# Generating Synthetic Data Mimicking the Real Cancer Population

Sepideh Mosaferi

National Cancer Institute (NCI) Presentation
University of Maryland College Park

July 19, 2016

**Intro**

**What Is Synthetic Data?**
Synthetic data are "any production data applicable to a given situation that are not obtained by direct measurement". The synthesized data might be directly obtained by perturbing the real data or based on some assumptions, distributions and frequency tables which might follow a particular real situation.

Why We Produce It?

- For the confidentiality purpose to protect disclosure risk
- For the Research Purpose
- For uniformity of evaluating some procedures or software

### Synthetic Data for the De-duplication & Record Linkage Purpose

The real data set is dirty. Some entities exist multiple times. The formats of fields are not the same, and there might be variety kinds of typos. Therefore if we want to synthesize a real data for the de-duplication or record linkage purpose, we need to consider all of these situations.
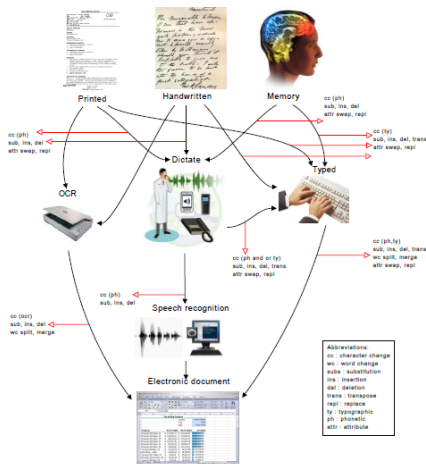
### Purpose & Goal of This Research:

Generating A Synthetic Data Mimicking the Real Cancer Population

- We need to consider different possibilities of errors that might happen in the real situation for the fields
- We need to incorporate frequency distributions or tables obtained from registered data or online reliable sources.

**Cont'd Purpose & Goal of This Research:**

- To have different scenarios we make fields dependent. For example a specific kind of cancer could be dependent to the gender. Or cancer and blood pressure could vary based on the age.

- We consider Zipf, Uniform, or Poisson distributions for randomly selecting units that need to be duplicated.

- Matching variables that we considered are: First name, Last name, SSN, DOB, State, Zipcode, Age, Phone Number, Cancer, Blood pressure, etc.

**Software**

1. DBGen (C programming) developed around 1990s; producing US mailing address
2. GeCo (Python programming) developed from 2002-2008; for the special purpose of bio-medical studies

Both programs are capable of producing one synthesized data set with unique IDs, duplication and errors.
We should give them the frequency tables not the full raw data set.

**DBGen**

1. We cannot incorporate as many errors as we want

2. We can generate fields such as: First Name, Middle Initial, Last Name, Street Number, P.O. Box, Apartment Number, City, State, Zipcode.

3. We can consider Insertions, deletions, replacements, swapping, and other errors.

4. The program is old and not handy. For changing it, we need to write a lot of codes

## Cont'd DBGen

**An Empirical Study (DBGen)**

To test whether the DBGen (since it is old) is good for generating data, we implemented an experimental study. We assumed n (number of units)=c (number of clusters)

I generated an error free data set based on the frequency tables coming from Automated Vital Statistics System under the Institute for Social, Behavioral, and Economic Research at the University of California, Santa Barbara, CA 93106. Number of girls=240311 and boys=252972.

**DBGen Evaluation**

To evaluate the duplicated data set, we considered %5, %10, %20, and %40 duplication. Then we selected six attributes with their frequency distributions, Average Frequencies, and Standard Deviations.
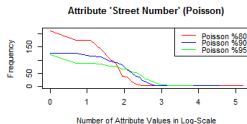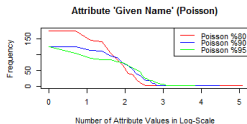
# DBGen Evaluation
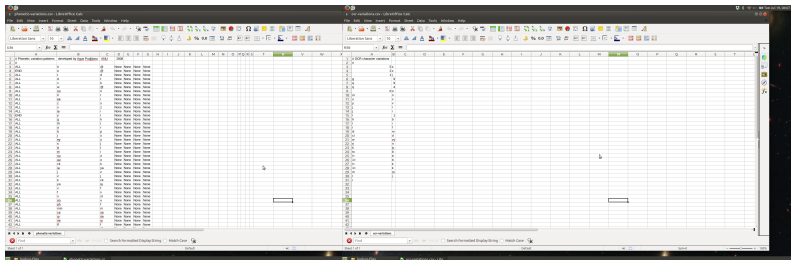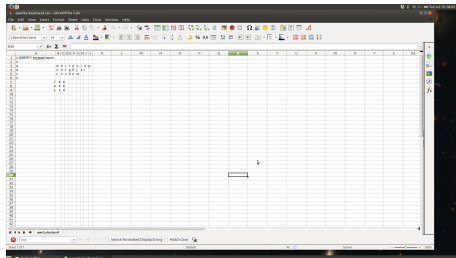
# DBGen Evaluation

# DBGen Evaluation

**GeCo**

1. Capable of incorporating variety kinds of errors, missing values, etc.

2. We can consider dependent attributes.

3. We need to define the number of duplication; otherwise, the program does not work.

4. Writing codes in Python is usually easier rather than C.

**GeCo (As Our Final Focus)**

We wrote codes and merged them with the original GeCo codes. I produced a big data set with variety of errors. Since the ids are unique just per data set not across the data set (generating the data is based on random process), we split the data set into 2 data set with different sizes after permutation for the purpose of record linkage. Unique IDs make it easy to evaluate the matching process since we have the true situations.

**GeCo (Some Possible Errors)**

## How the Produced Data Sets Look Like?

**Febrl (for de-duplication and record linkage)- Beyond the Contract**

Febrl is a program in Python which needs different packages and features on the computer for its implementation. Its installation is time-consuming and different per PC. It has different versions. I considered Febrl 0.4.1 (Dec 2008). It also has a data generator module called dsgen which is not error free and handy.
It contains a GUI. But we are considering the main Python codes to add some other criteria.
1. Indexing 2. Standardization 3. De-duplication, 4. Record Linkage, 5. Evaluation

# Febrl (For Adding Features)

**Febrl (GUI)**

**Some Future Steps**

- Considering more other specific scenarios for the cancer population

- Availability of real data set or frequency tables for some cancer population

- Febrl applications (Beyond the contract)

- Time and Energy

**Thank You!**

Questions: smosafer@umd.edu