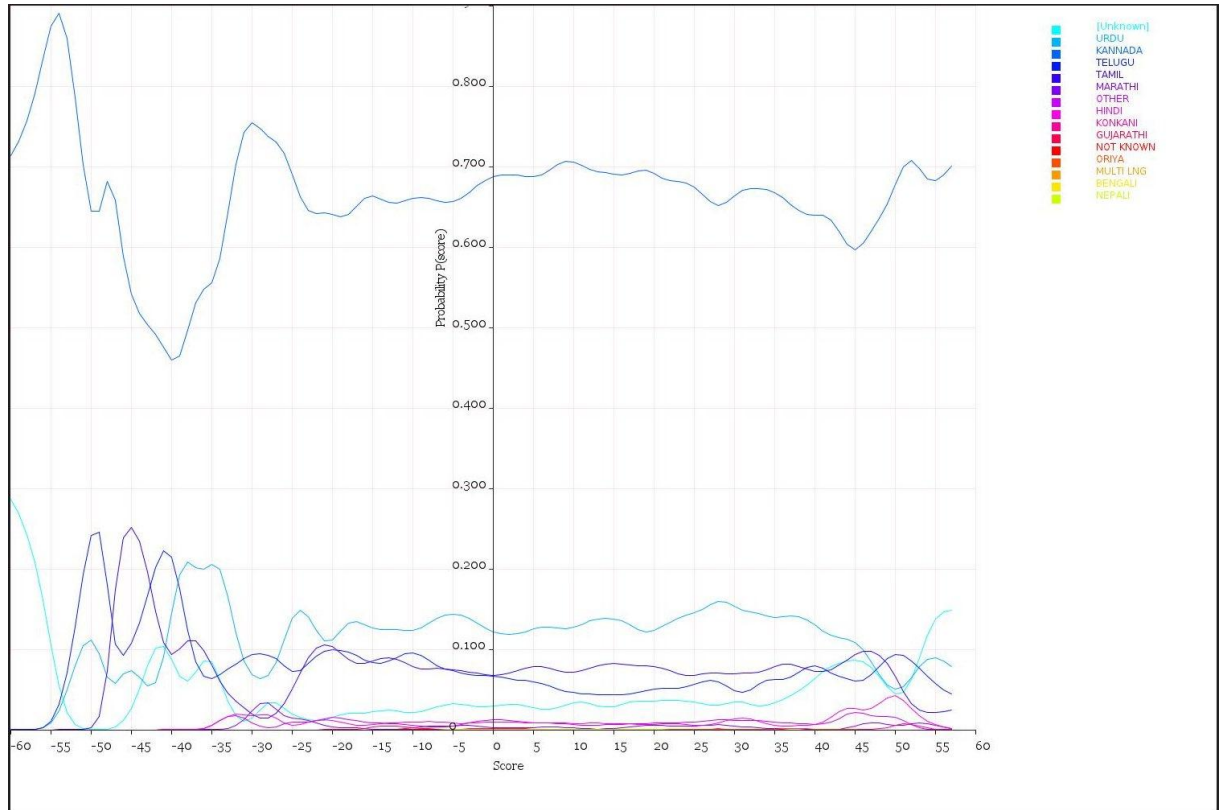
Null hypothesis rejected

8 Prediction and Classification

8.1 Decision Trees

8.2 Bayes classifier

Figure 12: Bayes posterior distribution of language from score improvement



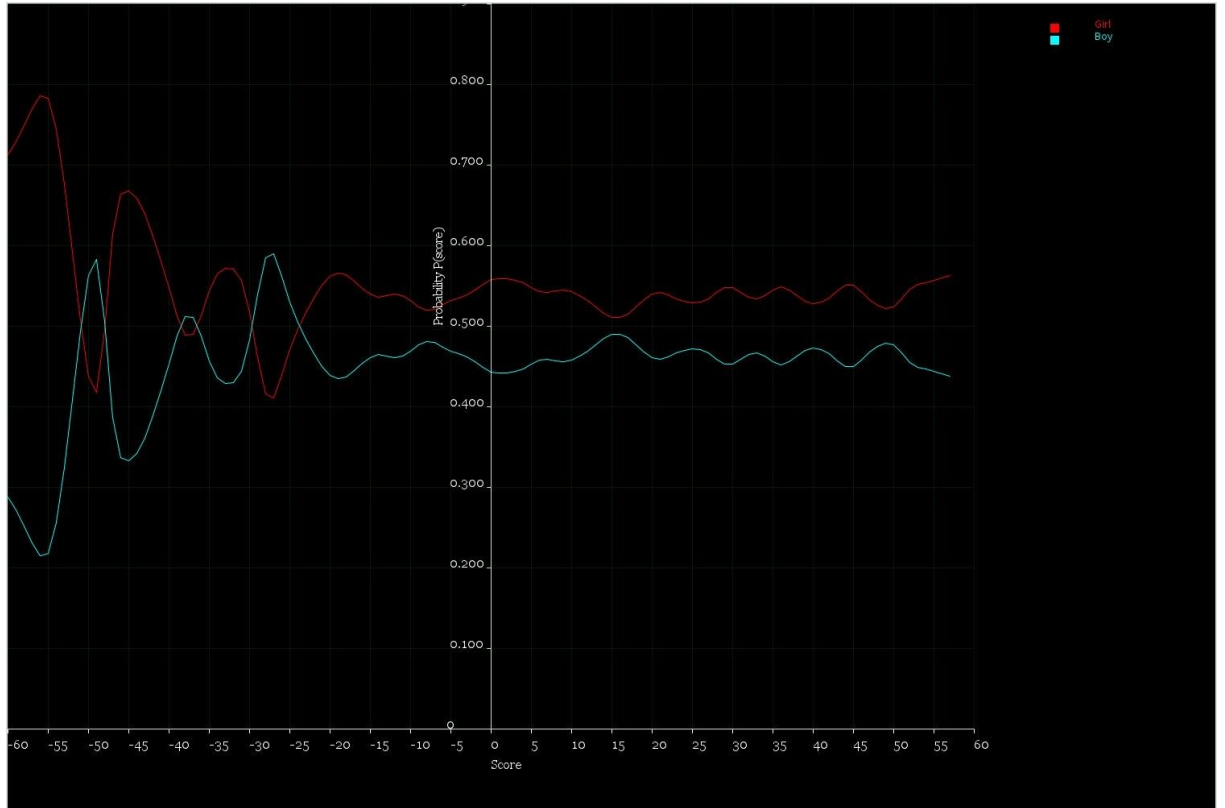
8.3 Density estimators

8.3.1 Naive Bayes density

8.3.2 Kernel density estimation

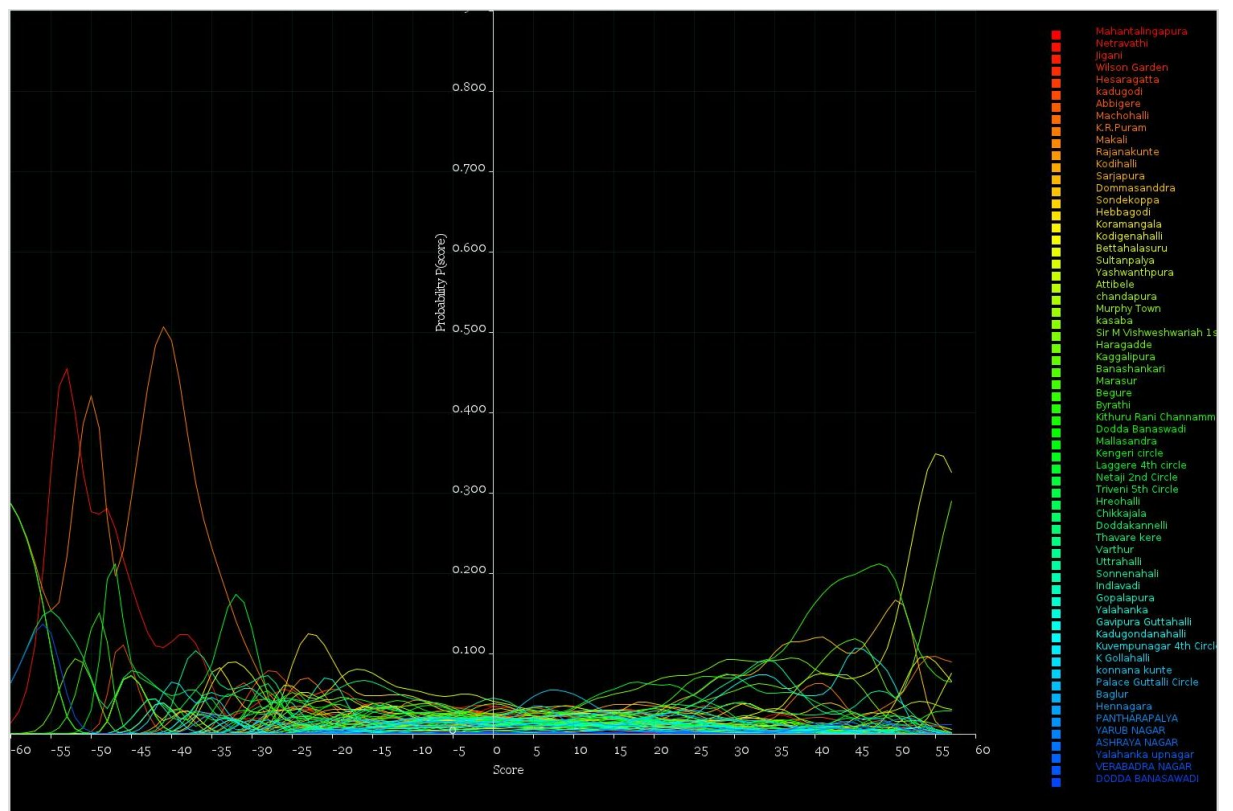
8.3.3 Results

Figure 13: Bayes posterior distribution of gender from score improvement



9 Dimension reduction/Factor analysis

Figure 14: Bayes posterior distribution of geocluster from score improvement



9.1 Principal Component Analysis

10 Technical notes

For base visualisation and interactivity, Processing was used through its Ruby bindings (Ruby-Processing). All plots were done using a coordinate plotting gem called Basis-Processing.

Some calculations like Principal Component Analysis were done using the Stat-sample gem from the SciRuby project.

JRuby 1.6.4 was used for all code needing visualisation; Ruby 1.9.2 was used for everything else.

Data was stored in a MySQL database.