

# Cataloging ~ G-Drive

## Contents

<b>1</b>	<b>Introduction</b>	<b>1</b>
1.1	Opening the files: . . . . .	2
<b>2</b>	<b>The G-Drive Contents</b>	<b>2</b>
2.1	Block Level analysis . . . . .	2
2.2	Census . . . . .	2
2.3	The ramosRivera Folder . . . . .	2
2.3.1	bg06_d00_shp Folder . . . . .	2
2.3.2	DART_IRL Scans Folder . . . . .	3
2.3.3	Summer2021_Dylan Folder . . . . .	3
2.3.4	Summer2021_Kelly Folder . . . . .	3
2.3.4.1	Task From Antonio 1 folder . . . . .	3
2.3.4.1.1	Original Material Folder . . . . .	3
2.3.4.2	Task From Antonio 2 Folder . . . . .	3
2.3.4.3	Task From Antonio 3 Folder . . . . .	4
2.3.4.4	Not in a folder . . . . .	4
2.3.5	T-Burk Folder . . . . .	4
2.3.5.1	ArcMap Folder . . . . .	4
2.3.5.2	BlockData Folder . . . . .	4
2.3.5.2.1	nhgis0003_shapefiles_tl2000_560_block_2000 . . . . .	4
2.3.5.2.2	Not in a folder . . . . .	6
2.3.5.3	k-function_local_results Folder . . . . .	7
2.3.5.4	LabData Folder . . . . .	8
2.3.5.5	PatentData . . . . .	8
2.3.5.6	PngData Folder . . . . .	9
2.3.5.7	Python Scripts Folder . . . . .	10
2.3.5.8	Tables folder . . . . .	11
2.3.5.9	ZipData folder . . . . .	11
2.3.5.9.1	tl_2010_us_ZCTA500 Folder . . . . .	11
2.3.5.9.2	Not in folder . . . . .	11
2.3.5.10	nhgis0005_shape Folder . . . . .	12
2.3.5.10.1	nhgis0004_shapefile_tl2000_330_block_2000 (zipped Folder) . . .	12
2.3.5.10.2	nhgis0004_shapefile_tl2010_110_block_2000 (zipped Folder) . . .	12
2.3.5.10.3	nhgis0004_shpaefile_tl2010_250_block_2000 (zipped Folder) . . .	12
2.3.5.11	nhgis0005_csv Folder . . . . .	12
2.3.6	tl_2010_06_zcta500 Folder . . . . .	12
2.3.7	~ Not in a Folder ~ . . . . .	12
2.4	ramosRivera - Backup062122 Folder . . . . .	13
2.5	Deleted or not in the Admin Folder . . . . .	13

## 1 Introduction

This Document catalogs all the files related to the labs— heterogeneity project.

## 1.1 Opening the files:

File type	How to Open
<b>.shp</b>	open ArcGIS, open a blank map, add a layer, and open the file in the program. (These files might also show as adobe acrobat files in the G-drive.)
<b>.cpg, .prj, .shx, .sbn, .sbx</b>	helper files to the shape file, cannot be opened independently.
<b>.xml</b>	Helper file to shape files, can open using excel.
<b>.dbf</b>	Helper file to shape files, can open using excel.
<b>.py</b>	python scripts, open using Spyder (Anaconda3)
<b>.pdf</b>	adobe Acrobat
<b>.csv</b>	excel
<b>.sas</b>	Stata
<b>.sas7bdat</b>	import into Stata
<b>.gz</b>	Use this link for more information on how to open in Python. (These zipped files support a python script, will not need to open)

## 2 The G-Drive Contents

For each folder in Labs-kbuzard-S18 that has files that are used in the heterogeneity project, there is a separate sub-heading below.

Copy and Paste this short cut to find the folder on <https://rds.syr.edu/rdweb/webclient/> in the search bar.

“G:/MAX-Fileer/Collab/Labs-kbuzard-S18/Admin” (you may have to change to backslashes)

### 2.1 Block Level analysis

- **CA\_Block\_Data.shp** - inputs for countSim\_speedUp, CountSim\_tester, and multiprocessing\_test2
- **CA\_ZCTA\_Data.shp** - inputs for countSim\_speedUp, CountSim\_tester, and multiprocessing\_test2
- **CA\_Labs\_Data.shp** - inputs for countSim\_speedUp, CountSim\_tester, and multiprocessing\_test2

### 2.2 Census

- **1998DART32.pdf** - input for pdf2Jpg.py
- **Labs1998.csv** - input for Geo\_coder, pngwork, and prep\_Labs

### 2.3 The ramosRivera Folder

#### 2.3.1 bg06\_d00\_shp Folder

has three documents - all of them make up the map of California used in paper.

- **bd06\_d00 (.shp)** - shows the map of California broken down by zipcodes
- **bg06\_d00 (.dbf, .shx)** - 9 observation: Area, Perimeter, BG06\_D00, BG06\_D00\_I, State, County, Tract, BLKGROUP, and NAME. These make up the information needed to recreate the California map (opened in excel)

### 2.3.2 DART\_IRL Scans Folder

Has two Pdfs copies of information on R&D labs and their location

- **1979IRL16** - Industrial Research Laboratories of the US, 16th Edition 1979
  - Original document containing information on the Industrial Research Laboratories of the US. Has information on 9,907 R&D facilities belonging to 6,323 organizations in 1979.
- **1989DART23** - Directory of American Research and Technology 1989, 23rd Edition
  - Original document containing information on organizations active in product development for business in American. Content includes information on 11,275 organizations in alphabetical order.

### 2.3.3 Summer2021\_Dylan Folder

Has eight documents all pertaining to the 1979 and 1989 data from IRL and DART pdfs

- **1979\_Digitized.txt** - digitized version of the 1979 IRL pdf
- **1979IRL16.pdf** - copy of the 1979IRL16 pdf in DART\_IRL Scans folder.
- **1989\_Digitized.txt** - digitized version of the 1989 DART pdf
- **1989DART23.pdf** - copy of the 1989DART23 pdf in DART\_IRL Scans folder. R
- **corr\_cattLabs97\_Wgeocode 1-6200.cvs** - excel file with lines 1-6200 corrected by Dylan. excel file with lines 1-6200 corrected by Dylan.
- **corr\_cattLabs97\_Wgeocode.cvs** - original excel file before Dylan and Kelly worked on it.
- **corr\_cattLabs97\_Wgeocode\_Line 6200 to Line 12765.cvs** - excel file with lines 6200-12765 corrected by Kelly.
- **OCR\_Result\_NO\_user.txt** - Antonio's initial OCR scan.
  - This data was input into the corr\_cattLabs97\_Wgeocode excel sheets

### 2.3.4 Summer2021\_Kelly Folder

#### 2.3.4.1 Task From Antonio 1 folder

##### 2.3.4.1.1 Original Material Folder

- **corr\_cattLabs97\_Wgeocode.cvs** - original excel file before Dylan and Kelly worked on it.
- **letter\_I\_cattell.txt** - digitized version of the research labs starting with the letter "I."
- **letter\_O\_cattell.txt** - digitized version of the research labs starting with the letter "O."
- **letter\_S\_cattell.txt** - digitized version of the research labs starting with the letter "S."  
this separation by letter section was done to make digitization process faster.
- **OCR\_Result\_NO\_user.txt** - Antonio's initial OCR scan.
  - This data was input into the corr\_cattLabs97\_Wgeocode excel sheets
- **corr\_cattLabs97\_Wgeocode\_Line 6200 to Line 12765.cvs** - excel file with lines 6200-12765 corrected by Kelly. Refer to Original Admin/ramosRivera/T-Burk/PngData/corr\_cattLabs97\_Wgeocode

#### 2.3.4.2 Task From Antonio 2 Folder

- **1979\_Digitized.pdf** - digitized version of the 1979 IRL pdf. For original refer to

#### 2.3.4.3 Task From Antonio 3 Folder

- **1989\_Digitized.txt** - digitized version of the 1989 DART pdf Admin/ramosRivera/Summer2021\_Dylan/1989\_Digitized

#### 2.3.4.4 Not in a folder

- **1989\_OCR\_Digitized.txt** - digitized version of the 1989 DART pdf from the OCR machine (unedited)
- **corr\_cattLabs97\_Wgeocode\_Line 6200 to Line 12765.csv** - excel file with lines 6200-12765 corrected by Kelly.Refer to original Admin/ramosRivera/T-Burk/PngData/corr\_cattLabs97\_Wgeocode
- **OneDrive\_2021-08-27** - zip drive that leads to the original material folder in Task From Antonio 1 folder.

#### 2.3.5 T-Burk Folder

It has 11 folders:

##### 2.3.5.1 ArcMap Folder

- **Converted\_Graphics (.cpg, .dbf, .prj, .shp, .shx)** - it only shows a green rectangle
- **Textile Labs (.cpg, .dbf, .prj, .sbn, .sbx, .shp, .shx)** -
- **ZCTAs (.cpg, .dbf, .prj, .sbn, .sbx, .shp, .shx)** - Opening the files in ArcMap it shows California in the ZCTAs areas and the location of the labs (dots).

##### 2.3.5.2 BlockData Folder

**2.3.5.2.1 nhgis0003\_shapefiles\_tl2000\_560\_block\_2000** together these files make up the the US by census block. These files is most likely the input for USA\_block\_emp.

- **AK\_block\_2000 (.dbf, .prj, .sbn, .sbx, .shp, .shx)** - Alaska by census block
- **AL\_block\_2000 (.dbf, xml, .sbn, shx, dbf, prj, adobe, shp)** - Alabama by census block
- **AR\_block\_2000 (.dbf, xml, .sbn, shx, dbf, prj, adobe, shp)** - Arkansas by census block
- **AZ\_block\_2000 (.dbf, xml, .sbn, shx, dbf, prj, adobe, shp)** - Arizona by census block
- **CA\_block\_2000 (.dbf, xml, .sbn, shx, dbf, prj, adobe, shp)** - California by census block
- **CO\_block\_2000 (.dbf, xml, .sbn, shx, dbf, prj, adobe, shp)** - Colorado by census block
- **CT\_block\_2000 (.dbf, xml, .sbn, shx, dbf, prj, adobe, shp)** - Connecticut by census block
- **DC\_block10\_2000 (.dbf, xml, .sbn, shx, dbf, prj, adobe, shp)** - Washington DC by census block
- **DE\_block\_2000 (.dbf, xml, .sbn, shx, dbf, prj, adobe, shp)** - Delaware by census block
- **FL-block\_2000 (.dbf, xml, .sbn, shx, dbf, prj, adobe, shp)** - Florida by census block
- **GA\_block\_2000(.dbf, xml, .sbn, shx, dbf, prj, adobe, shp)** - Georgia by census block
- **HI\_block\_2000 (.dbf, xml, .sbn, shx, dbf, prj, adobe, shp)** - Hawaii by census block
- **IA\_block\_2000 (.dbf, xml, .sbn, shx, dbf, prj, adobe, shp)** - Iowa by census block
- **ID\_block\_2000 (.dbf, xml, .sbn, shx, dbf, prj, adobe, shp)** - Idaho by census block
- **IL\_block\_2000 (.dbf, xml, .sbn, shx, dbf, prj, adobe, shp)** - Illinois by census block
- **IN\_block\_2000 (.dbf, xml, .sbn, shx, dbf, prj, adobe, shp)** - Indiana by census block
- **KS\_block\_2000 (.dbf, xml, .sbn, shx, dbf, prj, adobe, shp)** - Kansas by census block

- **KY\_block\_2000** (.dbf, xml, .sbn, shx, dbf, prj, adobe, shp) - Kentucky by census block
- **LA\_block\_2000** (.dbf, xml, .sbn, shx, dbf, prj, adobe, shp) - Louisiana by census block
- **MA\_block\_2000** (.dbf, xml, .sbn, shx, dbf, prj, adobe, shp) - Massachusetts by census block
- **MD\_block\_2000** (.dbf, xml, .sbn, shx, dbf, prj, adobe, shp) - Maryland by census block
- **ME\_block\_2000** (.dbf, xml, .sbn, shx, dbf, prj, adobe, shp) - Maine by census block
- **MI\_block\_2000** (.dbf, xml, .sbn, shx, dbf, prj, adobe, shp) - Michigan by census block
- **MN\_block\_2000** (.dbf, xml, .sbn, shx, dbf, prj, adobe, shp) - Minnesota by census block
- **MO\_block\_2000** (.dbf, xml, .sbn, shx, dbf, prj, adobe, shp) - Missouri by census block
- **MS\_block\_2000** (.dbf, xml, .sbn, shx, dbf, prj, adobe, shp) - Mississippi by census block
- **MT\_block\_2000** (.dbf, xml, .sbn, shx, dbf, prj, adobe, shp) - Montana by census block
- **NC\_block\_2000** (.dbf, xml, .sbn, shx, dbf, prj, adobe, shp) - North Carolina by census block
- **ND\_block\_2000** (.dbf, xml, .sbn, shx, dbf, prj, adobe, shp) - North Dakota by census block
- **NE\_block\_2000** (.dbf, xml, .sbn, shx, dbf, prj, adobe, shp) - Nebraska by census block
- **NH\_block\_2000** (.dbf, xml, .sbn, shx, dbf, prj, adobe, shp) - New Hampshire by census block
- **NJ\_block\_2000** (.dbf, xml, .sbn, shx, dbf, prj, adobe, shp) - New Jersey by census block
- **NM\_block\_2000** (.dbf, xml, .sbn, shx, dbf, prj, adobe, shp) - New Mexico by census block
- **NV\_block\_2000** (.dbf, xml, .sbn, shx, dbf, prj, adobe, shp) - Nevada by census block
- **NY\_block\_2000** (.dbf, xml, .sbn, shx, dbf, prj, adobe, shp) - New York by census block
- **OH\_block\_2000** (.dbf, xml, .sbn, shx, dbf, prj, adobe, shp) - Ohio by census block
- **OK\_block\_2000** (.dbf, xml, .sbn, shx, dbf, prj, adobe, shp) - Oklahoma by census block
- **OR\_block\_2000** (.dbf, xml, .sbn, shx, dbf, prj, adobe, shp) - Oregon by census block
- **PA\_block\_2000** (.dbf, xml, .sbn, shx, dbf, prj, adobe, shp) - Pennsylvania by census block
- **RI\_block\_2000** (.dbf, xml, .sbn, shx, dbf, prj, adobe, shp) - Rhode Island by census block
- **SC\_block\_2000** (.dbf, xml, .sbn, shx, dbf, prj, adobe, shp) - South Carolina by census block
- **SD\_block\_2000** (.dbf, xml, .sbn, shx, dbf, prj, adobe, shp) - South Dakota by census block
- **TN\_block\_2000** (.dbf, xml, .sbn, shx, dbf, prj, adobe, shp) - Tennessee by census block
- **TX\_block\_2000** (.dbf, xml, .sbn, shx, dbf, prj, adobe, shp) - Texas by census block
- **UT\_block\_2000** (.dbf, xml, .sbn, shx, dbf, prj, adobe, shp) - Utah by census block
- **VA\_block\_2000** (.dbf, xml, .sbn, shx, dbf, prj, adobe, shp) - Virginia by census block
- **VT\_block\_2000** (.dbf, xml, .sbn, shx, dbf, prj, adobe, shp) - Vermont by census block
- **WA\_block\_2000** (.dbf, xml, .sbn, shx, dbf, prj, adobe, shp) - Washington by census block
- **WI\_block\_2000** (.dbf, xml, .sbn, shx, dbf, prj, adobe, shp) - Wisconsin by census block
- **WV\_block\_2000** (.dbf, xml, .sbn, shx, dbf, prj, adobe, shp) - West Virginia by census block
- **WY\_block\_2000** (.dbf, xml, .sbn, shx, dbf, prj, adobe, shp) - Wyoming by census block

**2.3.5.2.2 Not in a folder** first 51 files are 7.zip files. There is one for each state + DC. You can open these files in python use this link for how. All of the .gz files are input files for `usa_block_emp.py`.

- `ak_wac_5000_JT00_2002.cvs.gz`
- `al_wac_5000_JT00_2002.cvs.gz`
- `ar_wac_5000_JT00_2002.cvs.gz`
- `az_wac_5000_JT00_2002.cvs.gz`
- `ca_wac_5000_JT00_2002.cvs.gz`
- `co_wac_5000_JT00_2002.cvs.gz`
- `ct_wac_5000_JT00_2002.cvs.gz`
- `dc_wac_5000_JT00_2002.cvs.gz`
- `de_wac_5000_JT00_2002.cvs.gz`
- `fl_wac_5000_JT00_2002.cvs.gz`
- `ga_wac_5000_JT00_2002.cvs.gz`
- `hi_wac_5000_JT00_2002.cvs.gz`
- `ia_wac_5000_JT00_2002.cvs.gz`
- `id_wac_5000_JT00_2002.cvs.gz`
- `il_wac_5000_JT00_2002.cvs.gz`
- `in_wac_5000_JT00_2002.cvs.gz`
- `ks_wac_5000_JT00_2002.cvs.gz`
- `ky_wac_5000_JT00_2002.cvs.gz`
- `la_wac_5000_JT00_2002.cvs.gz`
- `ma_wac_5000_JT00_2002.cvs.gz`
- `md_wac_5000_JT00_2002.cvs.gz`
- `me_wac_5000_JT00_2002.cvs.gz`
- `mi_wac_5000_JT00_2002.cvs.gz`
- `mn_wac_5000_JT00_2002.cvs.gz`
- `mo_wac_5000_JT00_2002.cvs.gz`
- `ms_wac_5000_JT00_2002.cvs.gz`
- `mt_wac_5000_JT00_2002.cvs.gz`
- `nc_wac_5000_JT00_2002.cvs.gz`
- `nd_wac_5000_JT00_2002.cvs.gz`
- `ne_wac_5000_JT00_2002.cvs.gz`
- `nh_wac_5000_JT00_2002.cvs.gz`
- `nj_wac_5000_JT00_2002.cvs.gz`
- `nm_wac_5000_JT00_2002.cvs.gz`
- `nv_wac_5000_JT00_2002.cvs.gz`

- ny\_wac\_5000\_JT00\_2002.cvs.gz
- oh\_wac\_5000\_JT00\_2002.cvs.gz
- ok\_wac\_5000\_JT00\_2002.cvs.gz
- or\_wac\_5000\_JT00\_2002.cvs.gz
- pa\_wac\_5000\_JT00\_2002.cvs.gz
- ri\_wac\_5000\_JT00\_2002.cvs.gz
- sc\_wac\_5000\_JT00\_2002.cvs.gz
- sd\_wac\_5000\_JT00\_2002.cvs.gz
- tn\_wac\_5000\_JT00\_2002.cvs.gz
- tx\_wac\_5000\_JT00\_2002.cvs.gz
- ut\_wac\_5000\_JT00\_2002.cvs.gz
- va\_wac\_5000\_JT00\_2002.cvs.gz
- vt\_wac\_5000\_JT00\_2002.cvs.gz
- wa\_wac\_5000\_JT00\_2002.cvs.gz
- wi\_wac\_5000\_JT00\_2002.cvs.gz
- wv\_wac\_5000\_JT00\_2002.cvs.gz
- wy\_wac\_5000\_JT00\_2002.cvs.gz
- USA\_block (.cpg, .dbf, .prj, .shp, .shx) - input and output for shapstich. Input file for USA\_block\_emp
- usa\_blockEmp (.cpg, .dbf, .prj, .shp, .shx) - I do not know what this is showing. Output file for USA\_block\_emp.

### 2.3.5.3 k-function\_local\_results Folder

- Manufa\_Emp\_C000\_0.5\_Buffers\_2 (.cpg, .dbf, .shp, .shx) -
- Manufa\_Emp\_C000\_0.25\_Buffers\_2 (.cpg, .dbf, .shp, .shx) - input file for start\_calc
- Manufa\_Emp\_C000\_0.75\_Buffers\_2 (.cpg, .dbf, .shp, .shx) -
- Manufa\_Emp\_C000\_1\_Buffers\_2 (.cpg, .dbf, .shp, .shx) -
- Manufa\_Emp\_C000\_2\_Buffers\_2 (.cpg, .dbf, .shp, .shx) - input file for start\_calc
- Manufa\_Emp\_C000\_5\_Buffers\_2 (.cpg, .dbf, .shp, .shx) -
- Manufa\_Emp\_C000\_10\_Buffers\_2 (.cpg, .dbf, .shp, .shx) - input file for start\_calc

All the previous files have missing spatial reference information. The data can be drawn in ArcMap , but not projected. ArcMap doesn't show anything

- Manufa\_Emp\_C000\_Points\_2 (.cpg, .dbf, .shp, .shx) - input file for start\_calc
- Manufa\_Emp\_C000\_local.txt - This is a log file with the date (04/05/2021) and time slot of some code running.

The files show the location of Manufacturing employment clusters I believe in California. They should correspond to Figure 1 and 2 of the draft.

#### 2.3.5.4 LabData Folder

- **cal\_lab\_fields** (.cpg, .dbf, .prj, .shp, .shx) – A folder for 34 different industries i.e. AERO –AERO . Output file for field\_org.
  - **Cal\_Labs.shp** (.cpg, .dbf, .prj, .shp, .shx) - input for field\_org, stat\_calc. Output for firm\_struc and shapify.
  - **comb\_emp\_C000\_local.txt** - log file
  - **Manufa\_Emp\_C000\_0.5\_Buffers\_cal0** (.cpg, .dbf, .shp, .shx) –
  - **Manufa\_Emp\_C000\_0.25\_Buffers\_cal0** (.cpg, .dbf, .shp, .shx) –
  - **Manufa\_Emp\_C000\_0.75\_Buffers\_cal0** (.cpg, .dbf, .shp, .shx) –
  - **Manufa\_Emp\_C000\_1\_Buffers\_cal0** (.cpg, .dbf, .shp, .shx) –
  - **Manufa\_Emp\_C000\_2\_Buffers\_cal0** (.cpg, .dbf, .shp, .shx) –
  - **Manufa\_Emp\_C000\_5\_Buffers\_cal0** (.cpg, .dbf, .shp, .shx) –
  - **Manufa\_Emp\_C000\_10\_Buffers\_cal0** (.cpg, .dbf, .shp, .shx) –
- All the previous files have missing spatial reference information. The data can be drawn in ArcMap , but not projected. ArcMap doesn't show anything
- **Manufa\_Emp\_C000\_Points\_cal0** (.cpg, .dbf, .shp, .shx) –
  - **Manufa\_Emp\_C000\_local.txt** - This is a log file with the date (04/18/2021) and time slot of some code running.
  - **USA\_labs\_2000** (.cpg, .dbf, .prj, .shp, .shx) - output file form prep\_Labs

#### 2.3.5.5 PatentData This is probably used to replicate Buzard 2017.

- **.RData** -
- **Rhistory** -
- **CA Control\_1\_ALT\_amos** (SAS Program)
- **cite\_same** (Excel) - input file for start\_calc
- **cite76\_06** (SAS Data set)
- **clustpatents** (SAS Data set)
- **columnsEFI\_CAbaseline** -
- **columnsEFI\_NEbaseline** -
- **LA5A\_ALT** (.cpg, .dbf, .shp, .shx) –
- **LA5B\_ALT** (.cpg, .dbf, .shp, .shx) –
- **LA5C\_ALT** (.cpg, .dbf, .shp, .shx) –
- **LA10A\_ALT** (.cpg, .dbf, .shp, .shx) –
- **LA10B\_ALT** (.cpg, .dbf, .shp, .shx) –
- **list\_of\_matches\_CAbaseline\_amos** –
- **originating** (SAS Data set) –
- **pat76\_06** (SAS Data set) –
- **replications\_CAbaseline** (Excel) –



- **SASclustpatentsCA** (Excel) –
- **SASoriginatingCA** (Excel) –
- **SASpossiblenclassCA** (Excel) –
- **SB5\_\_ALT** (.cpg, .dbf, .shp, .shx) –
- **SB10\_\_ALT** (.cpg, .dbf, .shp, .shx) –
- **SD5A\_\_ALT** (.cpg, .dbf, .shp, .shx) –
- **SD5B\_\_ALT** (.cpg, .dbf, .shp, .shx) –
- **SD10\_\_ALT** (.cpg, .dbf, .shp, .shx) –
- **SF5A\_\_ALT** (.cpg, .dbf, .shp, .shx) –
- **SF5B\_\_ALT** (.cpg, .dbf, .shp, .shx)–
- **SF10\_\_ALT** (.cpg, .dbf, .shp, .shx) –
- **tables** (word) –Table 2a is Table 3 is the draft –Table 2b is Table 4 in the draft –Table 3b is Table 5 in the draft. The draft only uses 5 and 10 miles ratio

#### 2.3.5.6 PngData Folder

- **OCR\_\_Output\_\_1998** –
- **letter\_\_I\_\_cattell.cvs** (excel) – not the same format as documents in Kelly
- **letter\_\_O\_\_cattell.cvs** (excel) –
- **letter\_\_S\_\_cattell.cvs** (excel) –
- **OCR\_\_Result** (text) – input for Address\_ID, output for OCR
- **OCR\_\_Result\_\_NO\_\_user** (text) All this files looks like the registry of labs.
- **1979 Digitized** (text): Registry of labs
- **1989\_\_OCR\_\_Digitized** (text): Registry of labs
- **calLabs97** (Excel): File with company name, facility name, state, ID and address for 1997
- **cattell\_\_1997\_\_raw** (STATA) - input file for state\_code\_rep
- **Cattell\_\_corr\_\_list** (STATA) - input file for pngwork.
- **cattell-all** (STATA)
- **cattLabs97** (Excel) - input file for Address\_ID, Geo\_coder, and state\_code\_rep. Output file for pngwork
- **CattwithBuzID** (Excel)
- **corr\_\_cattLabs97** (Excel) - input file for field\_org, firm\_struc, and start\_calc. Output files for state\_code\_rep.
- **corr\_\_cattLabs97\_\_Wgeocode** (Excel): This one has a column counting the observations. Only difference with the file below.
- **corr\_\_cattLabs97\_\_Wgeocode** (Excel)
- **field** (STATA) - input file for field\_org, pngwork, start\_calc
- **field\_\_lab\_\_counts** (EXCEL): count by sector. There are no differences with the file below
- **field\_\_lab\_\_counts2** (EXCEL): count by sector

- **field-master** (STATA) - input file for field\_org
- **geocoded\_facilities (EXCEL)**: has 8,737 observations. Input file for prep\_Labs. Input file for shapify. Output file for Geo\_coder
- **geocoded\_facilities\_cal (EXCEL)**: has 1,728 observations
- **geocoded\_facilities\_I (EXCEL)**: has 394 observations
- **geocoded\_facilities\_O (EXCEL)**: has 198 observations
- **geocoded\_facilities\_S (EXCEL)**: has 886 observations
- **id\_dataString (EXCEL)**: has the id, the full address and the buzzID
- **matched\_data (EXCEL)**: has 8,941 obs. input file for Geo\_coder.
- **matched\_data\_I (EXCEL)**: has 199 obs. Not sure what is matching or with which file.
- **matched\_data\_O (EXCEL)**: has 199 obs. Not sure what is matching or with which file.
- **matched\_data\_S (EXCEL)**: has 890 obs. Not sure what is matching or with which file.
- **newData (EXCEL)**: has 28,515 obs. 39 variables. information from the entire US (by looking at the states)
- **pngbuzz (EXCEL)**: has 2,951 obs. 39 variables. information from the entire US (by looking at the states)
- **pngCatIDList (EXCEL)**: has 11,313 obs. 5 variables. Output for Address\_ID.
- **single\_lab\_firm (EXCEL)**: has 7,430 obs. 21 variables. Input file for firm\_struc. Out file for pngwork.

Next step is to go to png website and see which files are downloaded from there and which ones were created by Antonio.

### 2.3.5.7 Python Scripts Folder

- **.pylint.d** -
- **stat\_calc1.stats (STATS)**:
- **Address\_ID** – Preparing and cleaning addresses
- **clust\_pat\_maker** – Python Script to read in patent data and conduct a spatial join with clusters then keep the patents that fall into those clusters as geodataframes and export them
- **countSim\_speedUP** – Point Count Simulation Computation. It turns the dictionary back into a dataframe.
- **countSim\_tester** – Same as “countSim\_speedUp:” but measure the time for each individual loop and the entire system (time it takes to perform the simulation).
- **field\_org** – Read in 1997 cattell lab data as well as my geocoded data and combine the two to produce geodataframes (gdf) for each technology field in the cattell directory which then get saved as shapefile currently the program is set in such a way so as to produce gdfs for the country wide data and gdfs for california and the NE corridor the script will then take the dictionary containing the california labs by field and save each gdf as a shapefile.
- **firm\_struc** – Reads in the 1997 Cattel directory data produced by Ivan Png 2016. It takes the data in this data set of American R&D Labs and organizes it based on firm structure. The final product is an Excel File that gets outputted.
- **GeoCode\_OCR**

- **GeoCoder** – Iterate through the address data and geocode each input address.
- **multiprocess\_test2** – processing time information
- **multiprocessing\_tester** –
- **OCR** – imports images, edit them to use with “tesseract”
- **Pdf2Jpg.py** –
- **pngwork** – create dataframe for firms that have at most 2 establishments, this will become the dataframe for firms with only one research establishment.
- **Prep\_Labs** – python script which reads in two csv of geocoded labs and joins them resulting in a pandas dataframe. Then it takes the coordinates for the labs in the dataframe and creates a geometry column to turn the df into a geodataframe it then saves the resulting geodataframe as a shapefile.
- **Prep\_ZBP** – python script that uses pandas and geopandas packages to read in census manufacturing employment data at the ZCTA level and shapefile of all ZCTA boundaries in the contiguous US. The employment data is prepared and the merged into the shapefiles data table resulting in a geopandas geodataframe which gets saved as a shapefile.
- **shapeStich** – create the new geodataframe by appending all state level gdfs.
- **shapify** – Reads in a csv file that contains point data in latitude and longitude form and converts them into geopandas geodataframe the result is then exported as a shapefile.
- **stat\_calc** – Read in the point file and associated cluster files produced after running the 3Stage\_Local program and calculates various statistics from it. Different sections of this program produce different stats and have been partitioned and commented accordingly.
- **state\_code\_rep** – Fix the state\_code column of the cattell png data for 1997
- **usa\_block\_emp** - reads in a list of csv files containing employment data at the block level. It then reads in a shape file of all US census block boundaries and merges the employment data into the shapefiles data table and saves the resulting geodataframe as a shapefile.

#### 2.3.5.8 Tables folder excel tables used in paper

- **5\_mile\_LDS** - Shows Originating Patents, Citing Patents, From Same Cluster, Percent (C/B), Treatment Patents, Treatment Citing For Same Cluster, Percent (F/E), Control Patents, Control Citing From Same Cluster, Percent (I/H), Location Differential (G/J), and P-values for 5-mile cluster in California. (excel) Output file for stat\_calc
- **10\_mile\_LDS** - Shows Originating Patents, Citing Patents, From Same Cluster, Percent (C/B), Treatment Patents, Treatment Citing For Same Cluster, Percent (F/E), Control Patents, Control Citing From Same Cluster, Percent (I/H), Location Differential (G/J), and P-values for 10-mile cluster in California. (excel) Output file for stat\_calc.
- **Spatial\_LDS** - Table that compares the 5 and 10 mile clusters (excel) Output file for stat\_calc.

#### 2.3.5.9 ZipData folder first folder is a duplicate folder of the “nhgis0005\_csv” folder found below

##### 2.3.5.9.1 tl\_2010\_us\_ZCTA500 Folder

- **tl\_2010\_us\_zcta500** (DBF, PRJ, Adobe, XML, SHX) - input file for prep\_ZBP

##### 2.3.5.9.2 Not in folder

- **OSF3\_geo\_header** - Data dictionary, explains U.S. Abbreviations, Geographic Area Codes by region, divisions, state (census, state (FIPS), county size code, FIPS County Subdivisions Class Code, Place Size Code, etc. (Word document)

- **Employment** - SAS Graph document created to collect ZIP code employment data for California
- **USA\_ZCTA\_emp(CPG, DBF, PRJ, SBN, SHX)** - Map of the US separated by ZCTA or zipcodes, output for prep\_ZBP.

**2.3.5.10 nhgis0005\_shape Folder** Files will not open because the folders are compressed.

**2.3.5.10.1 nhgis0004\_shapefile\_tl2000\_330\_block\_2000 (zipped Folder)**

- **NH\_block\_2000** - (DBF, PRJ, SHX, SBN, SBX, XML) -

**2.3.5.10.2 nhgis0004\_shapefile\_tl2010\_110\_block\_2000 (zipped Folder)**

- **DC\_block10\_2000** (DBF, PRJ, SHX, SBN, SBX, XML) -

**2.3.5.10.3 nhgis0004\_shapefile\_tl2010\_250\_block\_2000 (zipped Folder)**

- **MA\_block10\_2000** (DBF, PRJ, SHX, SBN, SBX, XML) -

**2.3.5.11 nhgis0005\_csv Folder** NHGIS data from 2000

- **nhgis0005\_ds151\_2000\_zcta** - excel files with 55 variables but only 32 variables have observations. Contains GISJOIN from the year 2000. Input file for prep\_ZBP.
- **nhgis0005\_ds151\_2000\_zcta\_codebook** - describes the variable labels in the nhgis0005\_ds151\_2000\_zcta excel file and where the data was ciphered from.
  - ex: GISJOIN: GIS Join Match Code
  - It also contains what the NHGIS codes are (ex: GMH001: Male » Agriculture, forestry, fishing and hunting, and mining)

**2.3.6 tl\_2010\_06\_zcta500 Folder**

Has five documents in different formats, builds map of California by census block.

- **tl\_2010\_06\_zcta500.dbf** - 11 observations: STATEFP00, ZCTA5CE00, GEOID00, CLASSFP00, MTFCC00, FUNCSTAT00, ALAND00, AWATER00, INTPTLAT00, INTPTLON00, PARTFLG00 (opened in excel)
- **tl\_2010\_06\_zcta500 (.prj, .shp, .shx, .xml)** Map of California by census block (opened with arcGIS), U.S Department of Commerce, U.S. Census Bureau, Geography Division 2010. (xml File)
  - Vector digital data from <http://www.census.gov/geo/www/tiger>

**2.3.7 ~ Not in a Folder ~**

- **cattell-all.sas** - 18 variables: Parent ID (new) (*referring to parent facility*), year, parent ID (Cattell original) (*referring to pdf scans parent facility ID*), Parent name, Facility name, Facility ID (Cattell original) (*referring to pdf scans facility ID*), Facility ID (new), zipcode (*the zipcode the facility is in*), Facility level, user, prof, doct, tech, parent name (alternative 2), parent name (alternative 3), parent name (alternative 4), state. (Stata file)
  - figure out what new vs cattell original is?
- **Dylan & Kelly notes from Summer 2021.pdf** - (5/30/2022) Dylan and Kelly's documentation on their work
- **Dylan & Kelly notes from Summer 2021.pdf** - (6/3/2022) Dylan and Kelly's documentation on their work with notes from Prof. Buzard

- **field.sas** - Shows Stata data on the Cattell ID for R&D fields and R&D sub-fields, the year the data was on, and the facilities (Stata file).
- **pngwork.py** - python script that uses the cattel-all.dta, field.dta, and a file called “GoodLabs.shp” for points (*this file is from Prof. Buzard’s earlier work*)

## 2.4 ramosRivera - Backup062122 Folder

This folder is a duplicate of the ramosRivera Folder created on June 21, 2022 as a backup to the original ramosRivera folder.

## 2.5 Deleted or not in the Admin Folder

- **OCR\_output** - input for OCR
- **ScanData** - input for pdf2Jpg.py