# Generalizing from Purposive Surveys
# How large a Sample is Needed

Richard Garfield
Columbia University School of Nursing

Jared P. Lander
JP Lander Consulting

December 20, 2012

## Contents

# 1  Distribution Functions

These are the functions used to calculate the distribution of each answer. They are general and should work with any question.

```
getwd()
```

```
[1] "C:/Users/Jared/week2/writeup/distFuncs"
```

```
# Distribution functions
require(useful)
require(plyr)
## builds the distribution for a given question
build.dist <- function(data, lhs, group, question)
{
    theFormula <- build.formula(lhs = lhs, rhs = c(group,
```

1

```r
        question))
    agg <- aggregate(theFormula, data, length)
    agg <- ddply(agg, .variables = group, .fun = function(x)
    {
        x$Percent <- x[[lhs]]/sum(x[[lhs]])
        return(x)
    })
    agg
}
## get random tehsils from a province
village.list <- function(x, num = 5, unit = "Tehsil")
{
    # get list of units
    units <- unique(x[, unit])

    # sample num of those without replacement
    keepers <- sample(x = units, size = min(num, length(units)),
        replace = FALSE)

    return(as.character(keepers))
}
# function to make names of dist's better
change.names <- function(names, include = names, prefix = "")
{
    theOnes <- which(!names %in% include)
    names[theOnes] <- sprintf("%s.%s", prefix, names[theOnes])
    return(names)
}
## function to impute missing
impute.col <- function(col, value = 0)
{
    col[is.na(col)] <- value
    return(col)
}
## this compares two distributions and computes an MSE
compare.dist <- function(full, partial, compare = "Percent",
    by = intersect(names(full), names(partial)))
{
    # prepend Pull onto certain names in full
    names(full) <- change.names(names = names(full), include = by,
        prefix = "Full")

    # prepend Partial onto certain names in full
    names(partial) <- change.names(names = names(partial),
        include = by, prefix = "Partial")
```

```r
    full.compare <- sprintf("Full.%s", compare)
    partial.compare <- sprintf("Partial.%s", compare)

    # join the two together
    both <- join(x = full, y = partial, by = by, type = "left")

    rm(full, partial)

    ## fill in any NA's with zero
    both[[full.compare]] <- impute.col(col = both[[full.compare]],
        value = 0)
    both[[partial.compare]] <- impute.col(col = both[[partial.compare]],
        value = 0)

    both$.Diff <- both[[full.compare]] - both[[partial.compare]]

    both$.MSE <- mean(both$.Diff^2)

    # attr(x=both, which='MSE') <- mean(both$.Diff^2)

    # aggregate(build.formula(lhs='.Diff', rhs=

    return(both)
}
```

## 2 Initial Stuff

The data is as described in Section 3.

We examined the answer to the question "What percentage of rice crops were lost due to the flood?" We then randomly chose five Tehsils from each province, then 10, then 15 and performed the same analysis on the reduced data.

In situations where a province has fewer than five, 10 or 15 Tehsils sampled, all were used.

## 3 The Data

The data were collected following the floods in Pakistan in 2010. Small Changes.

It surveyed affected villages in GB, KPK, Punjab and Sindh.

The distribution of villages within Tehsils within Provinces is seen in Figure 1.

```
ggplot(vills, aes(x = Tehsil)) + geom_bar(aes(y = Village),
    stat = "identity") + opts(axis.text.x = theme_text(angle = 270,
    hjust = 0)) + facet_wrap(~Province, scales = "free_x")
```



Figure 1: Distribution of villages within Tehsils within the four Provinces.

The analysis begins in Section 4.

# 4  Analyzing All Data

Here we analyze all of the data.

First we load the data and view a portion of it. Some more details.

These are the necessary packages.

```
require(useful)
require(plyr)
require(ggplot2)
```

```
load("../../data/pakistan/pak.rdata")
source("../../R/distFuncs.r")
corner(pak, c = 15)

  New_ID Age  Sex     Date Province District Tehsil
1   1288  26 Male 29082010      KPK  Shangla Besham
2   1290  30 Male 29082010      KPK  Shangla Besham
3   1370  54 Male 28082010      KPK  Shangla Besham
4   1372  53 Male 28082010      KPK  Shangla Besham
5   1371  64 Male 28082010      KPK  Shangla Besham
        Village Latitude Longitude Total Urban Rural
1 abaseen colony    34.94     72.88  90.6     -  90.6
2 abaseen colony    34.94     72.88  90.6     -  90.6
3 abaseen colony    34.94     72.88  90.6     -  90.6
4 abaseen colony    34.94     72.88  90.6     -  90.6
5 abaseen colony    34.94     72.88  90.6     -  90.6
                              Accommodation
1 Collective centers (school/Public building)
2                               Host family
3           On the site of the house (Damaged)
4           On the site of the house (Damaged)
5           On the site of the house (Damaged)
  StagnantWater
1          Few
2          Few
3          Few
4         None
5         None
```

Now we build a distribution for all the data and visualize it in Figure 2 with the code here:.

```
ricePerc <- build.dist(data = pak, lhs = "New_ID", group = "Province",
    question = "RiceLost")
ricePerc$Size <- "All"
ggplot(ricePerc, aes(x = RiceLost, y = Percent)) + geom_bar(stat = "identity") +
    facet_wrap(~Province) + opts(axis.text.x = theme_text(angle = 90))
```

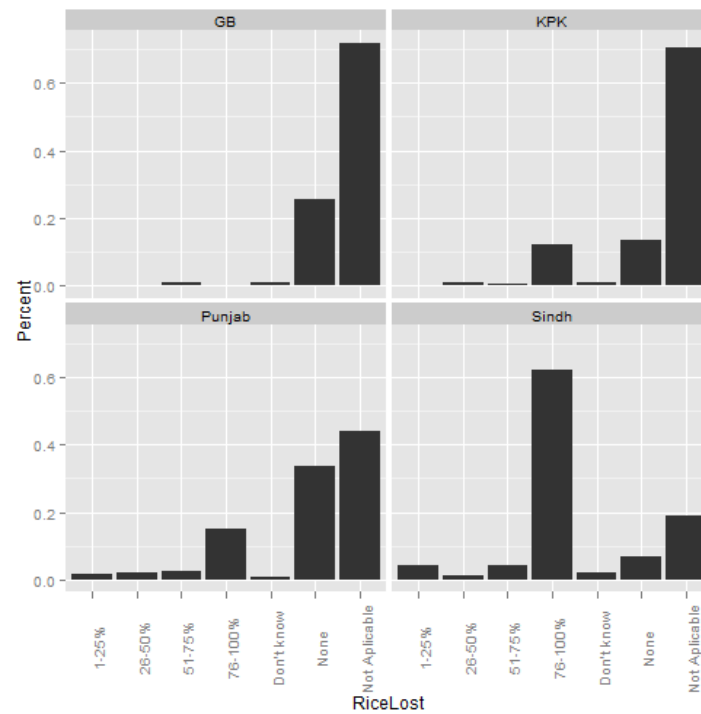In Section **??** we analyze the distribution of responses for samples of fewer Tehsils.

Figure 2: Graphical view of the distribution of responses for all the data.

# List of Figures