# Big Data: Data Wrangling Boot Camp
# Assumed File Structure

## Chuck Cartledge, PhD

## 23 February 2018

# Table of contents (1 of 1)

What are we going to cover?

Most of the R scripts used in this boot
camp assume a certain directory
structure, and specific locations for
scripts, data files, and images. To wit:

1. Scripts directory – where R scripts
   "live"
2. Data directory – where data
   usually comes from and goes to

If images are created, they are
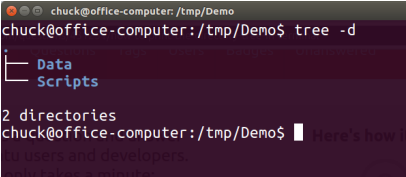considered a type of data.

## Relationship of directories.

There are two directories, relative
to one another, and anywhere in
the file system.

- Data – where data files are
  usually read from or written
  to
- Scripts – where R scripts are
  executed and "sourced"
  from

The file and directory names are
Unix case sensitive and should be
safe in a Windows environment.

# Where programs execute.

- "Base" R scripts may `source()` other script files. All files are assumed to live in the scripts directory.

- Access to data files within the R scripts is via the `file.path()` function.



All file accesses should be Operating System agnostic.

## Sometimes, things need to be persistent.

- Data files live, and die in the data directory.
- The Data directory is one "up" from the Scripts directory.
- All accesses to the data directory are via the `file.path()` function.



All file accesses should be Operating System agnostic.

Where temporary things live.

- There are times when data does not need to be persistent.
- Generally these data are stored wherever `tempfile()` or `tempdir()` put them.
- Be aware that data stored, either directly or indirectly using these functions may be removed when the R session completes.

All file accesses should be Operating System agnostic.