

# TEA for survey processing

Center for Statistical Research and Methodology Division  
U.S. Census Bureau

March 15, 2013

## Abstract

TEA is a system designed to unify and streamline survey processing, from raw data to editing to imputation to dissemination of output. Its primary focus is in finding observations that are missing data, fail consistency checks, or risk the disclosure of sensitive information, and then using a unified imputation process to fix all of these issues. Beyond this central focus, it includes tools for reading in data, generating fully synthetic data, and other typical needs of the survey processor.

## 1 Overview

WE INTEND to implement the many steps of survey processing into a single framework, where the interface with which analysts work is common across surveys, the full description of how a survey is processed is summarized in one place, and the code implementing the procedures are internally well-documented and reasonably easy to maintain.

Raw data is often rife with missing items, logical errors, and sensitive information. To ignore these issues risks alienating respondents and data users alike, and so data modification is a necessary part of the production of quality survey and census data. TEA is a statistical library designed to eradicate these errors by creating a framework that is friendly to interact with, effectively handles Census-sized data sets, and does the processing quickly even with relatively sophisticated methods.

This paper gives a detailed overview of the TEA system, its components, and how they're used. If you have this tutorial you should also have a working installation of TEA and a sample directory including `demo.spec`, `demo.R`, and `dc_pums_08.csv`. Basic versions of all of the steps described below are already implemented and running, though the system will evolve and grow as it is applied in new surveys and environments.

## 2 System basics

TEA implements a two step process for addressing issues with raw data. The first is to identify those failures listed above (missing data, logical errors, or sensitive information), and then, having identified problem data, impute new values to replace the old

ones. Although the term *imputation* is typically used only to describe filling in missing data, we use it broadly to mean any modification of a data item that involves choosing among alternatives, regardless of which of the above failures prompted the fill-in.

In terms of how TEA is implemented, we break the process into two parts: the specification of variable details—such as which file to read in, what values should be top-coded, or the full description of the model used for filling in missing data—and the actual procedure to be run by the computer based on those inputs. The specification of details as mentioned above will go into a plain text file, herein called the spec file. Based on your inputs to the spec file (which we will explain later as to what those inputs are/can be), you then run a script in a user-friendly statistical computing framework called **R**. This is where the computing (editing, imputation, etc) takes place. We will explain this in more detail later as well. For now, let's look closer at the spec file:

## 2.1 The spec file

The full specification of the various steps of TEA, from input of data to final production of output tables, that are performed during your implementation of TEA is specified in a single file, the *spec file*. There are several intents to this setup. First, because the spec file is separate from programming languages like R or SAS, it is a simpler grammar that is easier to write, and so analysts can write technical specifications without the assistance of a programmer or a full training in a new programming language. In other words, the spec file allows users whose areas of expertise are not in programming to customize and use TEA in a standardized and accessible environment. To see why this is the case, observe the following script that is taken from the `demo.spec` file (which you can open using Vi or any other text editor):

```
database: test.db

input {
    input_file :dc_pums_08.csv
    output_table: dc_pums
}
```

In this snippet of the `demo.spec` file, we specified a database to use, an input file to be parsed, and an output table to write our imputations to. Behind the scenes, SQL scripts and C functions are being executed. As we will see, other scripts that are run from the spec file perform more complicated behind-the-scenes algorithms. However, before we go through an example of a full spec file, let's take a look at the environment and underlying systems in which TEA runs to gain a better understanding of the processes taking place in the spec file:

## 2.2 Environment and underlying systems

TEA is based on three systems: C, R, and SQL. Each provides facilities that complement the others:

**SQL** is designed around making fast queries from databases, such as finding all observations within a given age range and income band. Any time we need a subset of the data, we will use SQL to describe and pull it. SQL is a relatively simple language, so users unfamiliar with it can probably learn the necessary SQL in a few minutes—in fact, a reader who claims to know no SQL will probably already be able to read and modify the SQL-language conditions in the checks sections below.

The TEA system stores data using an SQL database. The system queries the database as needed to provide input to the statistical/mathematical components (R, C function libraries, etc.). Currently, TEA is written to support SQLite as its database interface; however it would be possible to implement other interfaces such as Oracle or MySQL.

Output at each step is also to the database, to be read as input into the next step. Therefore, the state of the data is recorded at each step in the process, so suspect changes may be audited.

Outside the database, to control what happens and do the modeling, the TEA package consists of roughly 5,000 lines of R and C code.

**R** is a relatively user-friendly system that makes it easy to interact with data sets and write quick scripts to glue together segments of the survey-processing pipeline. R is therefore the interactive front-end for TEA. Users will want to get familiar with the basics of R. As with SQL, users well-versed in R can use their additional knowledge to perform analysis beyond the tools provided by the TEA system.

**C** is the fastest human-usable system available for manipulating matrices, making draws from distributions, and other basic model manipulations. Most of the numerical work will be in C. The user is not expected to know any C at all, because R procedures are provided that do the work of running the underlying C-based procedures.

Now that we have a better idea of the environments in which TEA is run, let's take a look at an example of a full spec file: `demo.spec`:

The configuration system (spec file) is a major part of the user interface with TEA. As is evident from `demo.spec`, there are many components of the spec file that all perform certain functions. We begin by explaining the concept of *keys*:

## 2.3 Keys

Everything in the spec file is a key/value pair (or, as may be more familiar to you, a tag: data definition). Each key in the spec file has a specific purpose and will be outlined in this tutorial. To begin, we start in the header of `demo.spec`:

```
[ database: demo.db
  id: SSN
```

Here, `database: demo.db` and `id: SSN` are examples of key: value pairs. As is the case in `demo.spec`, you will always need to begin your spec file by declaring the database (`database: your_database.db`) key and the unique identifier (`id: your_unique_identifier`) key. The `database:` key identifies the relational database where all of your data will be manipulated during the various processes of TEA. The `id` key provides a column in the database table that serves as the unique identifier for each set of data points in

your data set. Though the `id` key is not strictly necessary, we strongly advise that you include one in your `spec` file to prevent any unnecessary frustration as most of TEA's routines require its use. More information on both of these keys can be found in the appendix of this tutorial.

You may have noticed that keys are assigned values with the following syntax:

```
key:  value
```

This syntax is equivalent to the following form with curly braces:

```
database{
    demo.db
}

id {
    SSN
}
```

Clearly, this form is not as convenient as the `key: value` form for single values. However, it allows us to have multiple values associated with a single line of data, and even subkeys. For example, take the next line in `demo.spec` (the computer will ignore everything after a `#`, so those lines are comments for your fellow humans):

```
input {
    #As above, if you have a key with one value, you can use a colon:
    input file: ss08pdc.csv
    overwrite: no
    output table: dc
    #However, with curly braces we can associate multiple values with a single key
    types {
        AGEp: integer
        CATAGE: integer
    }
}
```

In the database, here is what the above will look like:

```
database      ourstudy.db
checks        age > 100 => age = 100
checks        status = 'married' and age < 16
input/input file  text_in.csv
input/output table dataset
```

Observe that the subkeys are turned into a slash-separated list. It is worth getting familiar with this internal form, because when you've made a mistake in your `spec` file, the error message that gets printed in R will display your keys in the above form. We will discuss this more later in the tutorial when we talk about running your `spec` file in R. Here is a succinct summation of the syntax rules for keys:

- You can have several values in curly braces, each on a separate line, which are added to the key/value list. Order is preserved.
- If there is a subkey, then its name is merged with the parent key via a slash; you can have as many levels of sub-sub-sub-keys as you wish.
- As a shortcut, you can replace key {single value} with key: single value.
- Each key takes up exactly one line (unless you explicitly combine lines; see below). Otherwise, white space is irrelevant to the parser. Because humans will also read the spec file, you are strongly encouraged to use indentation to indicate what is a member of what group.
- If you need to continue a value over several lines, put a backslash at the end of each line that continues to the next; the backslash/newline will be replaced with a single space.

We now continue through our spec file to the `input` key:

## 2.4 Input declarations

The `input` key is an important key in your spec file that specifies:

1. The `csv` file from which you will draw your data.
2. The option to *overwrite* the `csv` file currently written into the specified database with a new `csv` file that you have specified.
3. The output `table` key that specifies the table where you would like to write the data that is read in from the `csv` file.
4. The `types` key that specifies what type of variables are to be written into the output table.

If this all seems confusing, do not fret. The layout of the spec file will make more sense as we continue to explain its various features. Again, more information about these keys can be found in the appendix.

We now discuss the `fields` key in more detail.

## 2.5 Field declarations

The edit-checking system needs to know the type and range of every variable that has an associated edit. If a variable does not have an associated edit, there is no need to declare its range.

The declaration of the edit variables contained in the `fields` key consists of the field name, an optional type (`int`, `text`, or `real`) to be discussed further below, and a list of valid values. Here is a typical example:

```
fields {
  age: 0–100
  sex: M, F
  hh_type: int 1, 2, 4, 8
  income: real
}
```

In the above code, we've declared four variables: `age`, `sex`, `hh_type`, and `income` and we've declared those four variables in different and valid ways. By doing this, the edit-checking system knows what to expect for each of these variables. You can see that the list of values may be text or numeric, and the range 0–100 will be converted into the full sequence 0, 1, 2, ..., 100. By declaring our variables with a type and range, we can pass the information to the edit-checking system so that it knows what to verify when running its checks for each of these variables.

When declaring a variable, the first word following the `:` may be a type. For instance, for the `hh_type: int 1, 2, 4, 8` field above, we used the word `'int'` to indicate that the data values of the field were of type integer. Notice as well that a type does not necessarily need to be declared; in which case the default action is to treat the entry as plain text (which means that numeric entries get converted to text categories, like `"1"`, `"2"`, ...). Keep in mind that if you declare the incorrect type for a field that the edit-checking system may not correctly verify the values of that field in your data set.

**Auto-detect** If your field list consists of a single star, `*` (the wild-card character in SQL that represents all possible inputs), then the data set will be queried for the list of values used in the input data. All values the variable takes on will be valid; all values not used by the variable will not be valid. Keep in mind that this can present problems if there are errors in your input data set. In any case, using `*` may be useful when quickly assembling a first draft of a spec, but we recommend giving an explicit set of values for production. You may precede the `*` with a type designation, like `age: real *`.

As you know, it is often necessary to impute data points in separate categories given the distribution of the data. For this, we use the `recodes` key.

## 2.6 Recodes

In short, recode keys are just new variables that are deterministic functions of existing data sets based on parameters given by you, the analyst. Based on the variables specified in the `fields` key, you can "recode" those fields into other variables (which can be thought of as categories) based on parameters given in a certain syntax. This is often necessary to ensure that you are imputing your data over an accurate distribution. Observe the following typical example of a `recodes` key:

```
recodes {
  CATAGE {
    0 | age between 0 and 15
    1 | age between 16 and 64
```

```

    2 | age > 64
  }
}

```

Here, we have declared a new variable CATAGE whose data points are based off of a deterministic function of the variable AGE. As we will see later, this recode will be called in the `categories` key during imputation so that the data points we are attempting to impute will be imputed in categories based on their recode values rather than collectively as a single set of data points. The `recodes` key is fairly straightforward, although we will learn about some of its more advanced features later. For now, however, we continue our walkthrough of the `spec` file by discussing the `checks` key.

## 2.7 Checks

The consistency-checking system is complex for efficiency purposes: there are typically a few dozen to a few hundred checks that every observation must pass, from sanity checks like *fail if age < 0* to real-world constraints like *fail if age < 16 and status='married'*. Further, every time a change is made (such as by the imputation system) we need to re-check that the new value does not fail checks. For example, an OLS imputation of age could easily generate negative age values, so the consistency checks are essential for verifying that the imputation process is giving us accurate and useable data points. In addition to error checking we can also use consistency checks for other purposes, such as setting *structural zeros* for the raking procedure.

All of the checks that are to be performed are specified here, in the `checks` key. The checks you specify here will be performed on all input variables, as well as on all imputed data values. Let's take a look at an example `checks` key:

```

checks {
  age < 0
  age > 95 => age = 95
}

```

In the above example, we've indicated that the consistency checking system should verify that age is not less than 0 and that age is not greater than 95. Notice as well that when specifying that `age > 95` we've also included the line `'age = 95'` to indicate that when an age data point has a value higher than 95 that we should simply top-code it as 95. We haven't included an auto-declaration for `age < 0` because if a data-point has a negative age value that it's indicative of a real error rather than that should be properly imputed. It is up to you to decide when it is appropriate to utilize the auto-declaration feature.

We've now introduced the keys that precede the `impute` key; and describe its functions below.

## 2.8 Imputation

The `impute` key is fairly comprehensive and has several sub-keys that fulfill various roles in the imputation process. Many of the values of these sub-keys are derived from

values found earlier in the spec file; such as the fact that `categories` is based off of the variables declared in `recodes`. Take a look at the following example of an `impute` key that is used to outline the imputation of the `age` variable described in the above keys:

```
impute {  
  input table: viewdc  
  min group size: 3  
  draw count: 3  
  seed:2332  
  
  categories {  
    CATAGE  
    SEX  
  }  
  
  method: hot deck  
  output vars: age  
}
```

As you can see, there is quite a bit going on here. Let's walk through each of the sub-keys above and see what they're doing:

- The `input table` sub-key specifies the name of a view table to where you will write your recode variables. Recall that a view table is a digital SQL table that takes up very little memory and is very quick to render. In our example, we've indicated that the name of the view table to where our recode variables will be written is called 'viewdc'. We'll be able to access this view table from R later when we actually compile our spec file.
- `min group size` indicates the minimum number of known data points that can be used to impute any unknown data points. For instance, it wouldn't make much sense to use a single data point to impute five other data points as they would all end up being the same value. To prevent this, we set `min group size: 3`.
- `draw count` determines the number of multiple imputations that should be performed. Because imputation is often done using stochastic models, it is often beneficial to obtain multiple imputed data sets. This will be explained in more detail later. For now, just know that the `draw count` sub-key determines the number of multiple imputations that are performed.
- `seed` specifies the pseudo-random number generator seed that is used to obtain random numbers for stochastic imputation models. Keeping this seed the same will ensure that multiple imputations of the same data set will return the same



outputs. Be sure to record the seeds you choose for your imputations so that other analysts can reference your data and replicate your imputations if necessary.

- `categories` specifies which recodes will be used to impute the data set for the given variable. For example, if `CATAGE` is listed in the `categories` key then the data points that correspond to each of the three possible `CATAGE` values will be imputed separately. While it's not an absolute requirement that the `categories` key be implemented, because data sets are typically imputed by categories you will most likely be using this key for almost all of your imputations.
- `method` specifies the type of model that is used to impute your data points. Choosing this model is itself a non-trivial task and is discussed more thoroughly later on in the tutorial.
- `output vars` specifies the variable whose data points are to be imputed. For this reason you could consider `output vars` as both the 'input' variable as well as the 'output' variable to the model specified in `method`.

More can be found on each of the above keys in the appendix.

## 2.9 More features

This concludes our walkthrough of a typical `spec` file. By now you should have a basic idea of how a `spec` file is implemented within TEA. Before we continue on to explaining more about imputation, the models available in TEA, and other features that are available within the TEA framework, we will mention two more features of the `spec` file.

**Group recodes** The first feature that we'll mention is that of `group recodes`. Often in an imputation model, after declaring your `recodes`, it is also necessary to declare sub-recodes that are function of those earlier recodes. For example, suppose you had a `spec` file in which you intended to impute data points for various members of households. Then to perform imputations on

```
recodes {  
    EARN: case when WAGP<=0 then 'none' when WAGP>0 then 'black' else NULL  
        end  
    MOVE: case MIG when 1 then 'moved' else 'stayed' end  
    DEG: case when SCHL in ('24','23','22','21','20') then 'degree' else 'non-  
        degreed' end  
    MF: case SEX when '1' then 'Male' when '2' then 'Female' else NULL end
```

```

        REL: case RELP when '00' then 'Householder' when '01' then 'Spouse' \
              when '02' then 'Child' \
              when '06' then 'Parent' else NULL end
    }
    group recodes {
        group id column: SERIALNO
        recodes {
            NP: max(SPORDER)
            NHH: sum(RELP='00')
            NSP: sum(RELP='01')
            NUP: sum(RELP='15')
            HHAGE: max(case RELP when '00' then AGE end)
            SPAGE: max(case RELP when '01' then AGE end)
            SPORD: max(case RELP when '01' then SPORDER end)
            HHSEX: cast(round(avg(case RELP when '00' then SEX end)) as integer)
            SPSEX: cast(round(avg(case RELP when '01' then SEX end)) as integer)
        }
    }
}

```

**Including** Up to this point, we've constructed the entire spec file in a single file. Though this approach may be appropriate in some scenarios, it's often the case that you and your colleagues would like to make edits to the same spec file and update or construct it concurrently. To do this, you can simply mark subsidiary files to be included in the main spec file. Doing so allows you to easily combine these subsidiary files into one project through which all of TEA's routines will run.

The syntax for including is quite simple. To include a subsidiary file at a certain point in the spec file, simply insert the key `include: subsidiary_file_name` at the line in the spec file where you would like the contents of the subsidiary file to be inserted. For example, if you have written the consistency checks in a file named `consistency`, and the entire rank swapping configuration was in a file named `swap`, then you could use this parent spec file:

```

database: test.db

include: swap

checks {
    include: consistency
}

```

Note that any subsidiary files that you choose to include in your spec file must be in the same directory as the spec file or it will not be able to find them.

This concludes our discussion of the spec file layout and syntax. At this point, you should be able to implement a basic spec file to impute any data set. More information about the keys discussed above and others that were not present in `demo.spec` can be found in the appendix and at the `r-forge` website.

We now continue our tutorial by discussing imputation in more detail.

### 3 Imputation

Thus far we've seen how to impute a data set using a single `impute` group. In this form, single imputation produces a list of replacement values for those data items that fail consistency checks or initially had no values. For the case of editing, the replacements would be for data that fails consistency checks; for the case of basic imputation, the replacements would be for elements that initially had no values. Recall that as we stipulated earlier in the tutorial, for our purposes we consider a looser definition of imputation that includes the replacement of data points that both initially had no value as well as those that fail consistency checks.

In a similar vein, while single imputation gives us a single list of replacement values, any stochastic imputation method could be repeated to produce multiple lists of replacement values. For instance, we could utilize multiple imputation to calculate the variance of a given statistic, such as average income for a subgroup, as the sum of within-imputation variance and across-imputation variance.

To give a concrete example, consider a randomized hot deck routine, where missing values are filled in by pulling values from similar records. For each record with a missing income:

1. Find the universe of which the record is a member.
2. For each variable in turn:
  - Make a draw from model chosen for the variable, based on the specific universe from (1).

The simplest and most common example is the randomized hot deck, in which each survey respondent has a universe of other respondents whose data can be used to fill in the respondent's missing values. The hot deck model is a simple random draw from the given universe for the given variable.

Given this framework, there are a wealth of means by which universes are formed, and a wealth of models by which a missing value can be filled in using the data in the chosen universe.

**Universes** The various models described above are typically fit not for the whole data set, but for smaller *universes*, such as a single county, or all males in a given age group. A universe definition is an assertion that the variable set for the records in the universe was generated in a different manner than the variables for other universes.

An assertion that two universes have different generative processes could be construed in several ways:

- different mathematical forms, like CART vs. logistic regression
- One broad form, like logistic regression, but with different variables chosen (that is, the stochastic form is the same but the structural form differs).
- A unified form, like a logistic regression with age, sex, and ancestry as predictors, with different parameter estimates in each universe.

Universe definitions play a central role in the current ACS edit and imputation system. Here, universes allow data analysts to more easily specify particular courses of action in the case of missing or edit-inconsistent data items. To give an extreme example, for the imputation of marital status in ACS group quarters (2008), respondents are placed in two major universes: less than 15 years of age (U1) and 15 or more years of age (U2). The assertion here, thus, is that people older than 15 have a marital status that comes from a different generative process than those people younger than 15. This is true: people younger than 15 years of age cannot be married! Thus in the system, any missing value of marital status in U1 can be set to “never married”, and missing values in U2 can be allocated via the ACS hot-deck imputation system.

Now that we are more familiar with universes, we can discuss the various models that are available in TEA and the methodology of choosing the one that is appropriate for the imputation being performed.

## 4 Models

Now that we’re more familiar with the concept of universes, we examine how to choose the model that will give us the best results when imputing the data points in a specific universe. Indeed, given an observation with a missing data point and a universe, however specified, there is the problem of using the universe to fill in a value. Randomized hot-deck is again the simplest model: simply randomly select an observation from the universe of acceptable values and fill in. Other models make a stronger effort to find a somehow-optimal value:

- One could find the nearest neighbor to the data point to be imputed (using any of a number of metrics).
- Income is typically log-Normally distributed, and log of this year’s income, last year’s income, and if available, income reports from administrative records (ADREC) may be modeled as Multivariate Normal (with a high correlation coefficient among variables). After estimating the appropriate Multivariate distribution for the universe, one could make draws from the distribution.
- If the Multivariate Normal seems implausible, one could simply aggregate the data into an empirical multivariate PDF, then smooth the data into a kernel density estimator (KDE). Draws from a specified KDE can be made as easily as draws from a standard Multivariate Normal.
- Count data, such as the number of hospital visits over a given period, are typically modeled as a Poisson distribution. In a manner similar to fitting a Normal, one could find the best-fitting Poisson distribution and draw from that distribution to fill in the missing observation.
- Bayesian methods: if ADREC are available, they may be used to generate a prior distribution on a variable, which is then updated using non-missing data from the sur-

vey.

- One could do a simple regression using the data on hand. For example, within the universe, regress income on age, race, sex, and covariates from ADREC. Once the regression coefficients are calculated, use them to impute income for the record as the predicted value plus an error term.

- Discrete-outcome models, like the Probit or Logit, could be used in a similar manner.

- To expand upon running a single regression, one can conduct a structured search over regression models for the best-fitting regression.<sup>1</sup> Such a search is currently in use for disclosure avoidance in ACS group quarters.

A unified framework would allow comparison across the various imputation schemes and structured tests of the relative merits of each. Though different surveys are likely to use different models for step (2) of the above process, the remainder of the overall routine would not need to be rewritten across surveys.

We now discuss each of the models that are available in TEA.

## 4.1 Models

This section describes the models available for use in imputing missing data.

**Hot deck** This is randomized hot deck. Missing values are filled in by drawing from nonmissing values in the same subset. The assumption underlying the model is *missing completely at random* (MCAR), basically meaning that nonrespondents do not differ in a systematic way from respondents.

This model has no additional keys or options, although the user will probably want an extensive set of category subsets. Example:

```
impute{
  input table: dc
  min group size: 3
  draw count: 5
  id: serialno

  categories {
    agecat
    sex
  }
  output vars: sex
  method: hot deck
}
```

905,1  
76%

---

<sup>1</sup>Because the imputation step is not an inference step, such a search presents few conceptual difficulties.

```

    }

    impute{
      input table: dc
      min group size: 3
      draw count: 5
      id: serialno

      categories {
        agecat
        sex
      }
      method: ols
    }

    output vars: agep
    input vars: race1p, nativity1sex
  }
}

```

923,1  
78%

**Ordinary least squares (aka regression)** This is the familiar  $y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \epsilon$  form. The parameters (including the variance of  $\epsilon$ ) are estimated for the subset of the category where all variables used in the regression are not missing. For a point where  $y$  is missing but  $(x_1, x_2, \dots)$  are not, find  $\beta_1 x_1 + \beta_2 x_2 + \dots$ , then add a random draw from the distribution of  $\epsilon$ .

The variables may be specified via the usual SQL, with two exceptions to accommodate the fact that so many survey variables are categorical. typical meaning.

**Probit and Logit** These models behave like OLS, but are aimed at categorical variables. The output variable is not restricted to two categories.

### seqRegAIC

**Distributions** The Normal (aka Gaussian) distribution is the archetype of this type of model. First, estimate the mean and standard deviation using the non-missing data in the category. Then fill in the missing data via random draws from the Normal distribution with the calculated parameters. You may use either `method: normal` or `method: gaussian`.

For other situations, other distributions may be preferable. For example, income is typically modeled via `method: lognormal`. Count data may best be modeled via `method: poisson`.

Hot deck is actually a fitting of the Multinomial distribution, in which each bin has elements in proportion to that observed in the data; `method: hot deck` and `method: multinomial` are synonyms.

The *method: multivariate normal* doesn't work yet.

The distribution models have no additional options or keys.

**Kernel smoothing** A kernel density estimate overlays a Normal distribution over every data point. For example, if a data set consisted of two data points at 10 and one at 12, then there would be a large hump in the final distribution at 10, a smaller hump at 12, and the distribution would be continuous over the entire space of reals.

Thus, kernel smoothing will turn a discrete distribution consisting of values on a few values into a continuous distribution.

Invoke this model using either method: `kernel` or method: `kernel density`.

## 5 Overlaying

Each record in a data set might have several data items which require imputation for different reasons:

- a record could have several inconsistent items that we must replace
- an otherwise consistent record could have a blank item that we wish to impute.
- a record could have a combination of consistent items that could lead to personal identification

Each scenario implies slightly different knowledge about the data, and thus each scenario might require a different imputation method to properly use this knowledge.

An overlay is a secondary data table (or set of tables) that gives information regarding the *reason* for imputation. Using missing data as an example, a simple overlay could have an entry for each item in the data, indicating if that item is missing or not. A more complicated overlay could delineate the type of non-response for each data item.

## 6 Raking

Raking fills in a table where the totals for each column or row are known, but individual cell values may vary. It can be shown to find the optimum given a log-linear model on certain contrasts (i.e., a series of sets of dimensions).

The contrasts indicate which values are to remain fixed (while others may vary). To give an example, let us say that our data set is (age  $\times$  sex  $\times$  race  $\times$  block), but our log linear model is over the (race  $\times$  block) contrast. A sample data set:

age	sex	race	block
10	M	A	1
20	F	A	2
30	M	B	3
40	F	B	3

The raking algorithm would never be able to fill in a value in the cell (40, F, C, 3), for example, because there are no Cs in block 3, and the (race  $\times$  block) contrast guarantees that the counts of each (race, block) total won't change. But that's our only

assurance: (50, M, B, 3) is possible, because age and sex totals could be raked into any value.

For this example, the set of cells that could have a value are:

(any age)  $\times$  (M, F)  $\times$  A  $\times$  1  
(any age)  $\times$  (M, F)  $\times$  A  $\times$  2  
(any age)  $\times$  (M, F)  $\times$  B  $\times$  3

The algorithm takes in the location of a data table in the database and a set of contrasts, then rakes the data to fit the given contrasts.

**Internals** The Census Bureau often deals with tables that would have 1.5 million cells if written out in full, but where only maybe 20,000 surveys are in hand. Therefore, the algorithm is heavily oriented toward sparse data.

It begins with a long and tedious routine to write SQL to generate the set of possibly-nonzero values, as per the example above. SQL is the appropriate language for generating this list because it is optimized for generating the cross of several variables and for pruning out values that match our criteria. The tedium turns out to be worth it: our test data set takes about 25 seconds to run using the original full-cross ‘72 algorithm, and runs in under two seconds using the SQL-pruned matrix.

Recall that we had briefly discussed multiple imputation in the `impute` section above. We now discuss this in more detail.

## 7 Multiple Imputation

A single imputation would produce a list of replacement values for certain data items. Any stochastic imputation method could be repeated to produce multiple lists of replacement values. Variance of a given statistic, such as average income for a subgroup, is then the sum of within-imputation variance and across-imputation variance.

- For each imputation:
  - fill in the data set using the given set of imputed values
  - calculate the within-imputation variance
- Sum the across-imputation and average within-imputation variances to produce an overall variance figure that takes into account uncertainty due to imputation.

The question of what should be reported to the public from a sequence of imputations remains open. The more extensive option would be to provide multiple data sets; the less extensive option would be to simply provide a variance measure for each data point that is not a direct observation.

Interactive R The specification file described to this point does nothing by itself, but provides information to procedures written in R that do the work. The analyst would therefore write a script that calls the needed procedures in sequence. For example, this script loads the TEA library into R’s memory, reads the spec file, fixes inconsistencies, and imputes missing values.



```
library("tea")  
read_spec("spec")  
doChecks()  
doMImpute()  
checkOutImpute()
```