# Introducing Very Large Data Sets into the Classroom

Dason Kurkiewicz

Department of Statistics, Iowa State University

March 9, 2012

# Outline

- Dason's Introduction
- The GUI
- The Usability Study
- Future Work
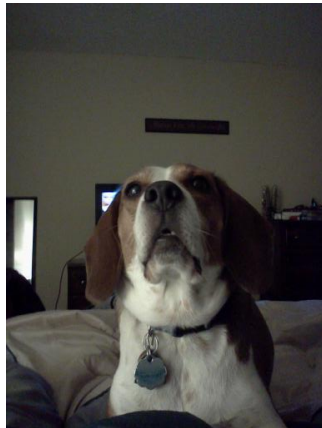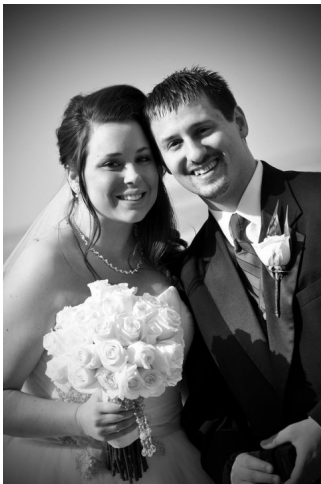- Conclusions

# Outline

# Outline

- Dason's Introduction

- The GUI

- The Usability Study

- Future Work

- Conclusions

# Outline

- Dason's Introduction
- The GUI
- The Usability Study
- Future Work
- Conclusions

# Outline

- Dason's Introduction
- The GUI
- The Usability Study
- Future Work
- Conclusions

# Outline

- Dason's Introduction
- The GUI
- The Usability Study
- Future Work
- Conclusions

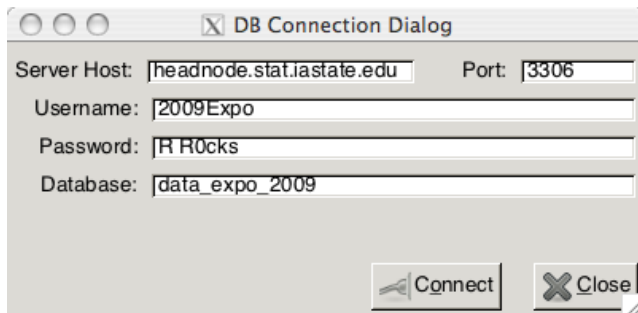# Introduction
Meet Dason



My Family

# dbConnectGUI
Installation

```
# install.packages("gWidgetsRGtk2")
# install.packages("DBI")
# install.packages("RMySQL")

install.packages("dbConnectGUI")
library(dbConnectGUI)
# Try it out
dbConnectGUI()
# Or
connect <- getConnection()
dbConnectGUI(connect$con)
```

# dbConnectGUI
Data Viewer

# dbConnectGUI
Variable Browser

# dbConnectGUI

Subsample

# dbConnectGUI

Query

# dbConnectGUI
Changes

## Learning Objectives

- Familiarizing students with databases and exploring the utility of databases for effcient data storage.
- Obtaining access to data stored in databases.
- Searching databases for specifc data records of interest as well as learning how to extract such data records for further statistical analysis in R or a different software package.
- Providing students with a first gentle and more paced introduction to the SQL language. The GUI can help facilitate the learning process and potentially improve student's attitude toward learning the database language SQL.

|                | 400-level  |                | 500-level  |                |
| -------------- | ---------- | -------------- | ---------- | -------------- |
|                | Statistics | non-Statistics | Statistics | non-Statistics |
| Undergraduate  | 23         | 6              | 1          | 0              |
| Graduate       | 1          | 8              | 21         | 15             |

Table: Overview of students participating in the usability study by course, status, and area.

# Usability Study

# Usability Study



Figure: User feedback: perceived importance of task (left) and usefulness of the GUI to help with the task (right). Undergraduates are red and Graduates are Blue.

# Usability Study
Students' opinions on task difficulty



Figure: User feedback: anticipated degree of difficulty is greater than experienced difficulty of the task. Undergraduates are Red and Graduates are Blue.

# Future Work

- Efficient Random Stratified Samples
- Joins
- Tutorial

# Future Work

- Efficient Random Stratified Samples
- Joins
- Tutorial

# Future Work

- Efficient Random Stratified Samples
- Joins
- Tutorial

# Future Work

- Efficient Random Stratified Samples
- Joins
- Tutorial

# Closing Remarks

- Big data is here to stay

- GUIs provide a level of comfort to some

- The usability study shows there is benefit to using the GUI

- Will keep working to make this better

# Closing Remarks

- Big data is here to stay
- GUIs provide a level of comfort to some
- The usability study shows there is benefit to using the GUI
- Will keep working to make this better

# Closing Remarks

- Big data is here to stay
- GUIs provide a level of comfort to some
- The usability study shows there is benefit to using the GUI
- Will keep working to make this better

# Closing Remarks

- Big data is here to stay
- GUIs provide a level of comfort to some
- The usability study shows there is benefit to using the GUI
- Will keep working to make this better

# Closing Remarks

- Big data is here to stay
- GUIs provide a level of comfort to some
- The usability study shows there is benefit to using the GUI
- Will keep working to make this better

Questions?