

CS 208 - Applied Privacy for Data Science Homework 1

Jason Huang
Spring 2019 - Harvard University

The public Github repo containing all work is at <https://github.com/TurboFreeze/cs208hw>. All code has also been included in the appendix of this PDF as specified.

Problem 1

The dataset was loaded into R, where preliminary data exploration took place. Most notably, there are 25766 entries and 18 variables in the dataset, with the variables being:

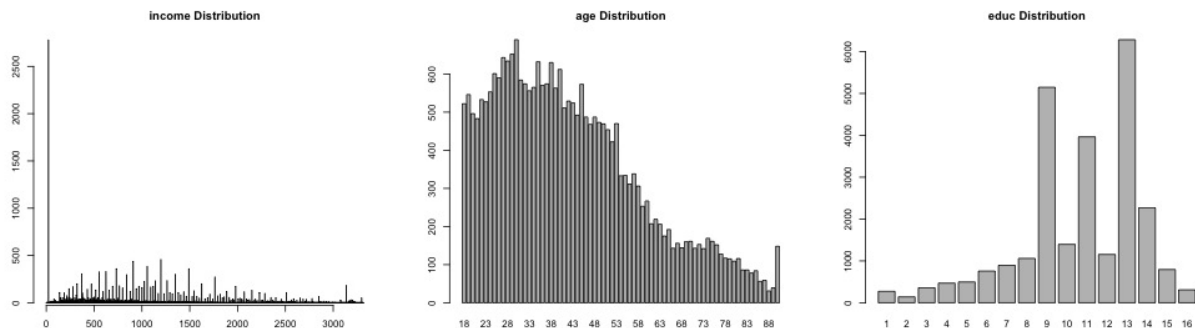
state, puma, sex, age, educ, income, latino, black, asian, married,
divorced, uscitizen, children, disability, militaryservice, employed,
englishability, fips

The naive and most effective starting point is tallying the unique values for each of these variables. Two of the variables, `state` and `fips`, have the same value for all rows and therefore will not be considered at all.

Now examine the most identifying variables (i.e. the variables with most unique values in the dataset, as determined by taking the length of its count table in R).

income	age	educ	puma
2763	73	16	7

All the other variables are binary variables (can only have two distinct values). `puma` is implicitly included (since the goal is to uniquely identify individuals *within a PUMA region*) and will be accounted for in the last step. However, these other three variables should be explored, particularly by examining their distributions, as plotted below.



The goal now is to quantify exactly how identifiable each variable is. Estimates for the probability of collision (of another individual having the exact same value) for each variable based on the largest bin (i.e. worst case) are made below. Actual probabilities for collision can also be calculated by the following formula:

$$p_{collision} = p(X_1 = X_2) = \sum_x p(X_1 = x)p(X_2 = x) = \sum_x \left(\frac{\text{count}(x)}{25766} \right)^2$$

where x is to take on all possible values for that variable.

income is the most identifiable variable but also raises an issue when examining the distribution histogram: a large number of individuals have zero income, and are therefore substantially more difficult to uniquely identify. Apart from this unique case, though, no bin has more than 500 individuals, which corresponds to approximately $500/25766 \approx 2\%$ of collision with another individual (actual: 0.01555964). **age** is distributed a lot better, with the worst case bin having no more than 800 individuals with the same age. When considering the overall size of 25766, that means that there is *at most* a $800/25766 \approx 3\%$ chance of having the same age as a randomly chosen person in the dataset (actual: 0.01831928). Lastly, for **educ**, with up to 6500 in the same bin, the probability of having the same education level as someone else would be $6500/25766 \approx 25\%$ (actual: 0.1416453).

Percentage estimates were given above to provide intuition and also as a sanity check, but the actual values calculated will be the ones used from here on.

When taking these three variables together, the probability of the collision is the intersection of all three variables colliding, which is the product of the three probabilities. This gives an overall probability of colliding to be $0.01555964 \times 0.01831928 \times 0.1416453 \approx 0.00004 = p_{collision}$.

Now, the geographic identifier of PUMA regions needs to be taken into account. Here is a summary of counts for each PUMA region.

PUMA	1101	1102	1103	1104	1105	1106	1107
Count	3215	5736	3728	3740	3128	3236	2983

These are very roughly equal (i.e. on the same order of magnitude), and these do appear to be 5% samples (with full populations of 50000-120000 in each PUMA, which is appropriate). Assume the average PUMA region to therefore be $1/7$ of the total dataset: $(1/7) \times 20 \times 25766 = 73617$.

To finally determine the percentage of individuals $p(U)$ within a PUMA that can be uniquely identified by the aforementioned variables, this calculation simply involves the probability of no collision with anyone in that region of size n , or:

$$p_{unique} = (1 - p_{collision})^n$$

In this situation, $p_{collision} = 0.00004$, $n = 73617$ as calculated above, so:

$$p_{unique} = 0.051183922689063$$

\therefore With the three variables of **income**, **age**, **educ**, approximately 5% of the population within a single PUMA can be identified.

Better reconstruction results can be achieved using more variables. The only remaining variables are binary, which should have roughly 50% chance of collision (though in actuality it will be higher due to uneven distribution). The lowest probabilities of collision are for **sex** (0.5014973) and **married** (0.5046546), two common and relatively evenly distributed binary indicators. When factoring these in, then:

$$p_{collision} = 0.00001 \implies p_{unique} = 0.471312198919265$$

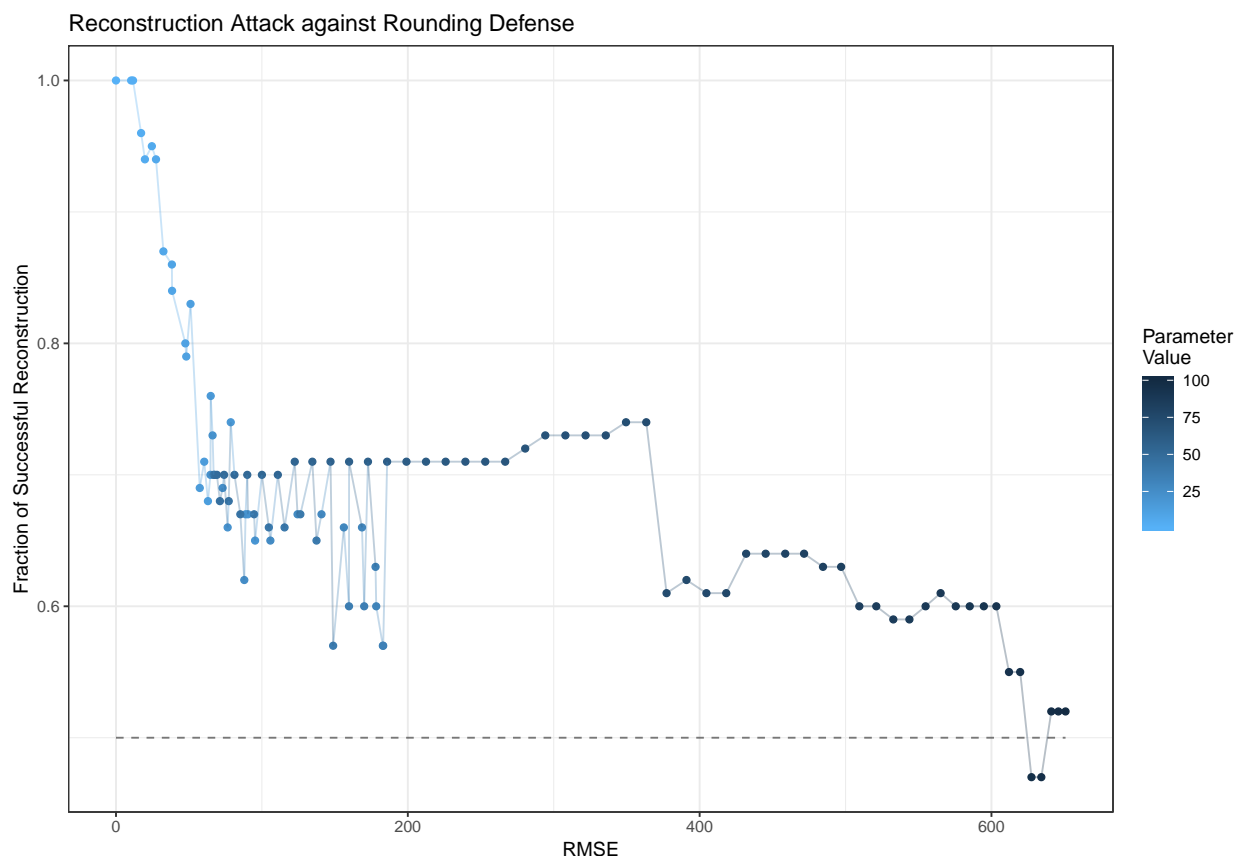
∴ With the five variables of **income**, **age**, **educ**, **sex**, **married**, approximately 47% of the population within a single PUMA can be identified.

While it is more difficult and less likely to orchestrate a reconstruction attack with large numbers of variables, out of interest and completeness, here are some further results. With a sixth variable of **black**, 68% of the population can be uniquely identified. With a seventh added variable of **employed**, 81% of the population can be uniquely identified. The remaining binary variables have sharply decreasing utility due to their uneven distribution (i.e. almost all individuals have the same value and therefore it is not very helpful in identifying someone).

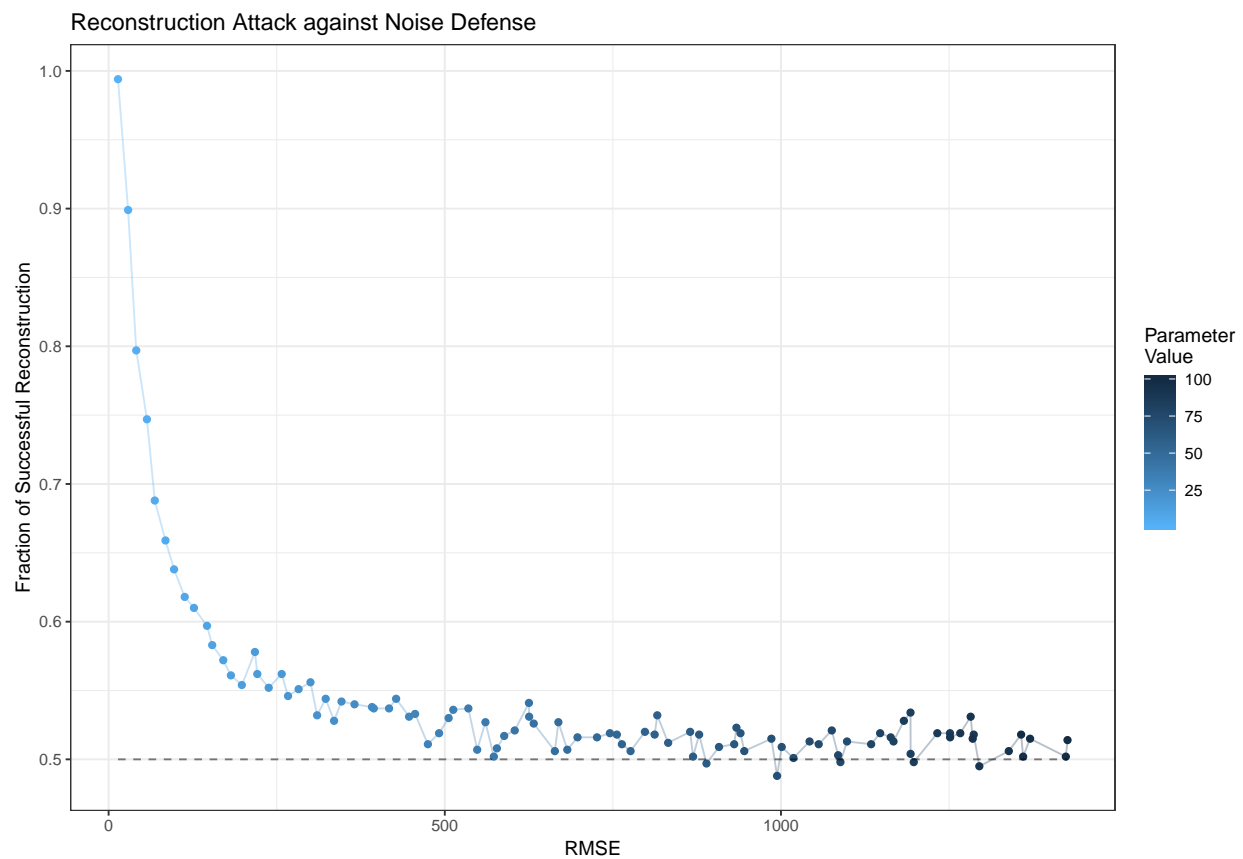
Finally, some small disclaimers. The results above are all approximate and derived from rough back-of-the-envelope calculations. They are also assuming that these variables are available in external sources for cross-referencing, which may present practical obstacles that may or may not be easy to handle. For example, a specific education encoding is used here. Other government databases may use the same schema, making cross-referencing trivial. Furthermore, knowing an individual's exact education level may also be sufficient to map it to one of the factors in this dataset's **educ** schema. However, if a different, less granular schema was used (i.e. only 5 levels instead of 16), then cross-referencing may not really be possible. Different levels of granularity would be a particularly common issue for **income**, with added concerns such as rounding. It may also affect **age**, though probably to a much lesser extent. Binary variables should otherwise be less problematic in this regard.

Problem 2

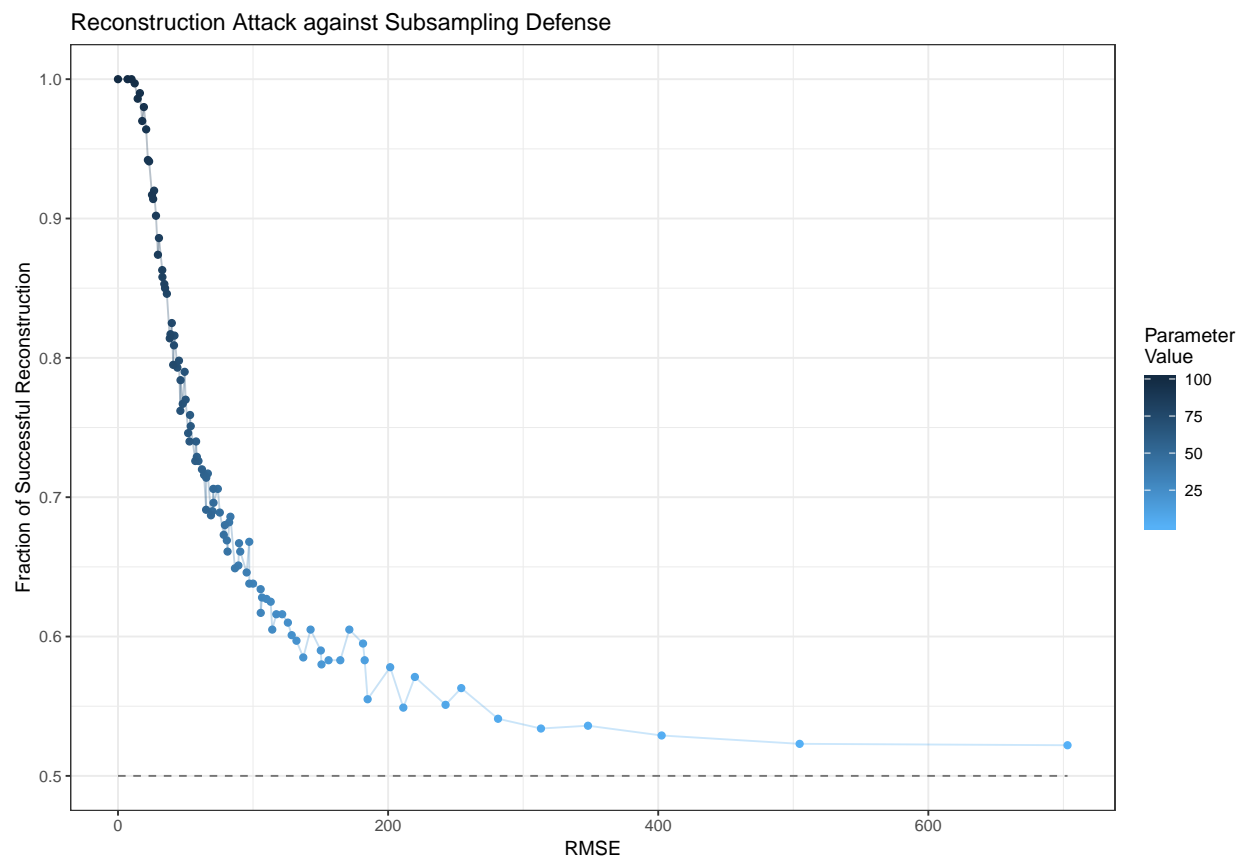
The code for this problem can be found in `q2reconstruction.R`. In this problem, a reconstruction attack was executed against three defense mechanisms (rounding, noise, subsampling) in a `R` script. Plots of the results follow, with the parameter values for each defense ranging from 1 to n (note that it is not really possible to round to the nearest multiple of 0, nor subsample 0 out of n rows). The resulting fraction of results that were successfully reconstructed were plotted against the root-mean-squared error (RMSE) metric. Each data point is the average of 10 trials in order to account for the randomness introduced by the defense mechanisms.



With the above attack against the rounding defense, after $R = 12$ success oscillates substantially before declining, suggesting this is roughly the rounding parameter threshold for success.



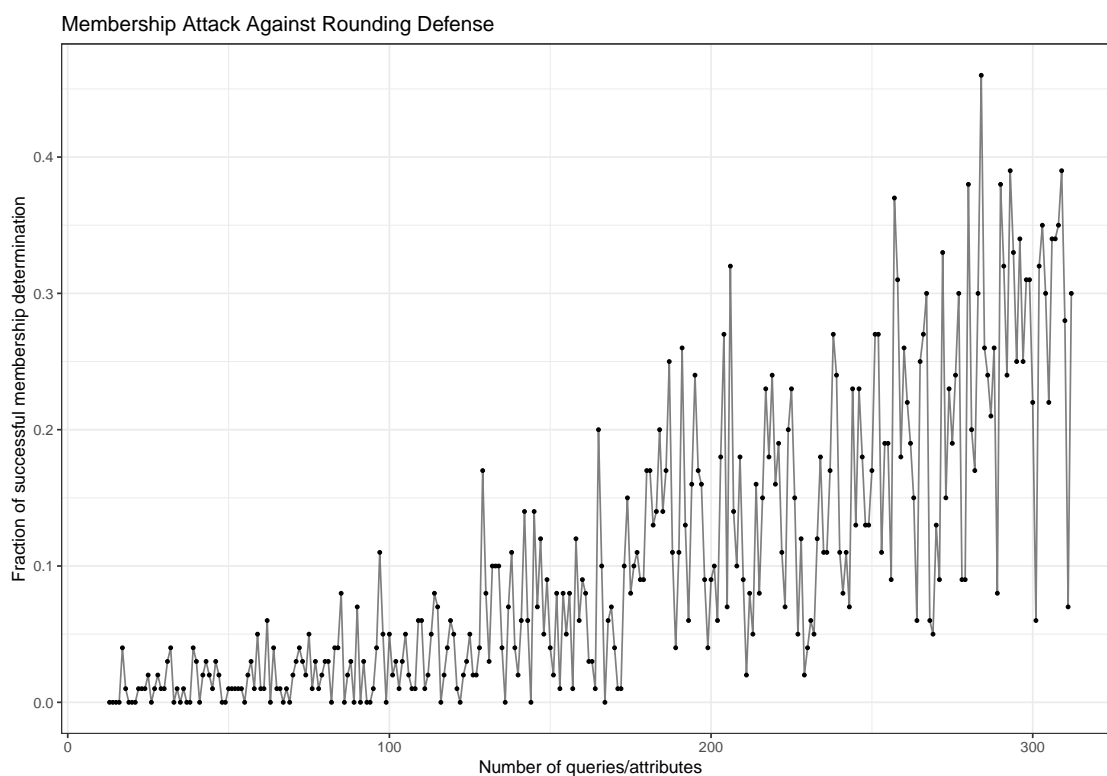
With the above attack against the noise injection defense, the results drop precipitously and approach failure with $\sigma \geq 14$ producing high enough noise to result in failure.

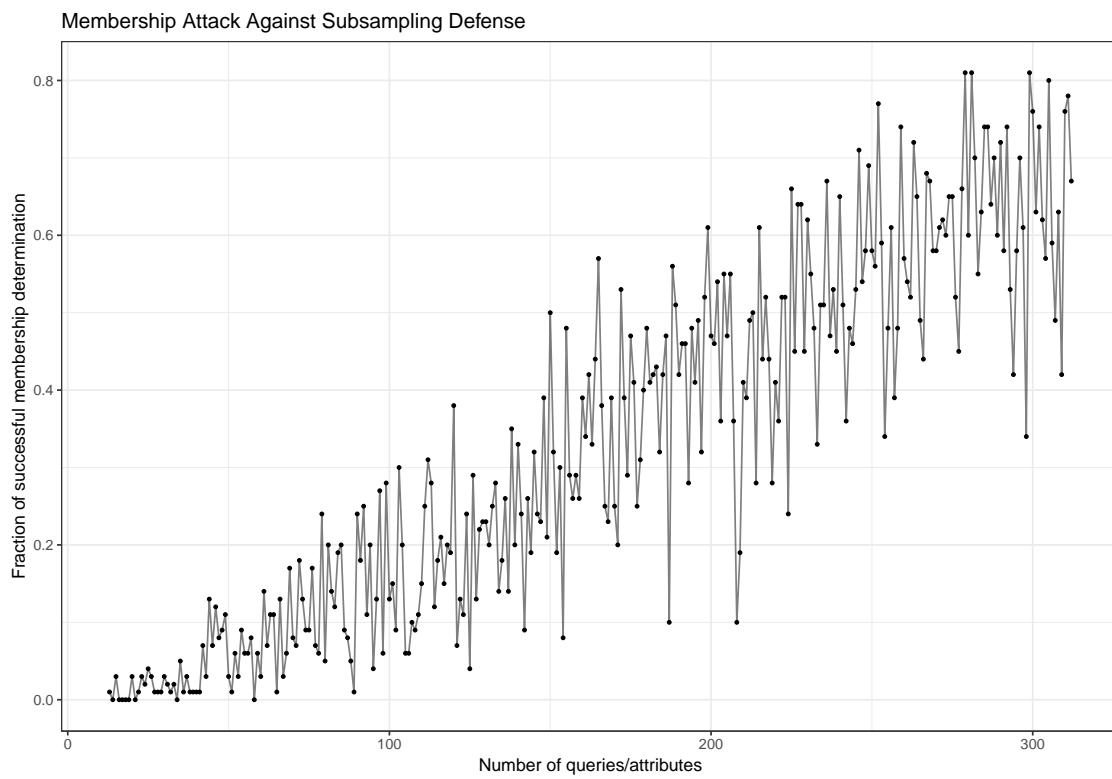
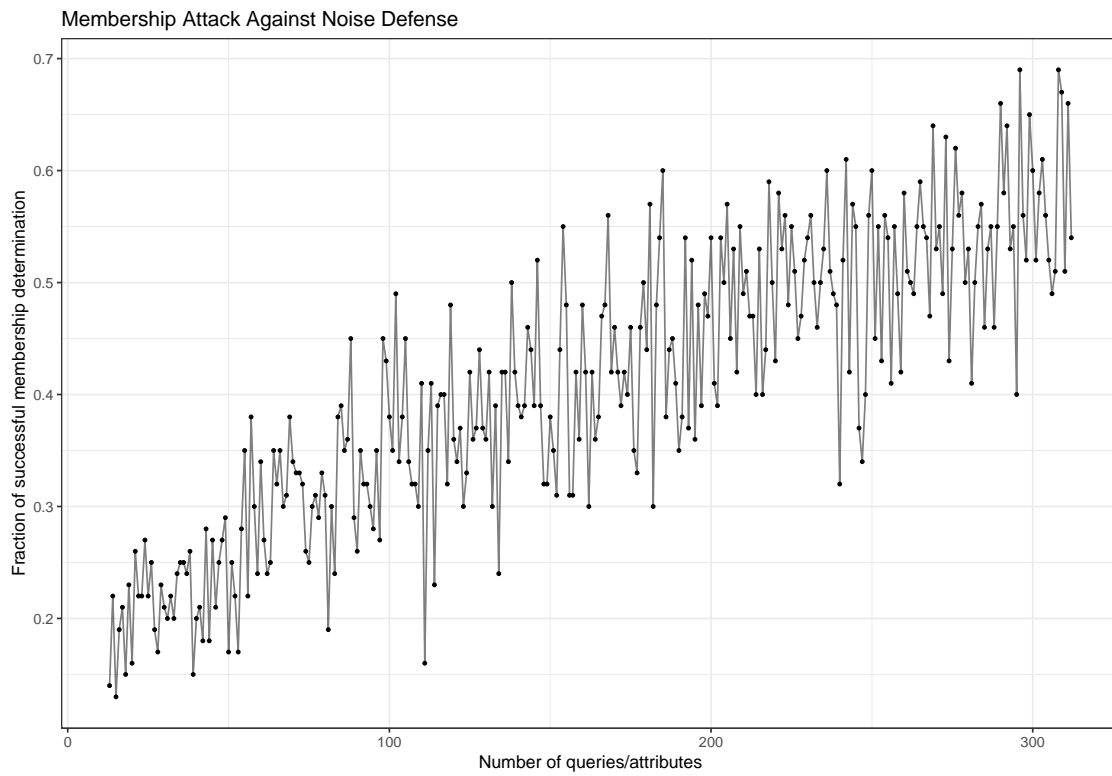


For the subsampling defense, larger subsamples result in worse privacy, with anything above subsamples of size around 13 generally providing sufficient protection against successful reconstruction.

Problem 3

The code for this problem can be found in `q3membership.R`. Plots of success fraction against number of queries executed have been included below for each of the three defense mechanisms, respectively. Note that the number of queries was limited to 300 for each as it provides a sufficiently clear illustration of convergence behavior, as well as for computational feasibility. The parameters for each defense were determined by the previous problem, corresponding to the reconstruction failure threshold. However, it can be seen that (with sufficiently large numbers of queries/attributes), membership attacks can still be performed with substantially high true positive probabilities.





Problem 4

Some tentative project topic ideas are listed and described as follow, roughly from most interest to least. Some of these are just starting points that require substantial research and reading to be a full project. The goal is to work on a theoretical project that can hopefully be parlayed down the road into a senior thesis or substantial research topic with potentially publishable results. However, some experimental projects like the case study may provide just as valuable insights and practice while being more feasible for the scope of this class.

- Cyberprivacy: Differential privacy foundations often focus on large datasets of personal information (i.e. census data), the ideal setting for evaluating and demonstrating privacy protection techniques. In practice, many pressing privacy concerns arise from the collection process and this proposed project idea focuses on this. This domain-specific application will partially be experimental, largely working with local differential privacy in the realm of cybersecurity. In the process, current implementations such as Apple's iOS differentially private collection mechanisms will likely be studied and discussed, with extensions to higher-profile situations such as social media and network traffic.
- Distributed DP: Many computational tasks and algorithms involve distributed computing. A decentralized differential privacy paradigm has already emerged with local DP. The privacy-utility tradeoff will be examined through this lense of centralized vs local DP. Distributed computing frameworks such as MapReduce and Spark will be evaluated for privacy, with considerations of how current local DP and DP algorithms in general can be extended to these distributed tasks.
- High-dimensional DP: High-dimensional problems are always a challenge for differential privacy given the inherently rich information that high-dimensional models often learn from the data. The challenge of creating private yet useful machine learning algorithms and, in particular, neural networks models will be studied, starting from reading the existing literature in the area before moving to experiments and possibly proposing new differentially private techniques in this area.
- Case Study: An experimental project involving a large, rich real-world dataset with the goal of evaluating privacy-utility tradeoff in this specific case. Various inference techniques and algorithms will be applied in the traditional manner assuming unfettered access to data, and again in a differentially private manner. One objective will be to determine a methodology towards selecting DP mechanisms and algorithms as to maximize utility with a given budget. Utility will also be compared against the non-differentially private results to understand the tradeoff being made in the name of privacy. Finally, attention will be drawn to privacy protection with attacks on the dataset using unprotected queries versus DP queries.
- Theory and Exposition: Given the somewhat limited experience in DP so far, yet the strong interest in in-depth theoretical projects, one possibility is the reading of relevant literature for a specific area of DP research and presenting an exposition of them. Some specific areas of interest include cryptography and developing new differentially-private algorithms, such as for time series models.

Appendix

Code for Problem 1

```
##
## q1.r
##
## Exploring Fulton PUMS data to identify variables for reidentification attack
##
## JH 2019/02/11
##

# read in the data
data <- read.csv('FultonPUMS5full.csv')
# summary of data
head(data)
nrow(data)
# variables
print(paste(names(data), collapse=", "))
print(length(names(data)))

# see how many unique values of each variable
vars <- c()
counts <- c()
for (n in names(data)) {
  vars <- c(vars, n)
  counts <- c(counts, nrow(unique(data[n])))
}
results <- data.frame(var=vars, count=counts)

# look at location identifiers
table(data$puma)
# unique(data$puma - data$jparow)
# unique(data$puma - data$X)
# # remove duplicates
# results <- results[!results$var %in% c("X", "jparow"),]

# remove variables with single count
results[results$count == 1,] # state and fips have one count
results <- results[results$count != 1,]

# plot distributions
jpeg('./figs/exploredist.jpg', width=1000, height=300)
par(mfrow=c(1,3))
barplot(table(data$income), xaxt='n', ylab='', main="income_Distribution")
```

```
axis(side=1)
barplot(table(data$age), ylab='', main="age_Distribution")
barplot(table(data$educ), ylab='', main="educ_Distribution")
par(mfrow=c(1,1))
dev.off()

# collision probabilities
collision <- function vardata) {
  sum((table(vardata) / nrow(data))^2)
}

p.collisions <- c()
for (v in as.character(results$var)) {
  p.collisions <- rbind(p.collisions, c(v, collision(data[v])))
}
p.collisions

# find average size of PUMA
avg.puma <- round(20 * nrow(data) / length(unique(data$puma)))

# calculate joint probability of collision
collision.joint <- function (vars) {
  joint <- 1
  for (v in vars) {
    joint <- joint * collision(data[v])
  }
  joint
}

# calculate percentage of average PUMA uniquely identifiable
identifiable <- function (p.collision) {
  (1 - p.collision)^avg.puma
}

# calculate results with proposed list of identifying variables
varlist <- c("income", "age", "educ")
identifiable(collision.joint(varlist))
varlist <- c(varlist, "sex", "married")
identifiable(collision.joint(varlist))
varlist <- c(varlist, "black")
identifiable(collision.joint(varlist))
varlist <- c(varlist, "employed")
identifiable(collision.joint(varlist))
```

Code for Problem 2

```
##  
## q2reconstructionattack.r  
##  
## Regression-based reconstruction attack against various defense mechanisms  
##  
## JH 2019/02/21  
##  
  
library(ggplot2)      # using ggplot graphing libraries (optional)  
  
##### Parameters  
  
set.seed(99)          # RNG seed  
n <- 100               # Dataset size  
k.trials <- 2 * n      # Number of queries  
exp.n <- 10            # Number of experiments at each step  
  
##### Data  
  
# read in the data  
sample.data <- read.csv('FultonPUMS5sample100.csv')  
sample.data.clean <- sample.data[3:17]  
attributes.n <- length(names(sample.data.clean))  
  
# get sensitive data (USCITIZEN)  
sensitive.var <- "uscitizen"  
sensitive.data <- sample.data[, sensitive.var]  
  
##### Querying Mechanisms (including defenses)  
  
# standard query of the dataset  
query <- function(data, pred, p) {  
  # extract data subset according to specified predicate  
  subset <- data[as.logical(pred)]  
  # calculate desired sum  
  sum <- sum(subset)  
  sum  
}  
  
# query w/ defense by rounding to nearest multiple of R  
query.rounding <- function(data, pred, r) {  
  # get actual query
```

```
    sum <- query(data, pred)
    # apply rounding defense
    sum.rounded <- round(sum / r) * r
    sum.rounded
  }

# query w/ defense by adding Gaussian noise
query.noisy <- function(data, pred, sigma) {
  # get actual query
  sum <- query(data, pred)
  # apply Gaussian noise as defense
  sum.noisy <- sum + rnorm(1, 0, sigma)
  sum.noisy
}

# query w/ defense by subsampling t out of n rows
query.subsampling <- function(data, pred, t) {
  # create subsample T with t out of n rows
  T.rows <- sample(length(data), t)
  # extract data subset from subsample according to specified predicate
  T.data <- data[T.rows]
  # get actual query on this subsample
  sum <- query(T.data, pred[T.rows])
  # scale up
  sum.sub <- sum * length(data) / t
  sum.sub
}

##### Define Predicates

# hashing to generate predicates
prime <- 563      # moderately large prime

predicate.single <- function(r.num, individual) {
  (sum(r.num * individual) %% prime) %% 2
}

# define predicates for the 2n queries
predicates <- matrix(NA, nrow = k.trials, ncol = n)

# generate 2n predicates
for (pred.index in 1:k.trials) {
  # generate random numbers for this hash
  r.num <- sample(0:(prime - 1), attributes.n, replace=TRUE)
```

```
# calculate for particular individual
pred.temp <- apply(sample.data.clean, MARGIN = 1, predicate.single, r.nums)
# store the predicate
predicates[pred.index,] <- pred.temp
}

##### Main Reconstruction Attack function
# takes a defense query mechanism and name of defense
# makes the attack and plots the results
reconstruction <- function(query.function, defense.name) {
  results <- data.frame(param=integer(), rmse=numeric(), frac=numeric())
  # loop through different parameter values
  for (param.value in 1:n) {
    # make 10 experiments each time, storing each RMSE and fraction of success
    experiments.rmse <- vector("numeric", exp.n)
    experiments.frac <- vector("numeric", exp.n)
    for (exp.index in 1:exp.n) {
      squared.errors <- c()
      history <- matrix(NA, nrow = k.trials, ncol = (n + 1))

##### Run Queries
# make the 2n queries using the generated predicates
for (pred.index in 1:k.trials) {
  # extract predicate
  pred.current <- predicates[pred.index,]
  # get standard query and query with defense
  q.standard <- query(sensitive.data, pred.current)
  q.defense <- query.function(sensitive.data, pred.current, param.value)
  # save this query (result w/ defense + predicate) to history table
  history[pred.index,] <- c(q.defense, as.numeric(pred.current))
  # calculate and store squared errors
  squared.errors <- c(squared.errors, (q.defense - q.standard) ** 2)
}
# calculate RMSE for this experiment
experiments.rmse[exp.index] <- sqrt(sum(squared.errors))

# convert into dataframe
xnames <- paste("x", 1:n, sep="")
varnames <- c("y", xnames)
release.data <- as.data.frame(history)
names(release.data) <- varnames

##### Regression Attack
# attack formula
```

```

    attack.formula <- paste(xnames, collapse=" + ")
    attack.formula <- paste("y ~ ", attack.formula, "-1")
    attack.formula <- as.formula(attack.formula)
    attack.output <- lm(attack.formula, data=release.data)
    # regression coefficient estimates
    attack.estimates <- attack.output$coef
    # estimates rounded to give reconstruction predictions
    attack.reconstructed <- as.numeric(attack.estimates > 0.5) #round(attack
    # calculate fraction successfully reconstructed
    experiments.frac[exp.index] <-
      sum(attack.reconstructed == sensitive.data) / n
  }
  # append the avg of the 10 experiments for this particular parameter value
  results <- rbind(results,
    c(param.value, mean(experiments.rmse), mean(experiments.
}

names(results) <- c("param", "rmse", "frac")

##### Visualization of results
# type to distinguish between successful (1) or not (0)
results$success <- as.numeric(results$frac > 0.5)
# plot the result
ggplot(results, aes(x=rmse, y=frac)) + # scatter plot
  geom_point(aes(color=param)) +
  # trend line
  geom_line(aes(color=param), alpha=0.3) +
  # success threshold
  geom_line(aes(y=0.5), alpha=0.5, size=0.5, linetype=2) +
  # labels and title
  ylab("Fraction of Successful Reconstruction") +
  xlab("RMSE") +
  ggtitle(paste("Reconstruction Attack against", defense.name)) +
  # legend formatting
  scale_color_continuous(name="Parameter\nValue", low="#56B1F7", high="#132
  theme_bw()
#plot(results$frac, results$rmse) # use this instead if no ggplot
}

# make the calls to the reconstruction attack function
reconstruction(query.rounding, "Rounding_Defense")
dev.copy2pdf(file="./figs/attackrounding.pdf")
reconstruction(query.noisy, "Noise_Defense")
dev.copy2pdf(file="./figs/attacknoise.pdf")
reconstruction(query.subsampling, "Subsampling_Defense")

```

```
dev.copy2pdf(file="./figs/attacksampling.pdf")
```


Code for Problem 3

```
##  
## q3membershipattack.r  
##  
## Membership attack  
##  
## JH 2019/02/23  
##  
  
library(ggplot2)      # using ggplot graphing libraries (optional)  
  
##### Parameters  
  
set.seed(99)          # RNG seed  
n <- 100              # Dataset size  
k.trials <- 300       # max number of additional attributes/predicates (m)  
n.samples <- 20  
n.sims <- 100  
delta <- 1 / (10 * n)  
# boundary parameters for reconstruction failure  
# obtained from reconstruction attack in q2 (see other script)  
param.rounding <- 13  
param.noisy <- 14  
param.subsampling <- 13  
  
##### Data  
  
# read in the data  
sample.data <- read.csv('FultonPUMS5sample100.csv')  
sample.data.clean <- sample.data[3:17]  
attributes.n <- length(names(sample.data.clean))  
  
sample.data.binary <- sample.data.clean[c(1, 5:15)]  
binary.names <- names(sample.data.binary)  
  
##### Querying Mechanisms (including defenses)  
  
# standard query of the dataset  
query <- function(data, p) {  
  # calculate the desired sum  
  sum(data) / length(data)  
}  
  
# query w/ defense by rounding to nearest multiple of R
```

```
query.rounding <- function(data, r) {  
  # apply rounding defense  
  sum.rounded <- round(sum(data) / r) * r  
  sum.rounded / length(data)  
}  
  
# query w/ defense by adding Gaussian noise  
query.noisy <- function(data, sigma) {  
  # apply Gaussian noise as defense  
  sum.noisy <- sum(data) + rnorm(1, 0, sigma)  
  sum.noisy / length(data)  
}  
  
# query w/ defense by subsampling t out of n rows  
query.subsampling <- function(data, t) {  
  # create subsample T with t out of n rows  
  T.rows <- sample(length(data), t)  
  # extract data subset from subsample according to specified predicate  
  T.data <- data[T.rows]  
  # get actual query on this subsample  
  sum <- sum(T.data)  
  # scale up  
  sum.sub <- sum * length(data) / t  
  sum.sub / length(data)  
}  
  
##### Predicates  
  
# hashing to generate predicates  
prime <- 563      # moderately large prime  
  
predicate.single <- function(individual, r.nums) {  
  (sum(r.nums * individual) %% prime) %% 2  
}  
  
# define predicates for the queries  
predicates <- matrix(NA, nrow = k.trials, ncol = n)  
  
# helper function to create data from population  
rmvbernoulli <- function(n=1, prob){  
  history <- matrix(NA, nrow=n, ncol=length(prob))  
  for(i in 1:n){  
    x<- rbinom(n=length(prob), size=1, prob=prob)
```

```

    x[x==0] <- -1
    history[i,] <- x
  }
  history
}

# test statistic (Dwork et al.)
test.dwork <- function(alice, sample.mean, population.mean){
  sum(alice * sample.mean) - sum(population.mean * sample.mean)
}

# Generate the critical value
generate.critical <- function(null.sims=1000, alpha, fun, population.prob){
  population.mean <- 2 * (population.prob - 0.5)
  hold <- rep(NA, null.sims)
  for(i in 1:null.sims){
    sample <- rmvbernoulli(n=n.samples, prob=population.prob)
    null.alice <- rmvbernoulli(n=1, prob=population.prob)
    sample.mean <- apply(sample, MARGIN=2, FUN=mean)
    hold[i] <- eval(fun(alice=null.alice, sample.mean=sample.mean, population
  )
  null.distribution <- sort(hold, decreasing=TRUE)
  null.distribution[round(alpha * null.sims)]
}

population.data <- read.csv('FultonPUMS5full.csv')
population.probs <- as.vector(apply(population.data[binary.names], 2, mean))
population.means <- 2 * (population.probs - 0.5)

# run simulation against provided query defense mechanism
simulate <- function(query.function, query.param) {
  predicates <- t(sample.data.binary)

  population.probs.all <- population.probs
  population.means.all <- population.means

  results <- matrix(NA, nrow = k.trials, ncol = 2)

  # loop through number of predicates/attributes
  for(pred.index in 1:k.trials) {
    # generate random numbers for this hash
    r.num <- sample(0:(prime - 1), attributes.n, replace=TRUE)
    # calculate for particular individual
    pred.temp <- apply(sample.data.clean, MARGIN = 1, predicate.single, r.num

```

```

# store the predicate
predicates <- rbind(predicates, pred.temp)
population.probs.all <- c(population.probs.all, 0.5)
population.means.all <- c(population.means.all, 0)

# critical value of the null distribution
null.critical <- generate.critical(alpha=delta, fun=test.dwork, populatio

history <- vector("numeric", n.sims)
for(i in 1:n.sims) {
  # Simulate data
  sample <- rmvbernoulli(n=n.samples, prob=population.probs.all)
  sample.mean <- apply(sample, MARGIN=2, FUN=query.function, query.param)
  alice <- sample[1,]

  # Conduct hypothesis test
  history[i] <- test.dwork(alice=alice, sample.mean=sample.mean, populatio
}

# fraction identified IN
accept.frac <- round(100 * mean(history > null.critical)) / 100

# storing results
results[pred.index, 1] <- nrow(predicates)
results[pred.index, 2] <- accept.frac

# helps visualize while generating more attributes
plot(results[, 1], results[, 2], type='l')
}
results
}

final.graph <- function(results, defense) {
  # final visualization with ggplot
  ggplot(data.frame(results), aes(x=X1, y=X2)) +
    geom_point(size=0.75) +
    geom_line(alpha=0.5) +
    xlab("Number_of_queries/attributes") +
    ylab("Fraction_of_successful_membership_determination") +
    ggtitle(paste("Membership_Attack_Against", defense)) +
    theme_bw()
}

```

```
# make simulations and visualizations
simulate.standard <- simulate(query, 0)
final.graph(simulate.standard, "Standard_Query_w/_No_Defense")
dev.copy2pdf(file="./figs/membership.pdf")

simulate.rounding <- simulate(query.rounding, param.rounding)
final.graph(simulate.rounding, "Rounding_Defense")
dev.copy2pdf(file="./figs/membershiprounding.pdf")

simulate.noisy <- simulate(query.noisy, param.noisy)
final.graph(simulate.noisy, "Noise_Defense")
dev.copy2pdf(file="./figs/membershipnoisy.pdf")

simulate.subsampling <- simulate(query.subsampling, param.subsampling)
final.graph(simulate.subsampling, "Subsampling_Defense")
dev.copy2pdf(file="./figs/membershipsubsampling.pdf")
```