Report Name: INF558 Homework 4 submission

Submitted By: Anand Prashar

Date: 04 October 2017

---

1. What are the fields in your data that require cleaning and what kind of data do these fields contain?

*Ans*: I downloaded some messy data from internet as my project did not have any. This is a tabular data on Africa'a provincial and district hospitals. My subset has 5 fields: Name, Owner, Category, District, Province.

3 fields needed cleaning:

Name (type: string data, some Unicode characters in hospital names, many whitespaces)

Owner (type: string data, different names for same entities in many fields)

Category (type: string data, uneven representation of letter casing, null cells)

---

2. Describe which cleaning operations should be performed on these fields?

*Ans*:

On Name field-

1. Need to remove unnecessary numbering labels .
2. Whitespaces varies in between words too. Need to have at max 1 whitespace in between.
3. Some names had weird characters. Need to map to Unicode to correct their state.

On Owner field-

1. Same words had . or abbreviation causing to show different instances in facets. Need clustering to resolve this.
2. In some cases, entries from Owner category has spread into adjacent column. Need to fix that with flagging and reference correction.

On Category field-

1. Same words are represented in uppercase and titlecases. Need to cluster them, thus reducing facet classes.
2. There are many null cells, need to at least replace with uniform whitespace.

---

Screenshots

Name field (Before and After)

| Name | Name (Cleaned) |
| --- | --- |
| 1. Chegutu | Chegutu |
| 2.Mhondoro | |
| 3. Gora | Gora |
| 4. Mbuyanehanda | Mbuyanehanda |
| 5.Monera RHC | |
| 6. Msengezi | Msengezi |
| 7. Musinami | Musinami |
| 8. Chegutu | Chegutu |
| 9. Chinengundu | Chinengundu |
| 10. Norton Selous | Norton Selous |
| 11. Chegutu | Chegutu |
| 12. Chikara | Chikara |
| | Chivero |
| | Dombwe |

Owner field (Before and After)

| Owner | Owner |
| --- | --- |
| Govt. | Govt. |
| Govt. | Govt. |
| GOVT | Govt. |
| Govt. | Govt. |
| Govt | Govt. |
| Govt. | Govt. |
| Govt. | Govt. |
| Municipality | Municipality |
| municipality . | Municipality |
| RDC | Rdc |
| RDC | Rdc |
| RDC | Rdc |

Category field (Before and After)

| Category | Category |
|---|---|
| District Hospital | District Hospital |
| Rural hospital | Rural Hospital |
| RHC | Rhc |
| RHC | Rhc |
|  |  |
| RHC .. | Rhc |
| RHC | Rhc |
| Clinic | Clinic |
| Clinic | Clinic |
| rural Hospit al | Rural Hospital |
| clinic | Clinic |
| Clinic |  |