# STRATEGIES FOR DATA ANALYSIS AND VISUALIZATION

For Increased Collaboration, Openness and Sharing

# This is Not A Talk About How To Analyze and Visualize your Data

- Your probably already better at that than me
- There are too many ways to analyze data
  - Project Specific
  - Domain Specific
- It would be really boring

# This Talk Is About

- ☐ Treating your analysis as a first class data object
- ☐ Maximizing your efficiency creating analyses and visualizations
- ☐ Increasing your ability to use your analyses and visualizations with collaborators
- ☐ Making your analyses and visualizations more open.

# General Outline

# Some Background

## The Changing Regulatory Landscape

# NSF Guidelines

- a. Investigators are expected to promptly prepare and submit for publication, with authorship that accurately reflects the contributions of those involved, all significant findings from work conducted under NSF grants. Grantees are expected to permit and encourage such publication by those actually performing that work, unless a grantee intends to publish or disseminate such findings itself.

# NSF Guidelines



- b. Investigators are expected to share with other researchers, at no more than incremental cost and within a reasonable time, the primary data, samples, physical collections and other supporting materials created or gathered in the course of work under NSF grants. Grantees are expected to encourage and facilitate such sharing. Privileged or confidential information should be released only in a form that protects the privacy of individuals and subjects involved.
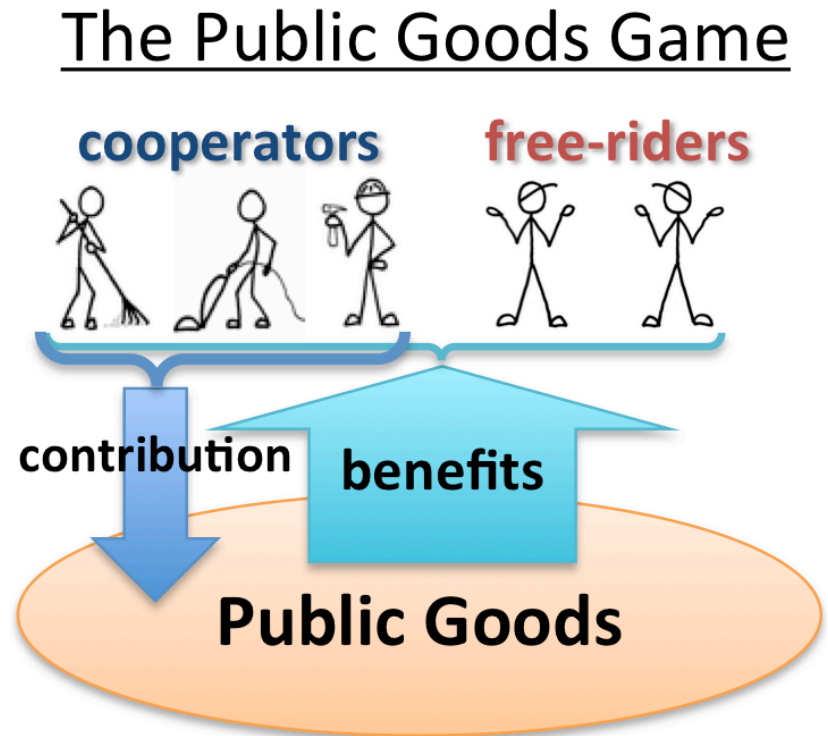
# NSF Guidelines

- c. Investigators and grantees are encouraged to share software and inventions created under the grant or otherwise make them or their products widely available and usable.

# NSF Guidelines

- d. NSF normally allows grantees to retain principal legal rights to intellectual property developed under NSF grants to provide incentives for development and dissemination of inventions, software and publications that can enhance their usefulness, accessibility and upkeep. Such incentives do not, however, reduce the responsibility that investigators and organizations have as members of the scientific and engineering community, to make results, data and collections available to other researchers.



The Public Goods Game

cooperators    free-riders

contribution    benefits

Public Goods

https://plektix.fieldofscience.com/2011/04/freedom-and-public-goods.html

# The Take Home Message

- ☐ Your grant funded work is a public good.
- ☐ Sharing the products of research benefits society.
- ☐ The works is still recognized as your work.
- ☐ But you still have to share.

# Why Do Funding Agencies Care?

- Good PR.

- Increased accountability.

- Benefits to the research community.

What's in it for me?

https://commons.wikimedia.org/wiki/File:Mad_scientist_transparent_background.svg

# Why You Should Care





- ☐ Increase your efficiency

- ☐ Facilitate collaboration.

- ☐ Maximize the usefulness of your data.

- ☐ Lack of Future Funding

- ☐ ???

# I thought it was just the data!

Nope. Your other products are first class objects as well.

# Data, Collection, Analyses and Visualizations are Inseparable

- The Analysis Visualizations Help Explain Your Data

- What your data is for.

- An example of how your data can be used.

# Your Analyses and Visualizations Are Important On Their Own

- Reuse
  - Increase your efficiency
  - Solve other peoples problems

*Good programmers write good code.*
*Great programmers steal great code.*
         *-unknown-*

# Functional Requirements

Things your analysis and visualization should do beyond analysis and visualization

# Analyze/Visualize Your Data

# Understandable



Gary Larson

- You should understand what you have done a year later
- Others should understand what you did without your direct explanation
- Don't assume we are all geniuses

# Acurate

- Your analyses should render the same results as those you reported
  - Every time

- This includes randomizations.

# Be Durable



- □ Platform independent

- □ Easy to migrate

- □ Easy to translate to other software

- □ Easy to see

# Best Practices

Fairly Easy Things To Improve Your Analyses and Visualizations
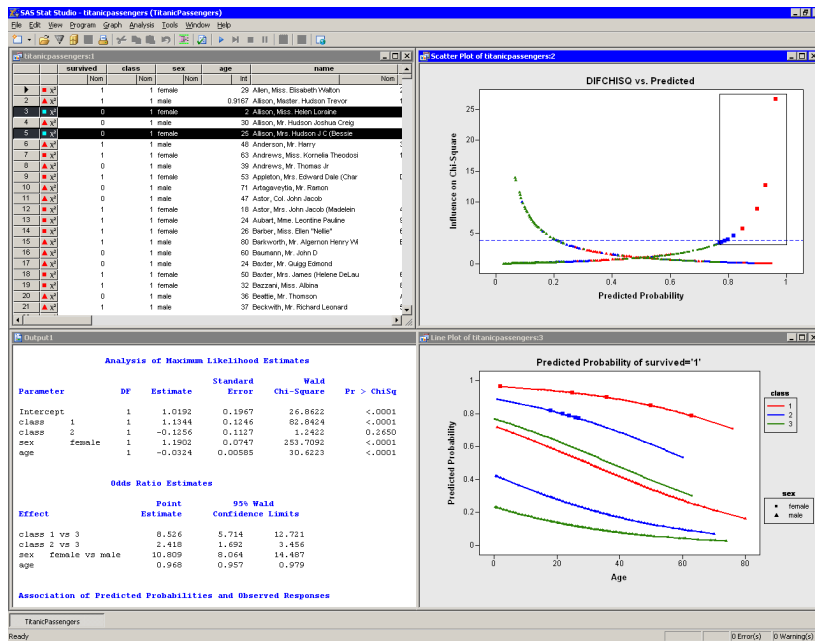
# Favor Text Over GUI

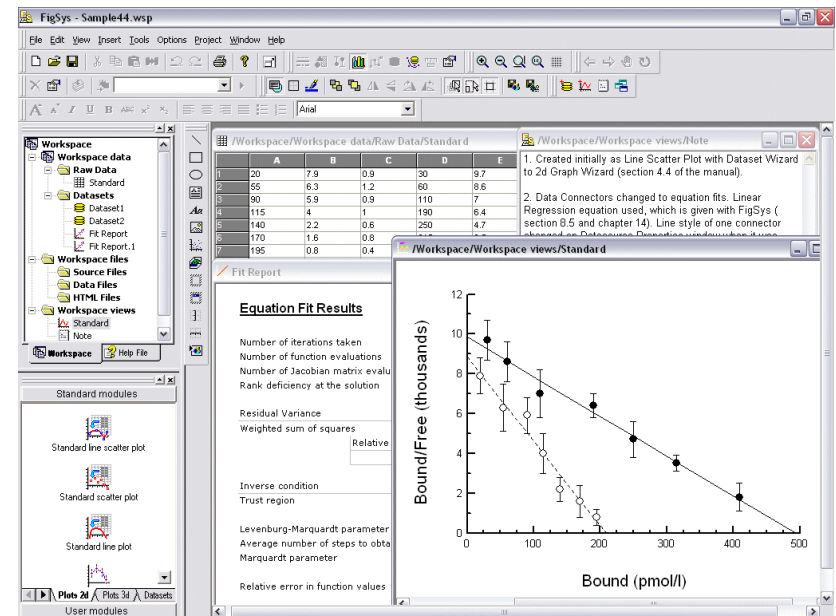| GUI | Text |
|---|---|
| ☐ Easy to use | ☐ Harder to use |
| ☐ Software Dependent | ☐ Less software dependent |
| ☐ Difficult to share | ☐ Easier to share |
| ☐ Difficult to modify | ☐ Easier to modify |
| ☐ Difficult to version | ☐ Easy to version |

# The Problem With GUIs



- Too easy to hide assumptions.

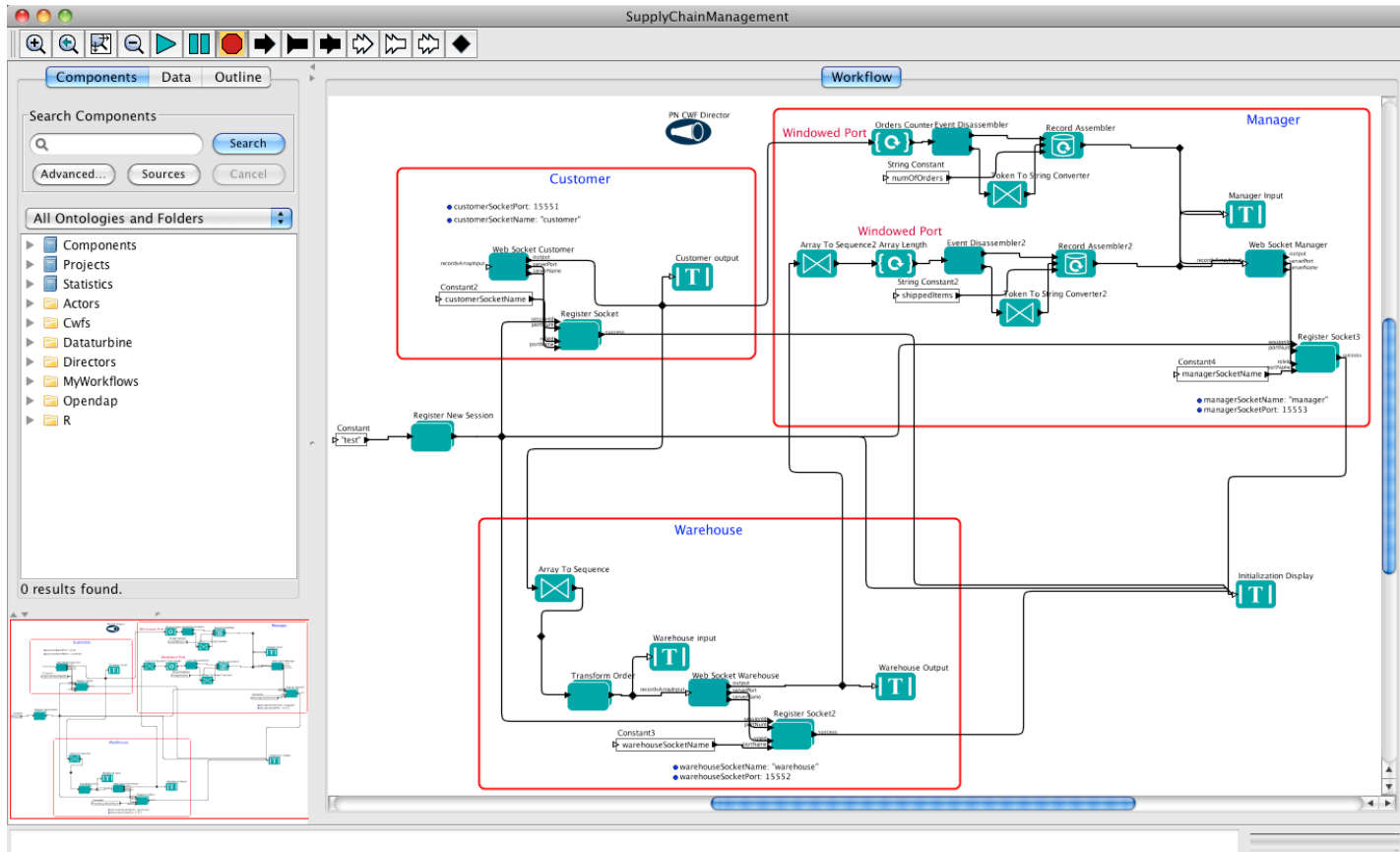- Difficult to document what you have done.

- Difficult to share.

# If You Use a GUI

- If you can, download the code.

- Document all settings.

- Save the file, open and rerun.



- Screenshots.
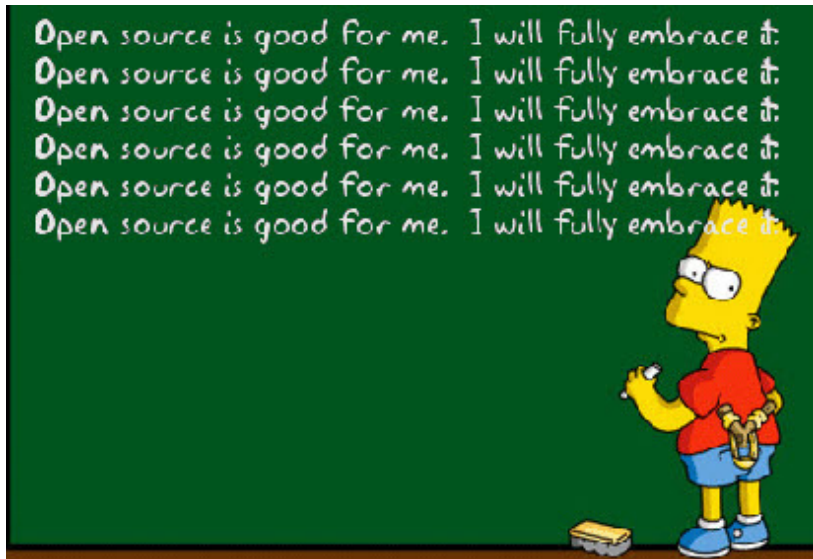
# Scientific Workflow Systems

# Analysis and Visualization as Software

Borrowing from AGILE development and other current coding best practices.

# Reuse Code and Share Code



- Speeds development

- Reduces errors

- Increases transparency

- Promotes Collaboration

# Test Frequently

- Write tests for your code and make them available

- Try to test each line of code


BUG FREE
**Guaranteed**

# Clean Code

- Be very organized
- Use descriptive variable and function names
- Keep code blocks small
- A non-program should be able to read and understand it.
- Comment but don't over comment your code.



Hey, Carl, can you look at this problem with me. I've been working on this for hours. You see the X variable clearly cannot be less than zero because Y has to be more than 20.... Oh wait. That's not right. OK, I've got it now. Thanks, Carl!

# Really Bad

```
function fetch(i) { int j=get_it(i);
  if(j>1)
  { do_something_important(i); } else
  { do_something_important(j); } }
  function go(i)
  { calculate_something(i); sleep(i);
  think_about_something(i); } function
  go_fetch(id) { int i=id; i++; go(i);
  fetch(i): return i; } function main()
  { int i=0; i=get_i(); go_fetch(i); }
```

# Getting Better

```
function fetch(i) {
    int j=get_it(i);
    if(j>1) {
        do_something_important(i);
    }
    else {
        do_something_important(j);
    }
}

function go(i) {
    calculate_something(i);
    sleep(i);
    think_about_something(i);
}
```

```
function go_fetch(id) {
    int i=id;
    i++;
    go(i);
    fetch(i);
    return i;
}

function main() {
    int i=0;
    i=get_i();
    go_fetch(i);
}
```
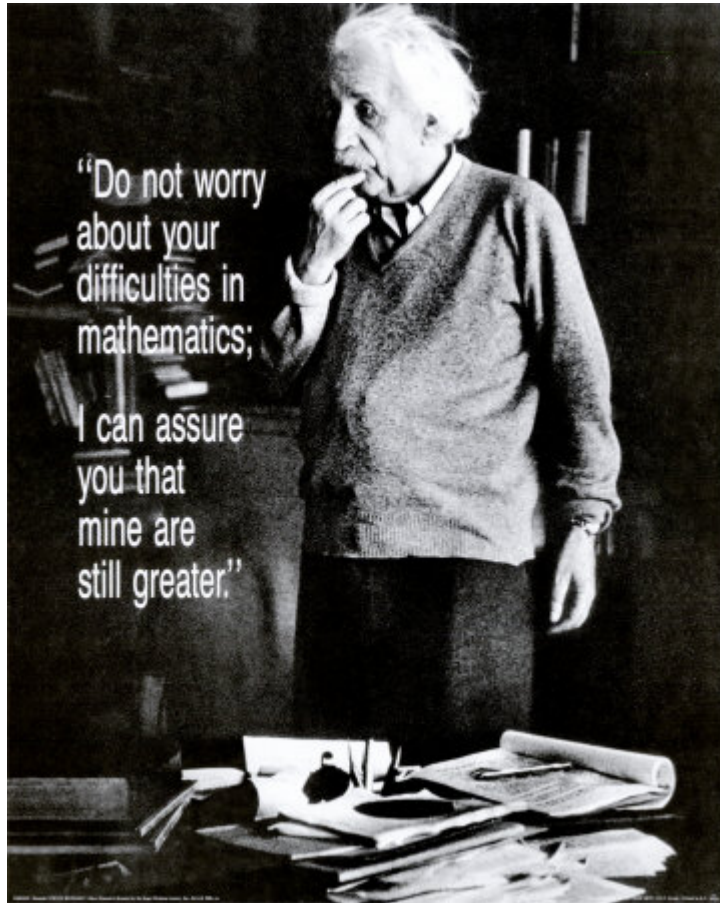
# Even Better

```
function main() {
    int i=0,j;
    i=get_i();
    i++;
    calculate_something(i);
    sleep(i);
    think_about_something(i);
    j=get_it(i);
    if(j>1) {

        do_something_important(i); }
    else {
        do_something_important(j);
    }
}
```

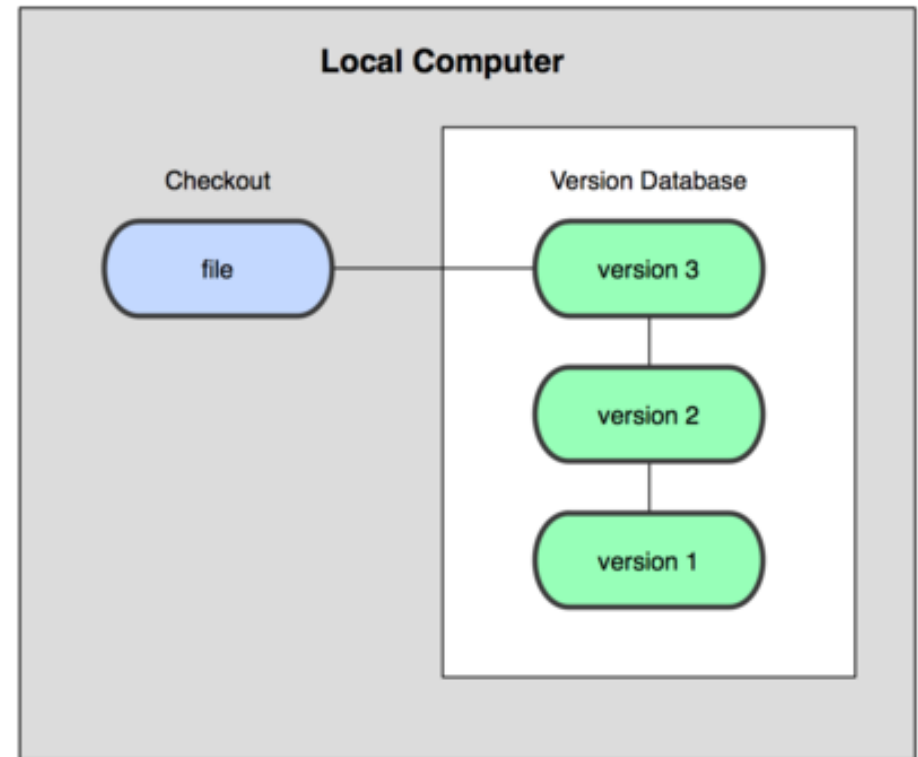- The variables are still poorly named.

# Why Worry About The Code?



"Do not worry about your difficulties in mathematics; I can assure you that mine are still greater."

- You will understand it later.
- Others will understand it.
- Increased reuse
- Decreased errors and bugs.

# Versioning

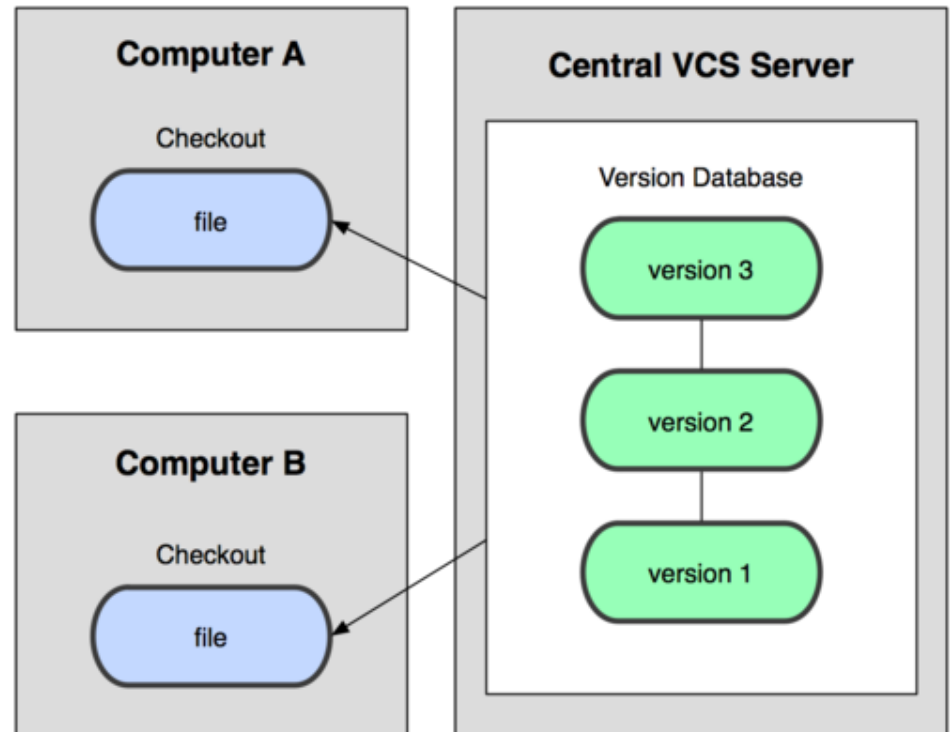- Developed for versioning code.

- You can version anything.

- Works best on text based objects.



**Local Computer**

Checkout — file

Version Database — version 3 — version 2 — version 1

# Versioning

- Allows you to roll back changes.

- Collaborate on files and merge changes.

- Branch off modifications.

# A Data Specific Use Case

**Version 1**
- Pub 1
- Pub 2

**Version 2**
- Pub 3
- Pub 4

# Dissemination

NSF Requires You To Be Open Source

# Where To Publish Your Analyses

- ☐ Open Source Repositories

- ☐ Statistical Repositories

- ☐ Institutional Repositories

- ☐ Personal Web Sites

# How To Publish

- ☐ Provide Adequate and Accurate Metadata

- ☐ Provide a License

- ☐ Use Social Networking

- ☐ Provide Links In Publications

# Licensing – Most Common Licenses

☐ GNU General Public License

☐ Apache Software License

☐ Creative Commons
  ▪ The most adaptable.

# Final Thoughts

- You are required to make them available

- You should make them available

- Reuse and Share

- Treat your analyses like data

- License your work

# Questions?

- The slides are available at https://github.com/olendorf/presentations/tree/master/ciday-2012