

Implementation for the Analysis of U.S. Incomes

Conrad Kosowsky

This file documents the code to fit various models to public-use CPS income data and produce associated results and figures. The actual code is contained in seven `py` files:

1. `gen_files.py`. The file to parse the data files. It assumes that the $\langle year \rangle$ CPS data file is contained in a directory “`../data/⟨year⟩ CPS/`” so that `data` is at the same level as the present working directory. This file will produce one data file in the present working directory called `data_⟨year⟩.txt` for each year, and the data files will contain four columns: household identifier, state, personal income, and survey weight.
2. `print_edges.py`. This file will read in all the data files and print the smallest and

Files

I	<code>gen_files.py</code>	3
II	<code>print_edges.py</code>	21
III	<code>bin.py</code>	22
IV	<code>estimate_parameters.py</code>	26
	Generalized Beta, Type II	32
	Dagum	39
	Burr	42
	Fisk	44
	Inverse-Gamma	47
	Davis	50
	Constant-Shift Inverse Gamma	56
	Constant-Shift-Scale Inverse Gamma	61
	Log-Normal, Pareto Cutoff	65
	Log-Normal, Pareto Mixture	73
V	<code>check_constants.py</code>	84
VI	<code>bootstrap.py</code>	87
VII	<code>make_figures.py</code>	93
VIII	<code>main.py</code>	101

Table 1: Variables in the Implementation

Purpose	Identifier	Meaning
From the Data	<code>household</code>	Household Identifier
	<code>income</code>	Individual Income
	<code>region</code>	Census Bureau Region Classification
	<code>state</code>	State Code
	<code>weight</code>	Survey Weight
Parameters (See Table 6)	<code>alpha, p, q</code>	Shape Parameters
	<code>beta, x0</code>	Scale Parameters
	<code>c</code>	Shift Parameter
	<code>gamma</code>	Mixture Parameter
	<code>k</code>	Cutoff Parameter
	<code>mu, sigma_sq</code>	Log-normal Parameters
In Binned Data (See Table 4)	<code>dens</code>	Density
	<code>freq</code>	Frequency
	<code>left</code>	Left Endpoint
	<code>mass</code>	Total Mass
	<code>mid</code>	Bin Midpoint
	<code>right</code>	Right Endpoint
	<code>width</code>	Bin Width

largest entries in the data. Used for manually determining the cutoff values before analysis.

3. `bin.py`. Code to bin data. It provides a single function that accepts a `DataFrame`, bins it, and returns the result.
4. `estimate_parameters.py`. The code to analyze the data. This file contains the routines to estimate parameters for all models of interest. It provides a number of estimation functions as well as other useful functions such as density and cumulative distribution. All user-level functions can be accessed through the `distribution`, `density`, `likelihood`, and `estimator` dictionaries.
5. `check_constants.py`. This file contains the code to run the linear regression on inverse-gamma parameters. It prints the results and writes them to a file.
6. `bootstrap.py`. This file contains the code to bootstrap standard errors using the synthetic survey design and estimation functions from `estimate_parameters.py`.
7. `make_figures.py`. This file contains functions that can produce figures using `pyplot`.
8. `main.py`. The file runs the analysis, produces parameter estimates, and makes the figures. This file loads `estimate_parameters.py`, `check_constants.py`, and `make_figures.py`. Switches at the start of the file determine which code executes when the user runs the file.

File I

gen_files.py

This file documents the code to parse the income information from the Current Population Survey public-use microdata files. We will store the information in a new set of smaller csv files, which reduces computing resources for repeated calculations. We will collect the following information: individual income record, household identifier, state, and survey weight. Table 2 lists the location of these records in the data files. For data files before 2019, the public-use data files encode each respondent's answers as a single string variable, where different responses occupy different positions in the string. For 2019 through 2023, the data files are separate csv files.

We begin by importing Pandas and defining a function to convert strings to integers if possible. We use `try_int` instead of `int` for earlier years of data because a very small number of the entries have errors.

```
1 import pandas as pd
2 def try_int(x):
3     try:
4         return int(x)
5     except:
6         return float("nan")
```

The dictionary `files` will store the filenames of all the files. The keys are years, and the values are filename strings.

```
7 files = {
8     #1964:"cpsmar64.dat",
9     #1966:"cpsmar66.dat",
10    1967:"cpsmar67.dat",
11    1968:"cpsmar68.dat",
12    1969:"cpsmar69.dat",
13    1971:"cpsmar71.dat",
```

Table 2: Locations of Variables in the Data Files

Year	Person Identifier	Income	Household Identifier	State	Survey Weight
1964–1967	starts with "00036"	170:180	10:17	43:45	20:30
1968–1971	first digit 1, 2, or 3	60:66	1:16	F37–38	204:216
1972–1975	first digit 1, 2, or 3	60:66	none	F37–38	204:216
1976	digits 7,8 are 1–39	246:253	0:6	H53–54	117:128
1977–1979	digits 7,8 are 1–39	246:253	0:6	H39–40	117:128
1980–1987	digits 7,8 are 1–39	247:254	0:6	H39–40	117:128
1988–2010	first digit is 3	439:447	1:6	H40–41	65:73
2011–2018	first digit is 3	579:587	1:6	H42–43	154:162
2019–2023	separate file	PTOTVAL	PH_SEQ	GESTFIPS	MARSUPWT

14 1972:"cpsmar72.dat",
15 1973:"cpsmar73.dat",
16 1974:"cpsmar74.dat",
17 1975:"cpsmar75.dat",
18 1976:"cpsmar76.dat",
19 1977:"cpsmar77.dat",
20 1978:"cpsmar78.dat",
21 1979:"cpsmar79.dat",
22 1980:"cpsmar80.dat",
23 1981:"cpsmar81.dat",
24 1982:"cpsmar82.dat",
25 1983:"cpsmar83.dat",
26 1984:"cpsmar84.dat",
27 1985:"cpsmar85.dat",
28 1986:"cpsmar86.dat",
29 1987:"cpsmar87.dat",
30 1988:"cpsmar88.dat",
31 1989:"cpsmar89.dat",
32 1990:"cpsmar90.dat",
33 1991:"cpsmar91.dat",
34 1992:"cpsmar92.dat",
35 1993:"cpsmar93.dat",
36 1994:"cpsmar94.dat",
37 1995:"cpsmar95.dat",
38 1996:"cpsmar96.dat",
39 1997:"cpsmar97.dat",
40 1998:"mar98pub.cps",
41 1999:"mar99pub.cps",
42 2000:"mar00supp.dat",
43 2001:"mar01supp.dat",
44 2002:"mar02supp.dat",
45 2003:"asec2003.pub",
46 2004:"asec2004.pub",
47 2005:"asec2005_pubuse.pub",
48 2006:"asec2006_pubuse.pub",
49 2007:"asec2007_pubuse_tax2.dat",
50 2008:"asec2008_pubuse.dat",
51 2009:"asec2009_pubuse.dat",
52 2010:"asec2010_pubuse.dat",
53 2011:"asec2011_pubuse.dat",
54 2012:"asec2012_pubuse.dat",
55 2013:"asec2013_pubuse.dat",
56 2014:"asec2014_pubuse_3x8_rerun.dat",
57 2015:"asec2015_pubuse.dat",
58 2016:"asec2016_pubuse_v3.dat",
59 2017:"asec2017_pubuse.dat",
60 2018:"asec2018_pubuse.dat",

Table 3: How the Microdata Stores State Identifiers

Year	Stored As	Corresponding Dictionary
1964–1967	(Modified) 1960 State Codes	st60_to_name1
1968–1972	(Modified) 1960 State Codes	st60_to_name2
1973–1976	(Modified) 1960 State Codes	st60_to_name3
1977–2010	(Modified) 1960 State Codes	st60_to_name4
2011–2023	Fips (GESTFIPS)	N/A

```

61 2019: ["hhpub19.csv", "pppub19.csv"],
62 2020: ["hhpub20.csv", "pppub20.csv"],
63 2021: ["hhpub21.csv", "pppub21.csv"],
64 2022: ["hhpub22.csv", "pppub22.csv"],
65 2023: ["hhpub23.csv", "pppub23.csv"]}

```

The dictionary `regions` stores the region information. It would be simpler to use the region entry from the microdata, but for historical reasons, we convert FIPS state codes to region numbers. The keys in `regions` are integer FIPS codes, and the values are region identifiers. (Northeast is 1; Midwest is 2; South is 3; and West is 4.) Storage of state codes varies from year to year, and Table 2 shows where the microdata contains the state code. Earlier years of microdata encode the state identifier as (a modified version of) the 1960 Census state codes, and the microdata don't contain explicit state FIPS codes until 2001. (Note that in 2005, the Census Bureau renamed HG-ST60 to GESTCEN.) Table 3 shows how each year of microdata stores the state. For our implementation, we use HG-ST60 through 2010 and then switch to GESTFIPS.

```

66 regions = {
67     1:3,    # Alabama
68     2:4,    # Alaska
69     # 3 is American Samoa
70     4:4,    # Arizona
71     5:3,    # Arkansas
72     6:4,    # California
73     # 7 is Canal Zone
74     8:4,    # Colorado
75     9:1,    # Connecticut
76     10:3,   # Delaware
77     11:3,   # DC
78     12:3,   # Florida
79     13:3,   # Georgia
80     # 14 is Guam
81     15:4,   # Hawaii
82     16:4,   # Idaho
83     17:2,   # Illinois
84     18:2,   # Indiana
85     19:2,   # Iowa
86     20:2,   # Kansas

```

```

87  21:3, # Kentucky
88  22:3, # Louisiana
89  23:1, # Maine
90  24:3, # Maryland
91  25:1, # Massachusetts
92  26:2, # Michigan
93  27:2, # Minnesota
94  28:3, # Mississippi
95  29:2, # Missouri
96  30:4, # Montana
97  31:2, # Nebraska
98  32:4, # Nevada
99  33:1, # New Hampshire
100 34:1, # New Jersey
101 35:4, # New Mexico
102 36:1, # New York
103 37:3, # North Carolina
104 38:2, # North Dakota
105 39:2, # Ohio
106 40:3, # Oklahoma
107 41:4, # Oregon
108 42:1, # Pennsylvania
109 # 43 is Puerto Rico
110 44:1, # Rhode Island
111 45:3, # South Carolina
112 46:2, # South Dakota
113 47:3, # Tennessee
114 48:3, # Texas
115 49:4, # Utah
116 50:1, # Vermont
117 51:3, # Virginia
118 # 52 is Virgin Islands
119 53:4, # Washington
120 54:3, # West Virginia
121 55:2, # Wisconsin
122 56:4, # Wyoming

```

Jolliffe (2003) doesn't specify how to handle fips codes that do not correspond to a state (and hence do not have an official region). We group them geographically with corresponding regions.

```

123  3:4, # American Samoa
124 60:4, # American Samoa
125 81:4, # Baker Island
126  7:3, # Canal Zone
127 64:4, # Federated States of Micronesia
128 14:4, # Guam
129 66:4, # Guam
130 84:4, # Howland Island

```

```

131 86:4, # Jarvis Island
132 67:4, # Johnston Atoll
133 89:4, # Kingman Reef
134 68:4, # Marshall Islands
135 71:4, # Midway Islands
136 74:4, # Minor Outlying Islands
137 76:3, # Navassa Island
138 69:4, # Northern Mariana Islands
139 70:4, # Palau
140 95:4, # Palmyra Atoll
141 43:3, # Puerto Rico
142 72:3, # Puerto Rico
143 52:3, # Virgin Islands
144 78:3, # Virgin Islands
145 79:4} # Wake Island

```

These non-state fips codes shouldn't appear in the data since the microdata says the FIPS codes run from 1 to 56. However, we code them just to be safe. Some of the earlier CPS datasets store the household state information using 1960 Census state codes (variable HG-ST60 instead of GESTFIPS), which are different from FIPS codes. Accordingly, we create several dictionaries to convert early-year state codes to FIPS values. These dictionaries use (modified versions of) 1960 state codes as keys and state names as values. See Table 3.

```

146 st60_to_name1 = {
147   # New England
148   1: "Maine",
149   3: "New Hampshire",
150   4: "Vermont",
151   2: "Massachusetts",
152   5: "Rhode Island",
153   6: "Connecticut",
154   # Middle Atlantic
155   10: "New York",
156   11: "New Jersey",
157   13: "Pennsylvania",
158   # East North Central
159   24: "Ohio",
160   23: "Indiana",
161   25: "Illinois",
162   26: "Michigan",
163   22: "Wisconsin",
164   # West North Central
165   31: "Minnesota",
166   32: "Iowa",
167   33: "Missouri",
168   34: "North Dakota",
169   35: "South Dakota",
170   36: "Nebraska",
171   37: "Kansas",

```

```

172 # South Atlantic
173 41: "Delaware",
174 42: "Maryland",
175 43: "DC",
176 44: "Virginia",
177 45: "West Virginia",
178 47: "North Carolina",
179 46: "South Carolina",
180 48: "Georgia",
181 49: "Florida",
182 # East South Central
183 51: "Kentucky",
184 52: "Tennessee",
185 53: "Alabama",
186 54: "Mississippi",
187 # West South Central
188 65: "Arkansas",
189 66: "Louisiana",
190 67: "Oklahoma",
191 68: "Texas",
192 # Mountain
193 71: "Montana",
194 72: "Idaho",
195 73: "Wyoming",
196 74: "Colorado",
197 75: "New Mexico",
198 76: "Arizona",
199 77: "Utah",
200 78: "Nevada",
201 # Pacific
202 87: "Washington",
203 88: "Oregon",
204 89: "California",
205 85: "Alaska",
206 86: "Hawaii"}

```

Another version. The 1969 microdata combines a number of states into single state codes. Fortunately, it won't affect our region assignment.

```

207 st60_to_name2 = {
208 ## New England
209 11: "Connecticut",
210 # 19 is also Massachusetts, New Hampshire,
211 # Rhode Island, and Vermont
212 19: "Maine",
213 ## Middle Atlantic
214 21: "New York",
215 22: "New Jersey",
216 23: "Pennsylvania",

```



```

217  ## East North Central
218  31: "Ohio",
219  32: "Indiana",
220  33: "Illinois",
221  # 39 also contains Wisconsin
222  39: "Michigan",
223  ## West North Central
224  # 41 also contains Minnesota
225  41: "Iowa",
226  43: "Missouri",
227  # 49 also contains Nebraska, Kansas, and South Dakota
228  49: "North Dakota",
229  ## South Atlantic
230  51: "DC",
231  52: "Maryland",
232  53: "West Virginia",
233  54: "Georgia",
234  55: "Florida",
235  # 57 also contains South Carolina
236  57: "North Carolina",
237  # 59 also contains Virginia
238  59: "Delaware",
239  ## East South Central
240  61: "Kentucky",
241  62: "Tennessee",
242  # 69 also contains Mississippi
243  69: "Alabama",
244  ## West South Central
245  71: "Louisiana",
246  72: "Texas",
247  # 79 also contains Oklahoma
248  79: "Arkansas",
249  ## Mountain
250  # 81 also contains Colorado and New Mexico
251  81: "Arizona",
252  # 89 also contains Montana, Nevada, Utah, and Wyoming
253  89: "Idaho",
254  ## Pacific
255  91: "Oregon",
256  92: "California",
257  # 99 also contains Hawaii and Washington
258  99: "Alaska"}

```

Another version. Same deal with the 1973 microdata

```

259 st60_to_name3 = {
260  ## New England
261  16: "Connecticut",
262  14: "Massachusetts",

```

```

263 # 19 is also New Hampshire,
264 # Rhode Island, and Vermont
265 19: "Maine",
266 ## Middle Atlantic
267 21: "New York",
268 22: "New Jersey",
269 23: "Pennsylvania",
270 ## East North Central
271 31: "Ohio",
272 32: "Indiana",
273 33: "Illinois",
274 # 39 also contains Wisconsin
275 39: "Michigan",
276 ## West North Central
277 # 49 also contains Iowa, North Dakota, South Dakota,
278 # Nebraska, Kansas, and Missouri
279 49: "Minnesota",
280 ## South Atlantic
281 53: "DC",
282 56: "North Carolina",
283 # 57 also contains Maryland, Virginia, and West Virginia
284 57: "Delaware",
285 # 58 also contains Georgia
286 58: "South Carolina",
287 59: "Florida",
288 ## East South Central
289 # 67 also contains Tennessee
290 67: "Kentucky",
291 # 69 also contains Mississippi
292 69: "Alabama",
293 ## West South Central
294 72: "Texas",
295 # 79 also contains Oklahoma and Louisiana
296 79: "Arkansas",
297 ## Mountain
298 # 89 also contains Idaho, Wyoming, Colorado, New Mexico,
299 # Utah, and Nevada
300 89: "Montana",
301 ## Pacific
302 92: "California",
303 # 99 also contains Hawaii, Washington, Oregon, and Alaska
304 99: "Alaska"}

```

Main version. This is the version that we use for most of the years of microdata with 1960 Census state codes.

```

305 st60_to_name4 = {
306 # New England
307 11: "Maine",

```

308 12: "New Hampshire",
309 13: "Vermont",
310 14: "Massachusetts",
311 15: "Rhode Island",
312 16: "Connecticut",
313 # Middle Atlantic
314 21: "New York",
315 22: "New Jersey",
316 23: "Pennsylvania",
317 # East North Central
318 31: "Ohio",
319 32: "Indiana",
320 33: "Illinois",
321 34: "Michigan",
322 35: "Wisconsin",
323 # West North Central
324 41: "Minnesota",
325 42: "Iowa",
326 43: "Missouri",
327 44: "North Dakota",
328 45: "South Dakota",
329 46: "Nebraska",
330 47: "Kansas",
331 # South Atlantic
332 51: "Delaware",
333 52: "Maryland",
334 53: "DC",
335 54: "Virginia",
336 55: "West Virginia",
337 56: "North Carolina",
338 57: "South Carolina",
339 58: "Georgia",
340 59: "Florida",
341 # East South Central
342 61: "Kentucky",
343 62: "Tennessee",
344 63: "Alabama",
345 64: "Mississippi",
346 # West South Central
347 71: "Arkansas",
348 72: "Louisiana",
349 73: "Oklahoma",
350 74: "Texas",
351 # Mountain
352 81: "Montana",
353 82: "Idaho",
354 83: "Wyoming",

```

355 84: "Colorado",
356 85: "New Mexico",
357 86: "Arizona",
358 87: "Utah",
359 88: "Nevada",
360 # Pacific
361 91: "Washington",
362 92: "Oregon",
363 93: "California",
364 94: "Alaska",
365 95: "Hawaii"}

```

The next dictionary uses state (and territory) names as keys and FIPS codes as values.

```

366 name_to_fips = {
367     # states plus DC
368     "Alabama": 1,
369     "Alaska": 2,
370     "Arizona": 4,
371     "Arkansas": 5,
372     "California": 6,
373     "Colorado": 8,
374     "Connecticut": 9,
375     "Delaware": 10,
376     "DC": 11,
377     "Florida": 12,
378     "Georgia": 13,
379     "Hawaii": 15,
380     "Idaho": 16,
381     "Illinois": 17,
382     "Indiana": 18,
383     "Iowa": 19,
384     "Kansas": 20,
385     "Kentucky": 21,
386     "Louisiana": 22,
387     "Maine": 23,
388     "Maryland": 24,
389     "Massachusetts": 25,
390     "Michigan": 26,
391     "Minnesota": 27,
392     "Mississippi": 28,
393     "Missouri": 29,
394     "Montana": 30,
395     "Nebraska": 31,
396     "Nevada": 32,
397     "New Hampshire": 33,
398     "New Jersey": 34,
399     "New Mexico": 35,
400     "New York": 36,

```

```

401 "North Carolina":      37,
402 "North Dakota":       38,
403 "Ohio":                39,
404 "Oklahoma":           40,
405 "Oregon":              41,
406 "Pennsylvania":       42,
407 "Rhode Island":       44,
408 "South Carolina":     45,
409 "South Dakota":       46,
410 "Tennessee":          47,
411 "Texas":               48,
412 "Utah":                49,
413 "Vermont":             50,
414 "Virginia":            51,
415 "Washington":         53,
416 "West Virginia":      54,
417 "Wisconsin":           55,
418 "Wyoming":             56,
419 # territories
420 "American Samoa":      3,
421 "American Samoa":      60,
422 "Baker Island":        81,
423 "Canal Zone":          7,
424 "Federated States of Micronesia": 64,
425 "Guam":                14,
426 "Guam":                66,
427 "Howland Island":      84,
428 "Jarvis Island":       86,
429 "Johnston Atoll":      67,
430 "Kingman Reef":        89,
431 "Marshall Islands":    68,
432 "Midway Islands":      71,
433 "Minor Outlying Islands": 74,
434 "Navassa Island":      76,
435 "Northern Mariana Islands": 69,
436 "Palau":               70,
437 "Palmyra Atoll":       95,
438 "Puerto Rico":        43,
439 "Puerto Rico":        72,
440 "Virgin Islands":      52,
441 "Virgin Islands":      78,
442 "Wake Island":         79}

```

Check that we spelled all the state names correctly.

```

443 for i in range(1,5):
444     temp = eval("st60_to_name{0}".format(i))
445     for k in temp:
446         if temp[k] not in name_to_fips:

```

```

447         raise KeyError(
448 "for st60_to_name{0}: {1} not in name_to_fips".format(i, temp[k]))
449 print("\nIndex check for dicts is good\n")

```

The dictionary `year_keys` stores the year information from Table 2. The keys are years, and the values are the first year with the same variable locations according to the table. However, it is more compact to tell Python to assemble the dictionary for us. We code the list `base_years`, which contains the first year associated with each row in Table 2.

```

450 base_years = [1964, 1968, 1973, 1976, 1977, 1980, 1988, 2011, 2019]

```

The list does not contain 1972 because we will use the 1968 file reader for the 1972 data files. We loop through years from 1964 to 2023 and populate `year_keys` with elements of this range as keys. The pointer `curr_id` keeps track of where we are in `base_years`, and `curr_year` and `next_year` keep track of current and successive entries in `base_years`. On each iteration of the loop, we first check if `i` equals the `next_year`. That means we are outside the target interval and need to update `curr_id`. The `next_year` becomes the `curr_year`, and we take a new `next_year` from `base_years` if possible. Then we set the value in `year_keys` to be `curr_year`.

```

451 year_keys = {}
452 curr_id = 0
453 curr_year = base_years[curr_id]
454 next_year = base_years[curr_id + 1]
455 for i in range(1964, 2024):
456     if i == next_year:
457         curr_year = next_year
458         try:
459             curr_id = curr_id + 1
460             next_year = base_years[curr_id + 1]
461         except:
462             next_year = False
463     year_keys[i] = curr_year

```

Now define the functions to parse the data files. Each function will accept a filename and return a formatted DataFrame. The DataFrame will contain four columns: one for income, one for survey weight, one for state code, and one for household identifier. For most years, we manually parse the file, but for 2019, where the data file becomes csv, we have Pandas parse the file contents. We call each function `read_income_⟨key⟩`, where the *key* is the first year associated with the current year according to the first column of Table 2.

We begin with 1964 and progress chronologically. For 1964, we make blank lists for income, identifier, state, and weight, and we loop through the lines of the file. For each line, we pull the relevant information if the entry corresponds to a person from the March supplement. We create a DataFrame called `temp` and convert the state codes to FIPS values. (We don't add the region until we're ready to write the data extract file.)

```

464 def read_income_1964(filename):
465     household = []
466     income = []
467     state = []
468     weight = []

```

```

469 f = open(filename)
470 for line in f:
471     if line[0:5] == "00036":
472         household.append(try_int(line[10:17]))
473         income.append(try_int(line[170:180]))
474         state.append(try_int(line[43:45]))
475         weight.append(try_int(line[20:30]))
476 f.close()
477 temp = pd.DataFrame(
478     {"household": household,
479      "income": income,
480      "state": state,
481      "weight": weight}).dropna()
482 temp["state"] = temp["state"].map(st60_to_name1).map(name_to_fips)
483 return temp

```

For 1968–1975, we add our own household identifier, which will be an integer starting from 1 that increments at each household record. (The earliest CPS data files do not make as conspicuous a distinction between household and family records as in later iterations.) The first family record in each household has nonzero entry for F47–48, so we increment our household identifier whenever we encounter such a record. Once again, we change the state identifier to a FIPS code.

```

484 def read_income_1968(filename):
485     household = []
486     income = []
487     state = []
488     weight = []
489     f = open(filename)
490     curr_household = 0
491     for line in f:
492         if line[0] == "4": # family record
493             if line[46:48] != "00": # household record
494                 curr_household = curr_household + 1
495                 try_int_trace = 1
496                 curr_state = try_int(line[36:38])
497                 try_int_trace = 0
498             else: # person record
499                 household.append(curr_household)
500                 income.append(try_int(line[60:66]))
501                 state.append(curr_state)
502                 weight.append(try_int(line[204:216]))
503     f.close()
504     temp = pd.DataFrame(
505         {"household": household,
506          "income": income,
507          "state": state,
508          "weight": weight}).dropna()
509     temp["state"] = temp["state"].map(st60_to_name2).map(name_to_fips)

```

```
510     return temp
```

Parsing the 1973–1975 years is exactly the same as 1968 except we use a different dictionary to get the FISP codes.

```
511 def read_income_1973(filename):
512     household = []
513     income    = []
514     state     = []
515     weight    = []
516     f = open(filename)
517     curr_household = 0
518     for line in f:
519         if line[0] == "4": # family record
520             if line[46:48] != "00": # household record
521                 curr_household = curr_household + 1
522                 curr_state = try_int(line[36:38])
523             else: # person record
524                 household.append(curr_household)
525                 income.append(try_int(line[60:66]))
526                 state.append(curr_state)
527                 weight.append(try_int(line[204:216]))
528     f.close()
529     temp = pd.DataFrame(
530         {"household": household,
531          "income": income,
532          "state": state,
533          "weight": weight}).dropna()
534     temp["state"] = temp["state"].map(st60_to_name3).map(name_to_fips)
535     return temp
```

For 1976–1987, the functions are very similar. The main difference between 1968 is that we look at digits 7 and 8 to determine household versus family versus person records. (The household records are 0 here.) At this point, the data is good enough that we can use `int` rather than `try_int`. Again, we convert the state codes to FIPS numbers.

```
536 def read_income_1976(filename):
537     household = []
538     income    = []
539     state     = []
540     weight    = []
541     f = open(filename)
542     for line in f:
543         type = int(line[6:8])
544         if type == 0: # household record
545             curr_state = int(line[52:54])
546         elif 1 <= type and type <= 39: # person record
547             household.append(int(line[0:6]))
548             income.append(int(line[246:253]))
549             state.append(curr_state)
```



```

550     weight.append(int(line[117:128]))
551 f.close()
552 temp = pd.DataFrame(
553     {"household": household,
554      "income": income,
555      "state": state,
556      "weight": weight})
557 temp["state"] = temp["state"].map(st60_to_name3).map(name_to_fips)
558 return temp

```

In 1977, the location of the state record is different, and we begin using the main (fourth) `st60_to_name` dictionary to get the FIPS codes. Otherwise the function is the same.

```

559 def read_income_1977(filename):
560     household = []
561     income = []
562     state = []
563     weight = []
564     f = open(filename)
565     for line in f:
566         type = int(line[6:8])
567         if type == 0:                                # household record
568             curr_state = int(line[38:40])
569         elif 1 <= type and type <= 39: # person record
570             household.append(int(line[0:6]))
571             income.append(int(line[246:253]))
572             state.append(curr_state)
573             weight.append(int(line[117:128]))
574     f.close()
575     temp = pd.DataFrame(
576         {"household": household,
577          "income": income,
578          "state": state,
579          "weight": weight})
580     temp["state"] = temp["state"].map(st60_to_name4).map(name_to_fips)
581     return temp

```

In 1980, the location of the income record changes.

```

582 def read_income_1980(filename):
583     household = []
584     income = []
585     state = []
586     weight = []
587     f = open(filename)
588     for line in f:
589         type = int(line[6:8])
590         if type == 0:                                # household record
591             curr_state = int(line[38:40])
592         elif 1 <= type and type <= 39: # person record
593             household.append(int(line[0:6]))

```

```

594     income.append(int(line[247:254]))
595     state.append(curr_state)
596     weight.append(int(line[117:128]))
597 f.close()
598 temp = pd.DataFrame(
599     {"household": household,
600      "income": income,
601      "state": state,
602      "weight": weight})
603 temp["state"] = temp["state"].map(st60_to_name4).map(name_to_fips)
604 return temp

```

The structure of the files changes in 1988. Now household records start with a 1, and person records start with 3.

```

605 def read_income_1988(filename):
606     household = []
607     income     = []
608     state      = []
609     weight     = []
610     f = open(filename)
611     for line in f:

```

The 1996 data file has an extra character at the very end of the file, which breaks the parser as written. So we check that the length of each line is greater than 2 before processing it. This problem appears to be confined to the 1996 and 1997 data files.

```

612     if len(line) > 2:
613         type = int(line[0])
614         if type == 1:    # household record
615             curr_state = int(line[39:41])
616         elif type == 3:  # person record
617             household.append(int(line[1:6]))
618             income.append(int(line[439:447]))
619             state.append(curr_state)
620             weight.append(int(line[65:73]))
621     f.close()
622     temp = pd.DataFrame(
623         {"household": household,
624          "income": income,
625          "state": state,
626          "weight": weight})
627     temp["state"] = temp["state"].map(st60_to_name4).map(name_to_fips)
628     return temp

```

In 2011, the location of the income and weight variables changes. At this point, we can stop manually converting the state codes to FIPS values since we now pull that information directly from the data file.

```

629 def read_income_2011(filename):
630     household = []
631     income    = []

```

```

632 state      = []
633 weight     = []
634 f = open(filename)
635 for line in f:
636     type = int(line[0])
637     if type == 1:    # household record
638         curr_state = int(line[41:43])
639     elif type == 3:  # person record
640         household.append(int(line[1:6]))
641         income.append(int(line[579:587]))
642         state.append(curr_state)
643         weight.append(int(line[154:162]))
644 f.close()
645 return pd.DataFrame(
646     {"household": household,
647      "income": income,
648      "state": state,
649      "weight": weight})

```

For 2019 and on, the data files are separate csv files, so we have Pandas parse them.

```

650 def read_income_2019(h_filename, p_filename):
651     household = pd.read_csv(h_filename, header=0)
652     person = pd.read_csv(p_filename, header=0)
653     temp1 = person[["PH_SEQ", "PTOTVAL", "MARSUPWT"]].rename(
654         {"PH_SEQ": "household",
655          "PTOTVAL": "income",
656          "MARSUPWT": "weight"}, axis=1)
657     temp2 = household[["H_SEQ", "GESTFIPS"]].rename(
658         {"H_SEQ": "household",
659          "GESTFIPS": "state"}, axis=1)
660     result = pd.merge(temp1, temp2, how="inner", on="household",
661                       sort=False, validate="many_to_one")
662     return result

```

Now call the parsing functions and write the new data files.

```

663 for i in files:
664     print("Year {0}".format(i))
665     if i < 2019:
666         data = eval("read_income_" + str(year_keys[i]))(
667             "../data/{0} CPS/{1}".format(i, files[i]))
668     else:
669         data = read_income_2019("../data/{0} CPS/{1}".format(i, files[i][0]),
670                                 "../data/{0} CPS/{1}".format(i, files[i][1]))

```

Use the `regions` dictionary to map state fips codes onto region codes. For observations where the fips code does not correspond to a state, this likely indicates some error in the data file itself, so we drop those observations.

```

671 data["region"] = data["state"].map(regions)
672 temp = data[~((data["region"] >= 1) | (data["region"] <= 4))]

```

```

673 if len(temp) > 0:
674     print("""
675 Unable to assign regions for some rows;
676 will drop them. Portion of data getting
677 dropped: """)
678     print(temp)
679     print()
680 data = data[data["region"].notna()]
681 data["state"] = data["state"].astype(int)
682 data["region"] = data["region"].astype(int)

```

Finally, we sort the DataFrame and write it to a txt file.

```

683 data.sort_index(axis=1, inplace=True)
684 data.sort_values("income", inplace=True)
685 data.to_csv("data_{0}.txt".format(i), index=False)
686 print()

```

File II

print_edges.py

We need code to print the smallest and largest incomes for manual examination. This will allow us to form the correct boundary incomes for the estimation file.

```
1 import pandas as pd
2 for i in range(1967, 2024):
3     if i != 1970:
4         print("YEAR", i)
5         d = pd.read_csv("data_{0}.txt".format(i), header=0)["income"]
6         print(d.head(30))
7         print(d.tail(30))
8         print("\n\n\n")
```

File III

bin.py

The file provides a single function, `bin_data`, which accepts a `DataFrame` and produces a sample density also containing a few extra pieces of information. See Table 4 for a specific list of columns contained in the return value. The general approach here is to first establish bin boundaries, then loop through rows in the income `DataFrame`. As we loop through incomes, we simultaneously loop through bins and increment the weight on the corresponding bin. Doing both loops at the same time is the smart way to bin the data! Running a loop within a loop—for each income value, loop through bins until we find the correct bin—is hugely inefficient. (It would be much more expedient to use `np.histogram()` here, but I didn’t know about that function when I first wrote this code. Live and learn.)

We use linear bins below some cutoff value and logarithmic bins above it. If we implement a combination linear-logarithmic binning scheme, we specify three parameters to determine the bins: linear-logarithmic cutoff, number of linear bins between 0 and the cutoff (which determines linear bin width), and the logarithmic binning factor. See Table 5 for a list of default values. The choice of factor makes the first logarithmic bin lightly larger than the linear bins. After that, we will continue binning using the same factor until we reach the end of the data. We set the parameters as function arguments with default values. That way we can change their values if need be when calling `bin_data` without having to modify the file.

We begin by loading Pandas.

```
1 import pandas as pd
```

We load the exponential and logarithm functions from Numpy.

```
2 from numpy import exp
3 from numpy import log
```

The `bin_data` function will have three arguments. The `data` argument should be a `DataFrame`, and we check the object type for simple error checking. The `var` argument is the name of a column in `data` that we want to use as the variable for binning, and the `wgt` is an optional argument for the name of a column in `data` that contains survey weights. If the user calls `bin_data` with `wgt`, we will use that column in determining the total bin mass. Otherwise, we weight each observation equally.

Table 4: Columns in Binned Data

Information	Key
Left Endpoint	<code>left</code>
Right Endpoint	<code>right</code>
Midpoint	<code>mid</code>
Length of Bin	<code>width</code>
Total Mass in Bin	<code>mass</code>
Fraction of Mass in Bin	<code>freq</code>
Sample Density	<code>dens</code>

```

4 def bin_data(data, var, wgt=None, *, cutoff=57000,
5             num_linear_from_0=10, factor=1.2):
6     if not isinstance(data, pd.DataFrame):
7         raise ValueError("First argument of bin_data should be a DataFrame")
8     if var not in data:
9         raise KeyError("Variable to bin is not a column in the data")
10    if isinstance(wgt, type(None)):
11        called_with_wgt = False
12    else:
13        if wgt not in data:
14            raise KeyError("Weight variable is not a column in the data")
15        called_with_wgt = True

```

To make sure the data is monotonic, we sort it ourselves. To avoid any potential issues with indexing, we also reset the index.

```

16    data_to_use = data.sort_values(var).reset_index(drop=True)
17    n = len(data_to_use)

```

We create the linear portion of the bins as a list. For the linear bins, we want to have roughly `num_linear_from_0` linear bins between 0 and the cutoff, and we subtract 0.5 in order to guarantee a bin midpoint at 0. We calculate the linear bin width by dividing the cutoff by the adjusted number of bins. The `min` and `max` are minimum and maximum incomes in the data.

```

18    lin_width = cutoff / (num_linear_from_0 - 0.5)
19    min = data_to_use.loc[(0, var)]
20    max = data_to_use.loc[(n-1, var)]

```

We calculate the number of linear bins by doing integer division. The variable `remainder` is the distance from the minimum income to the left endpoint of the first full linear bin. Once we know the `remainder`, we can easily calculate the linear bin endpoints.

```

21    num_linear_bins = int((cutoff - min) // lin_width)
22    remainder = (cutoff - min) % lin_width
23    bins = [min]
24    temp = min + remainder
25    while temp < max and temp <= cutoff:
26        bins.append(temp)
27        temp = temp + lin_width
28    temp = cutoff * factor
29    while temp < max:
30        bins.append(temp)
31        temp = temp * factor

```

Table 5: Default Bin Specifications

Specification	Keyword	Default Value
Linear-logarithmic Cutoff	<code>cutoff</code>	\$57,000
Linear Bins from 0 to Cutoff	<code>num_lin_from_0</code>	$10 - 0.5 = 9.5$ (width \$6000)
Logarithmic Factor	<code>factor</code>	1.2

We copy `bins`, remove the zeroth entry, and append `max` to the copy. We create a `DataFrame` where `bins` becomes the `left` column, and the updated `next_bins` becomes the `right` column. We use `binned` to refer to our eventual return value.

```

32 next_bins = bins.copy()
33 next_bins.pop(0)
34 next_bins.append(max)
35 binned = pd.DataFrame({"left":bins, "right":next_bins})
36 binned["mid"] = 0.5 * (binned.left + binned.right)
37 binned["width"] = binned.right - binned.left

```

Now we do the actual binning. We create a list of total weights in each bin. We keep track of:

- `i`: the current row in `data`
- `j`: the current index in the bin list
- `k`: the current total weight for bin `j`

First we make the list. We loop through `range(len(bins))`, and on each iteration, we set `k` to be 0, then progressively increase `i` until we reach either a row where the income entry equals the next entry in `bins` or the final row in the data. The last bin in the data is a bit odd in terms of the loop, so we code it manually instead.

```

38 mass_vals = [0 for _ in bins]
39 i = 0
40 for j in range(len(bins)):
41     k = 0
42     while data_to_use.loc[(i, var)] < next_bins[j]:
43         if called_with_wgt: # if weight was specified, use survey weights
44             k = k + data_to_use.loc[(i, wgt)]
45         else:                # otherwise, each point adds mass of 1
46             k = k + 1
47         i = i+1
48     mass_vals[j] = k

```

Now add the last data points. We could have multiple observations where the variable takes value `max`, so we can't simply add the final row information to `mass_vals`. Instead, we count backwards to the first data entry that is different from `max`. (This approach assumes `max` \neq `min`, which should not be a problem.)

```

49 k = 0
50 i = n-1
51 while data_to_use.loc[(i, var)] == max:
52     if called_with_wgt:
53         k = k + data_to_use.loc[(i, wgt)]
54     else:
55         k = k + 1
56     i = i-1
57 mass_vals[-1] = mass_vals[-1] + k

```

Finally, we add the `mass_vals` information to `binned`, calculate the sample density, and return the result.


```
58  binned["mass"] = mass_vals
59  binned["freq"] = binned.mass / binned.mass.sum()
60  binned["dens"] = binned.freq / binned.width
61  return binned
```

File IV

estimate_parameters.py

This file contains the code for estimating parameters. Big picture, this file is where the heavy lifting happens. Contents include: (1) the code for calculating the Kolmogorov-Smirnov statistic/objective function, (2) functions to estimate parameters for the models of interest, and (3) several other functions associated with different models. The setup will be as follows:

1. General setup such as loading modules and defining pointers.
2. Define the general functions we need. These are the method-of-moment estimators for the Fisk distribution and the Kolmogorov-Smirnov objective function calculation.
3. Define the estimation functions and other functions associated with specific models. This will be where the numerics happen and is the bulk of the code in this file.

In this section, we outline the structure of the file and functions encountered. See also Table 6 for a list of models used in this section.

We use a variety of functions in the estimation process. In general, the function argument `data` should be a DataFrame, and `var` should be a column name in the DataFrame that serves as the variable of interest. We set `wgt` to be `None` by default, but if present, it should be the name of a column of survey weights. Previous versions of this file assumed that `var` was "income" and `wgt` was "weight", but this approach provides a more portable and general interface that we can use for other datasets if need be.

Calculating the Kolmogorov-Smirnov statistic can be computationally intensive if we aren't careful. We will write all `cdf_` functions such that they accept a set of x -values in the form of a numpy array. Then we calculate the distribution using vectorized operations, which is much faster than using `Series.map()`. The `len_data` argument will store the length of the data array.

For each model, we estimate parameters in a multi-step process. Most parameters are easier to find if we know c , some much easier, so we estimate as follows:

Table 6: Candidate Distributions

Distribution (Shifted)	Parameters	How Many?	String Key
Generalized Beta, Type II	α, β, p, q, c	5	GB2
Dagum	α, β, p, c	4	Dagum
Burr	α, β, q, c	4	Burr
Fisk	α, β, c	3	Fisk
Inverse Gamma	α, β, c	3	InvG
Davis	α, β , and c	3	Davis
Constant-Shift Inverse Gamma	α, β	2	CS_InvG
Constant-Shift-Shape Inverse Gamma	α	smallest \rightarrow 1	CSS_InvG
Log-Normal/Pareto Cutoff	$\mu, \sigma^2, \alpha, x_0, k, c$	6	LogN_P_cut
Log-Normal/Pareto Mixture	$\mu, \sigma^2, \alpha, x_0, \gamma, c$	6	LogN_P_mix

General functions:

- **validate_var_wgt**($\langle data \rangle, \langle var \rangle, \langle wgt \rangle$)—checks that *data* is a DataFrame, if *var* is a column in *data*, and whether *wgt* is a column in *data* if present.
- **make_ecdf**($\langle data \rangle, \langle var \rangle, \langle wgt \rangle = \text{None}$)—creates a DataFrame containing the empirical cumulative distribution of the *data*.
- **kolmogorov_smirnov**($\langle cdf \rangle, *, \langle ecdf \rangle = \text{None}, \langle x \rangle = \text{None}, \langle y \rangle = \text{None}, \langle data \rangle = \text{None}, \langle var \rangle = \text{None}, \langle wgt \rangle = \text{None}$)—calculates the value of the objective function based on the Kolmogorov-Smirnov statistic, i.e. $\sup |F(x) - G(x)|$. Here *cdf* should be a (symbolic) cumulative distribution function. To evaluate the fit to data, the function accepts either an empirical cumulative distribution *ecdf* or a DataFrame *data*. In the *ecdf* case, the arguments *x* and *y* should be column names in *ecdf* corresponding to *x* and *y*-values, and in the *data* case, the *var* and *wgt* arguments work the same way as in other places. Specifying *ecdf* takes precedence over *data*.
- **fisk_moments**($\langle data \rangle, \langle var \rangle, \langle wgt \rangle = \text{None}$)—function that accepts data and returns the moment estimators for the Fisk distribution. Used as the initial guess for various estimators.

1. Assume some value of *c*. Transform incomes according to this test value of *c* (subtract *c* from all incomes), and calculate the other parameters conditional on this test value for *c*. In other words, apply the estimator for unshifted data to the set of transformed incomes.
2. In all cases, this calculation will involve some sort of maximum likelihood for the remaining parameters. For Burr, Dagum, and Fisk, we maximize the likelihood purely

Model-specific functions:

- **cdf_** $\langle dist \rangle(\langle x \rangle, \langle params \rangle, \langle len_data \rangle = 1)$ —function that accepts a value of *x* and a list of parameters and returns a cumulative distribution value. This function should treat *data* as a numpy array and vectorize the operations to find the cumulative distribution.
- **density_** $\langle dist \rangle(\langle x \rangle, \langle params \rangle)$ —function that accepts a value of *x* and a list of parameters and returns a density value.
- **L_** $\langle dist \rangle(\langle data \rangle, \langle params \rangle, \langle var \rangle, \langle wgt \rangle = \text{None})$ —function that accepts a DataFrame and a list of parameters and returns a log-likelihood value.
- **estimate_** $\langle dist \rangle(\langle data \rangle, \langle var \rangle, \langle wgt \rangle = \text{None})$ —for each distribution, the helper function that estimates parameters; accepts a DataFrame with 2 columns where one column is income and one column is normalized weights. THERE WILL also be other functions involved internally in the estimation for each distribution. For example, each model has a corresponding **estimate_** $\langle dist \rangle_fit_for_c$ function that calculates the value of the objective function.

numerically. For `Inv_G` and `LogN_P_cut`, we numerically solve a single equation that arises from algebraic calculation. For `GB2`, `Davis`, and `LogN_P_mix`, the estimation is more complicated.

3. Evaluate the test c and other parameters using the objective function derived from the Kolmogorov-Smirnov statistic. Numerically find the value of c that minimizes the objective function.
4. We may have to calculate other parameters in the same step as c . For example, the `LogN_P_cut` distribution has a cutoff k , and we simultaneously fix test values of c and k . Then we calculate the optimal c - k pair that minimizes the objective function.

The estimation functions in this file return dictionary containing the value of the objective function under the key `fit` and a list of parameter estimates under the key `parameters`.

Setup

First import Pandas, Numpy, and two Scipy submodules.

```
1 import numpy          as np
2 import pandas         as pd
3 import scipy.optimize as opt
4 import scipy.special  as spec
5 import time
```

Some debugging stuff. In previous versions of this code, numpy was giving a `VisibleDeprecationWarning`. We also disable the `SettingWithCopyWarning` from Pandas.

```
6 #np.warnings.filterwarnings('error', category=np.VisibleDeprecationWarning)
7 pd.options.mode.chained_assignment = None
8 def the_time():
9     print("The time is", time.asctime())
```

New pointers to important Numpy and Scipy variables.

```
10 pi    = np.pi
11 e      = np.e
12 exp    = np.exp
13 floor  = np.floor
14 log    = np.log
15 sqrt   = np.sqrt
16 min    = opt.minimize
17 min_s  = opt.minimize_scalar
18 root   = opt.root_scalar
19 B      = spec.beta
20 erf    = spec.erf
21 G      = spec.gamma
22 I      = spec.betainc
23 Phi    = spec.ndtr
24 Phinv  = spec.ndtri
25 psi    = spec.digamma
26 Q      = spec.gammaincc
27 def psi1(x):
```

```

28 return spec.polygamma(1,x)
29 def Phi_prime(x):
30     return exp(-x**2 / 2) / sqrt(2 * pi)
31 def zeta(x):
32     return 1 + sp.zetac(x)
33 def zeta(x):
34     return 1 + spec.zetac(x)

```

We use the the approximation for $\log \Gamma(z)$ from Chen (2016) when z is large to avoid any numerical problems.¹ The approximation is

$$\log \Gamma(x+1) \approx \log \sqrt{2\pi x} + x(\log x - 1) + \left(x^2 + \frac{53}{210}\right) \log \left(1 + \frac{1}{12x^3 + \frac{24}{7}x - \frac{1}{2}}\right)$$

When $x \geq 10$, this approximation is accurate to one part in a billion. It's probably easier (better?) to use the built-in gamma function in Scipy, but I didn't think about that when I wrote the code. I am confident that for practical purposes, the improvement in accuracy is at best tiny, so it wouldn't gain much to rewrite the code. Plus this way we get to use a fun approximation theorem.

```

35 Gamma_frac1 = 53/210
36 Gamma_frac2 = 24/7
37 def log_G(z):
38     return log(G(z))
39 def log_G_approx(z):
40     x = z - 1
41     term1 = 0.5 * log(2 * pi * x)
42     term2 = x * (log(x) - 1)
43     term3 = (x * x + Gamma_frac1)
44     term4 = log(1 + 1 / (12 * (x * x * x) + Gamma_frac2 * (x) - 0.5))
45     return term1 + term2 + term3 * term4

```

Now create several empty dictionaries. We will fill them up as we get to them.

```

46 distribution = {} # cumulative distribution functions
47 density      = {} # densities
48 likelihood   = {} # likelihoods
49 estimator    = {} # estimators

```

General Functions

The first function we code is to validate input. This function does the same thing as the error checking at the beginning of `bin.py`.

```

50 def validate_var_wgt(data, var, wgt):
51     if not isinstance(data, pd.DataFrame):
52         raise ValueError("First argument of function should be a DataFrame")
53     if var not in data:

```

¹Chen, Chao-Ping. 2016. "A More Accurate Approximation for the Gamma Function." *Journal of Number Theory* 164: 417–428.

```

54     raise KeyError("Variable to analyze is not a column in the data")
55     if wgt and wgt not in data:
56         raise KeyError("Weight variable is not a column in the data")

```

We will call `validate_var_wgt` at the start of each estimating function.

We need a function to make the empirical cumulative distribution. The result is a `DataFrame` with columns `x` and `y` containing x and y -values respectively. Because of the discrete nature of data, we will repeat each x -value twice, and the final empirical cumulative distribution `DataFrame` will contain two rows for each x -value. For a given x , the first corresponding y -value will be $F(x^-)$, and the second y -value will be $F(x^+)$, where F denotes the empirical cumulative distribution. This means the return `temp DataFrame` value will be twice as long as the original data.

```

57 def make_ecdf(data, var, wgt=None):
58     validate_var_wgt(data, var, wgt)
59     temp = data.copy()
60     if not wgt:
61         wgt = "y"
62         if var == "y":
63             wgt = "y1"
64         temp[wgt] = 1

```

Make sure everything is sorted before we do stuff with it.

```

65     if not temp[var].is_monotonic_increasing:
66         temp = temp.sort_values(var).reset_index(drop=True)

```

We have to calculate the total weight given to unique x -values. (In other words, account for repeated x -values.) The way to do this is to group by `var` and then reset the index. Note that the concatenation uses `axis=0`, i.e. is vertical concatenation.

```

67     temp = temp[[var, wgt]].groupby(var).sum() # get unique var values
68     temp.index.set_names(None, inplace=True)
69     temp = temp.cumsum() / temp[wgt].sum()      # turn wgt into cdf
70     temp[var] = temp.index                      # add var vals back into data
71     temp = pd.concat([temp, temp]).sort_values(var).reset_index(drop=True)

```

Now shift the weight column, manually add the weight in the zeroth row, and return.

```

72     temp[wgt] = temp[wgt].rename(lambda x: x + 1)
73     temp.loc[(0, wgt)] = 0
74     return temp

```

The resulting x -column will have name `var`, and the y column will have `wgt` as its name if specified in `make_ecdf`. Otherwise, it will be called `"y"` or `"y1"`.

We code the objective function to calculate the Kolmogorov-Smirnov statistic. For two distribution functions F and G , the statistic is

$$\sup |F(x) - G(x)| = \|F - G\|_\infty$$

In our case, F is the cumulative distribution under the model of interest, and G is the empirical cumulative distribution. We begin by validating the arguments and making sure the user specified the required information.

```

75 def kolmogorov_smirnov(cdf, *, ecdf=None, x=None, y=None,

```

```

76         data=None, var=None, wgt=None):
77     if isinstance(ecdf, type(None)) and isinstance(data, type(None)):
78         raise ValueError("Please specify ecdf or data to use Kolmogorov-Smirnov")
79     elif isinstance(ecdf, type(None)) and not isinstance(data, type(None)):
80         validate_var_wgt(data, var, wgt)
81         ecdf = make_ecdf(data, var, wgt)
82         x = var
83         if not wgt and var != "y":
84             y = "y"
85         elif not wgt and var == "y":
86             y = "y1"
87         elif wgt:
88             y = wgt
89     elif not isinstance(ecdf, type(None)):
90         validate_var_wgt(ecdf, x, None)
91         validate_var_wgt(ecdf, y, None)
92         ecdf = ecdf.copy()
93     else:
94         raise RuntimeError("Something weird happened")

```

Now we calculate and return the statistic.

```

95     #ecdf = ecdf[(ecdf[y] > 0) & (ecdf[y] < 1)]
96     vals = cdf(np.array(ecdf[x])) # model cdf values
97     emp_vals = np.array(ecdf[y]) # empirical cdf values

```

We can add weights for a weighted Kolmogorov-Smirnov statistic.

```

98     #avg = (vals + ecdf[y]) / 2
99     #weights = (1 / (avg * (1-avg))) ** (1/2)
100     #weights = 1
101     return np.max(abs(vals - emp_vals))

```

Before getting into individual distributions, we also code the Fisk moment estimators. The function `fisk_moments` accepts a DataFrame containing columns `income` and `weights` and calculates the moment estimators for the Fisk distribution as suggested by Graf and Nedyalkova (2014). The Fisk distribution has shape parameter α and scale parameter β , and the estimators are

$$\alpha = \frac{\pi}{\sqrt{3v}}$$

$$\beta = e^m$$

where

$$m = \frac{1}{n} \sum_{i=1}^n w_i \log(x_i) \qquad v = \frac{1}{n} \sum_{i=1}^n w_i (\log(x_i) - m)^2$$

Here n is the total normalized survey weights (total number of observations), and w_i is the (normalized) survey weight on respondent i . The function stores the length of the data in `n`, calculates the intermediate quantities, and returns a 2-list containing the estimates.

```

102 def fisk_moments(data, var, wgt=None):
103     if wgt:
104         n = data.weight.sum()
105         m = (1/n) * (data[wgt] * log(data[var])).sum()
106         v = (1/n) * (data[wgt] * (log(data[var]) - m)**2).sum()
107     else:
108         n = len(data)
109         m = (1/n) * log(data[var]).sum()
110         v = (1/n) * ((log(data[var]) - m)**2).sum()
111     a = pi / sqrt(3 * v)
112     b = exp(m)
113     return [a, b]

```

We will use this function in GB2, Burr, Dagum, and Fisk estimation.

Specific Distributions

GB2 Generalized Beta, Type 2. Parameters:

- $\alpha > 0$: shape. It's unclear if this parameter needs to be strictly positive or just non-zero. Same with the analogous parameter in child distribution such as Burr and Dagum. Wikipedia has positive values for α , but McDonald and Ransom (2008) say that α can also be negative. If $\alpha < 0$, we should replace α in the numerator by $|\alpha| = -\alpha$. In any event, testing and inspection suggest that, at least for GB2, $\alpha > 0$, so I think it's probably fine to use this formula.
- $\beta > 0$: scale
- $p > 0$: shape
- $q > 0$: shape
- c : shift (for shifted GB2)

The cumulative distribution is

$$F(x) = I\left(\frac{\xi}{1+\xi}; p, q\right)$$

where I is the regularized incomplete beta function and

$$\xi = \left(\frac{x-c}{\beta}\right)^\alpha$$

```

114 def cdf_GB2(x, params):
115     a, b, p, q, c = params
116     x = x.astype(np.float64)
117     mask = x <= c
118     x[mask] = 0
119     frac = exp(-a * (log(x[~mask] - c) - log(b)))
120     x[~mask] = I(p, q, 1/(1+frac))
121     return x

```


The density is

$$y = \frac{\alpha (x - c)^{\alpha p - 1}}{\beta^{\alpha p} B(p, q) \left(1 + \left(\frac{x - c}{\beta} \right)^{\alpha} \right)^{p + q}}$$

where B is the beta function. We rewrite the numerator and denominator in terms of a logarithm in case any factors are particularly small or particularly large. The beta function can be difficult for a computer to handle numerically depending on how it is implemented. We use the previous approximation for $\log \Gamma(z)$ to calculate $\log B(p, q)$ (specifically $\log \Gamma(p)$, $\log \Gamma(q)$, and $\log \Gamma(p + q)$) without any risk of numerical complications. Code:

```
122 def density_GB2(x, params):
123     a, b, p, q, c = params
```

We manually check if p , q , and $p + q$ are each less than 10. If yes, we can use the built-in gamma function, and otherwise, we approximate.

```
124     if p < 10:
125         log_G_p = log_G(p)
126     else:
127         log_G_p = log_G_approx(p)
128     if q < 10:
129         log_G_q = log_G(q)
130     else:
131         log_G_q = log_G_approx(q)
132     if p + q < 10:
133         log_G_pq = log_G(p + q)
134     else:
135         log_G_pq = log_G_approx(p + q)
```

And now calculate the density.

```
136         if x <= c:
137             return 0
138         else:
139             num = log(a) + (a*p-1) * log(x-c)           # numerator
140             denom = ( (a*p) * log(b) + log_G_p
141                     + log_G_q - log_G_pq
142                     + (p+q) * log(1+((x-c)/b)**a)      ) # denominator
143             return exp(num - denom)
```

The likelihood is

$$L = n \log \alpha + (\alpha p - 1) \sum_{i=1}^n w_i \log(x_i - c) - n \alpha p \log \beta \\ - n \log(B(p, q)) - (p + q) \sum_{i=1}^n w_i \log \left(1 + \left(\frac{x_i - c}{\beta} \right)^{\alpha} \right)$$

where w_i is the survey weights. The `L_GB2` calculates the likelihood under `params` with respect to the `data`. We calculate each term separately and return their sum. Code:

```

144 def L_GB2(data, params, var, wgt=None):
145     validate_var_wgt(data, var, wgt)
146     a, b, p, q, c = params

    Again we manually check if  $p$ ,  $q$ , and  $p + q$  are each less than 10.

147     if p < 10:
148         log_G_p = log_G(p)
149     else:
150         log_G_p = log_G_approx(p)
151     if q < 10:
152         log_G_q = log_G(q)
153     else:
154         log_G_q = log_G_approx(q)
155     if p + q < 10:
156         log_G_pq = log_G(p + q)
157     else:
158         log_G_pq = log_G_approx(p + q)

    And calculate the likelihood.

159     if wgt:
160         n = data[wgt].sum()
161         term1 = n * log(a)
162         term2 = (a*p - 1) * (data[wgt] * log(data[var] - c)).sum()
163         term3 = -n * a * p * log(b)
164         term4 = -n * (log_G_p + log_G_q - log_G_pq)
165         term5 = -(p+q) * (data[wgt] * log(1 + ((data[var] - c)/b)**a)).sum()
166     else:
167         n = len(data)
168         term1 = n * log(a)
169         term2 = (a*p - 1) * log(data[var] - c).sum()
170         term3 = -n * a * p * log(b)
171         term4 = -n * (log_G_p + log_G_q - log_G_pq)
172         term5 = -(p+q) * log(1 + ((data[var] - c)/b)**a).sum()
173     return term1 + term2 + term3 + term4 + term5

```

And the unshifted version is

```

174 def L_GB2_unshift(data, params, var, wgt=None):
175     a, b, p, q = params
176     return L_GB2(data, [a, b, p, q, 0], var, wgt)

```

Differentiating the likelihood function gives us

$$\frac{dL}{d\alpha} = \frac{n}{\alpha} + p \sum_{i=1}^n w_i \log(x_i - c) - np \log \beta - (p + q) \sum_{i=1}^n w_i \frac{\left(\frac{x_i - c}{\beta}\right)^\alpha \log\left(\frac{x_i - c}{\beta}\right)}{1 + \left(\frac{x_i - c}{\beta}\right)^\alpha}$$

$$\begin{aligned}
&= \frac{n}{\alpha} + p \sum_{i=1}^n w_i \log \left(\frac{x_i - c}{\beta} \right) - (p + q) \sum_{i=1}^n \frac{w_i \log \left(\frac{x_i - c}{\beta} \right)}{1 + \left(\frac{\beta}{x_i - c} \right)^\alpha} \\
\frac{dL}{d\beta} &= -\frac{n\alpha p}{\beta} - (p + q) \sum_{i=1}^n w_i \frac{(x_i - c)^\alpha (-\alpha) \beta^{-\alpha-1}}{1 + \left(\frac{x_i - c}{\beta} \right)^\alpha} \\
&= -\frac{n\alpha p}{\beta} - (p + q) \frac{\alpha}{\beta} \sum_{i=1}^n \frac{w_i}{1 + \left(\frac{\beta}{x_i - c} \right)^\alpha}
\end{aligned}$$

Setting both expressions to 0 and solving for p and q gives us

$$\begin{aligned}
p &= \frac{1}{\alpha} \left(\frac{1}{n} \sum_{i=1}^n \frac{w_i}{1 + \left(\frac{\beta}{x_i - c} \right)^\alpha} \right) \left(\frac{1}{n} \sum_{i=1}^n \frac{w_i \log \left(\frac{x_i - c}{\beta} \right)}{1 + \left(\frac{\beta}{x_i - c} \right)^\alpha} \right. \\
&\quad \left. - \left[\frac{1}{n} \sum_{i=1}^n \frac{w_i}{1 + \left(\frac{\beta}{x_i - c} \right)^\alpha} \right] \left[\frac{1}{n} \sum_{i=1}^n w_i \log \left(\frac{x_i - c}{\beta} \right) \right] \right)^{-1} \\
q &= \frac{1}{\alpha} \left(1 - \frac{1}{n} \sum_{i=1}^n \frac{w_i}{1 + \left(\frac{\beta}{x_i - c} \right)^\alpha} \right) \left(\frac{1}{n} \sum_{i=1}^n \frac{w_i \log \left(\frac{x_i - c}{\beta} \right)}{1 + \left(\frac{\beta}{x_i - c} \right)^\alpha} \right. \\
&\quad \left. - \left[\frac{1}{n} \sum_{i=1}^n \frac{w_i}{1 + \left(\frac{\beta}{x_i - c} \right)^\alpha} \right] \left[\frac{1}{n} \sum_{i=1}^n w_i \log \left(\frac{x_i - c}{\beta} \right) \right] \right)^{-1}
\end{aligned}$$

so if we know α and β , we can find p and q . Ostensibly. Previous versions of this document made use of these equations for p and q , but that does not work. We want the values of p and q that make

$$\frac{dL}{dp} = \frac{dL}{dq} = 0,$$

not the values that cause α and β to be maxima. Going back to the likelihood function and differentiating again gives us

$$\begin{aligned}
\frac{dL}{dp} &= \alpha \sum_{i=1}^n w_i \log(x_i - c) - n\alpha \log \beta - n[\psi(p) - \psi(p + q)] \\
&\quad - \sum_{i=1}^n w_i \log \left(1 + \left(\frac{x_i - c}{\beta} \right)^\alpha \right) \\
&= \alpha \sum_{i=1}^n w_i \log \left(\frac{x_i - c}{\beta} \right) - n[\psi(p) - \psi(p + q)] - \sum_{i=1}^n w_i \log \left(1 + \left(\frac{x_i - c}{\beta} \right)^\alpha \right)
\end{aligned}$$

$$\frac{dL}{dq} = -n[\psi(q) - \psi(p+q)] - \sum_{i=1}^n w_i \log \left(1 + \left(\frac{x_i - c}{\beta} \right)^\alpha \right),$$

where ψ is the digamma function. It follows that we are looking for p and q such that

$$\begin{aligned} \psi(p+q) - \psi(p) &= \frac{1}{n} \sum_{i=1}^n w_i \log \left(1 + \left(\frac{x_i - c}{\beta} \right)^\alpha \right) - \frac{\alpha}{n} \sum_{i=1}^n w_i \log \left(\frac{x_i - c}{\beta} \right) \\ \psi(p+q) - \psi(q) &= \frac{1}{n} \sum_{i=1}^n w_i \log \left(1 + \left(\frac{x_i - c}{\beta} \right)^\alpha \right) \end{aligned}$$

Numerically solving for p and q is more complicated than using the previous equations that give p and q exactly—hence why previous versions of this document did not do this. The `find_pq` function will perform the numerical root-finding. The `log_term1` is the sum that appears in both equations, and the `log_term2` is the sum from the $\psi(p)$ condition. We assume that the `data` has been standardized, so we can ignore β and c .

```

177 def estimate_GB2_find_pq(data, alpha, var, wgt):
178     if wgt:
179         n = data[wgt].sum()
180         log_sum1 = (data[wgt] * log(1 + data[var] ** alpha)).sum() / n
181         log_sum2 = alpha * (data[wgt] * log(data[var])).sum() / n
182     else:
183         n = len(data)
184         log_sum1 = log(1 + data[var] ** alpha).sum() / n
185         log_sum2 = alpha * log(data[var]).sum() / n

```

Now code the function to solve (`mle_foc`) and its Jacobian.

```

186 def mle_foc(p_and_q):
187     p, q = p_and_q
188     p_term = psi(p + q) - psi(p) - log_sum1 + log_sum2
189     q_term = psi(p + q) - psi(q) - log_sum1
190     return [float(p_term), float(q_term)]
191 def d_mle_foc(p_and_q):
192     p, q = p_and_q
193     return [[psi1(p + q) - psi1(p), psi1(p + q)
194             [psi1(p + q), psi1(p + q) - psi1(q)]]

```

We find p and q and return. The initial guess $p = 2$ and $q = 0.5$ seems to work well.

```

195 return opt.root(mle_foc, x0=[3, 1], jac=d_mle_foc).x

```

The `find_a` function will return the value of α that maximizes the likelihood using the p and q estimates conditional on α . The `test_L` function will be our objective function to maximize.

```

196 def estimate_GB2_find_a(data, var, wgt):
197     def test_L(a):

```

```

198     p, q = estimate_GB2_find_pq(data, a, var, wgt)
199     return L_GB2(data, [a, 1, p, q, 0], var, wgt)
200 alpha = min(lambda x: -test_L(x), x0=1, method="Nelder-Mead",
201             bounds=[(0.1, None)]).x[0]
202 p, q = estimate_GB2_find_pq(data, alpha, var, wgt)
203 return [alpha, p, q]

```

The `fit_for_c` function accepts c and β values. It standardizes the data, finds remaining parameters using `estimate_GB2_find_a`, and returns the objective function values and parameters in a dictionary.

```

204 def estimate_GB2_fit_for_c(data, c, b, var, wgt, ecdf, x, y):
205     stand_data = data.copy()
206     stand_data[var] = (stand_data[var] - c) / b
207     stand_data = stand_data[stand_data.income > 0]
208     #print("finding a, p, q:", time.asctime())
209     a, p, q = estimate_GB2_find_a(stand_data, var, wgt)

```

The `fit` is the objective function, and the `parameters` are the estimates.

```

210     #print("fitting model:", time.asctime())
211     def F(x):
212         return cdf_GB2(x, [a, b, p, q, c])
213     temp = {
214         "fit": kolmogorov_smirnov(F, ecdf=ecdf, x=x, y=y),
215         "parameters": [a, b, p, q, c]}
216     #print("done:", time.asctime())
217     return temp

```

Before we code the final estimating function, we need a housekeeping function to extract β and c values from the return value of the `fit_for_c` function.

```

218 def estimate_GB2_get_cb(dict):
219     return [dict["parameters"][-1], dict["parameters"][1]]

```

Finally we estimate β and c by minimizing the Kolmogorov-Smirnov statistic. Unfortunately, the `fit_for_c` function is not quite well-behaved (i.e. practically identifiable) enough to algorithmically minimize the objective function, at least not using the Nelder-Mead algorithm as coded in Scipy. Accordingly, we use a brute-force estimator to find the optimum (c, β) pair over three iterations. In the first iteration, we use a partition mesh size of \$500, and in the second iteration, we use a partition mesh size of \$100 around the ten best-fitting pairs from the first iteration. In the third iteration, we use a partition mesh size of \$10 around the ten best-fitting pairs from the second iteration.

```

220 def estimate_GB2(data, var, wgt=None, *, ecdf=None, x=None, y=None):
221     validate_var_wgt(data, var, wgt)
222     if isinstance(ecdf, type(None)):
223         ecdf = make_ecdf(data, var, wgt)
224         x = var

```

```

225     if wgt:
226         y = wgt
227     else:
228         if var == "y":
229             y = "y1"
230         else:
231             y = "y"

```

First iteration. The $\langle b \text{ or } c \rangle_vals$ lists contain the test values of the variables we are using, and the `best_fits` list stores corresponding values of the objective function.

```

232 b_vals = [10000 + 500*i for i in range(81)]
233 c_vals = [-15000 + 500*i for i in range(21)]
234 best_fits = []
235 the_time()
236 for b in b_vals:
237     for c in c_vals:
238         temp = estimate_GB2_fit_for_c(data, c, b, var, wgt, ecdf, x, y)

```

We modify `best_fits` dynamically. If the list has fewer than ten entries, we add the current results. Otherwise, we check whether the current result fits better than the last entry in `best_fits`. If no, we ignore. If yes, we insert it into `best_fits` at the first position where the current result fits worse than previous items in `best_fits`. In this way, we maintain a list of best-fitting parameter combinations ordered by the value of their objective function. The while-loop we're using now is probably not the fastest way to implement the ordered list. Oh well.

```

239     if len(best_fits) < 10 and temp["parameters"][0] > 0.1:
240         best_fits.append(temp)
241     else:

```

We make sure to exclude possibilities where the estimate of α is 0.1.

```

242         if temp["fit"] < best_fits[-1]["fit"] and temp["parameters"][0] > 0.1:
243             i = 0
244             while temp["fit"] > best_fits[i]["fit"]:
245                 i = i + 1
246             best_fits.insert(i, temp)
247             best_fits.pop()

```

For the second iteration, we store the (c, β) values in `pairs`. First extract the values from `best_fits`, and then we populate the grid of parameter values to test. We turn `pairs` into a set to avoid duplicate parameter pairs.

```

248 pairs = [estimate_GB2_get_cb(i) for i in best_fits]
249 pairs = [(p[0] - 500 + 100*i, p[1]) for p in pairs for i in range(11)]
250 pairs = [(p[0], p[1] - 500 + 100*i) for p in pairs for i in range(11)]
251 pairs = set(pairs)
252 best_fits = []

```

```

253 the_time()
254 for p in pairs:
255     temp = estimate_GB2_fit_for_c(data, *p, var, wgt, ecdf, x, y)
256     if len(best_fits) < 10 and temp["parameters"][0] > 0.1:
257         best_fits.append(temp)
258     elif temp["parameters"][0] > 0.1:
259         if temp["fit"] < best_fits[-1]["fit"] and temp["parameters"][0] > 0.1:
260             i = 0
261             while temp["fit"] > best_fits[i]["fit"]:
262                 i = i + 1
263             best_fits.insert(i, temp)
264             best_fits.pop()

```

Now do this again for the third iteration. This time, we save just the best-fitting entry, not the whole list.

```

265 pairs = [estimate_GB2_get_cb(i) for i in best_fits]
266 pairs = [(p[0] - 100 + 10*i, p[1]) for p in pairs for i in range(21)]
267 pairs = [(p[0], p[1] - 100 + 10*i) for p in pairs for i in range(21)]
268 pairs = set(pairs)
269 solution = {"fit":2}
270 the_time()
271 for p in pairs:
272     temp = estimate_GB2_fit_for_c(data, *p, var, wgt, ecdf, x, y)
273     if temp["fit"] < solution["fit"] and temp["parameters"][0] > 0.1:
274         solution = temp
275 the_time()
276 return solution

```

And add functions to dictionaries.

```

277 distribution["GB2"] = cdf_GB2
278 density["GB2"]      = density_GB2
279 likelihood["GB2"]   = L_GB2
280 estimator["GB2"]     = estimate_GB2

```

Done with GB2!

Dagum Dagum distribution. Parameters:

- $\alpha > 0$: shape.
- $\beta > 0$: scale.
- $p > 0$: shape.
- c : shift. (For shifted Dagum)

The cumulative distribution is

$$y = \left(1 + \left(\frac{x - c}{\beta}\right)^{-\alpha}\right)^{-p}$$

```

281 def cdf_Dagum(x, params):
282     a, b, p, c = params
283     x = x.astype(np.float64)
284     mask = x <= c
285     x[mask] = 0
286     frac = exp(-a * (log(x[~mask] - c) - log(b)))
287     x[~mask] = (1 + frac) ** (-p)
288     return x

```

The density is

$$y = \frac{\alpha p (x - c)^{\alpha p - 1}}{\beta^{\alpha p} \left(1 + \left(\frac{x - c}{\beta}\right)^{\alpha}\right)^{p+1}}$$

As with GB2, we calculate the numerator and denominator in logs to avoid any potential numerical issues.

```

289 def density_Dagum(x, params):
290     a, b, p, c = params
291     if x <= c:
292         return 0
293     else:
294         num = log(a) + log(p) + (a*p-1) * log(x-c)
295         denom = (a*p) * log(b) + (p+1) * log(1 + ((x-c)/b)**a)
296         return exp(num - denom)

```

The likelihood is

$$L = n \log \alpha + n \log p + (\alpha p - 1) \sum_{i=1}^n w_i \log(x_i - c) - \alpha p \log \beta - (p + 1) \sum_{i=1}^n w_i \log \left(1 + \left(\frac{x_i - c}{\beta}\right)^{\alpha}\right)$$

```

297 def L_Dagum(data, params, var, wgt=None):
298     validate_var_wgt(data, var, wgt)
299     a, b, p, c = params
300     if wgt:
301         n = data[wgt].sum()
302         term1 = n * log(a)
303         term2 = n * log(p)
304         term3 = (a*p-1) * (data[wgt] * log(data[var] - c)).sum()
305         term4 = -n * a * p * log(b)
306         term5 = -(p+1) * (data[wgt] * log(1 + ((data[var] - c)/b)**a)).sum()
307     else:
308         n = len(data)
309         term1 = n * log(a)
310         term2 = n * log(p)
311         term3 = (a*p-1) * log(data[var] - c).sum()

```



```

312     term4 = -n * a * p * log(b)
313     term5 = -(p+1) * log(1 + ((data[var] - c)/b)**a).sum()
314     return term1 + term2 + term3 + term4 + term5

```

Unshifted version is

```

315 def L_Dagum_unshift(data, params, var, wgt=None):
316     a, b, p = params
317     return L_Dagum(data, [a, b, p, 0], var, wgt)

```

Once again, we use the moment estimators of the Fisk distribution as the initial guess for the Dagum likelihood maximization.

```

318 def estimate_Dagum_initial(data, var, wgt):
319     a, b = fisk_moments(data, var, wgt)
320     return [a, b, 1]

```

The estimation function is similar to GB2.

```

321 def estimate_Dagum_unshift(data, var, wgt):
322     guess = estimate_Dagum_initial(data, var, wgt)
323     def neg_L(params):
324         return -L_Dagum_unshift(data, params, var, wgt)
325     sol = min(neg_L, x0=guess, method="Nelder-Mead")
326     a, b, p = sol.x
327     return [a, b, p]

```

The `fit_for_c` function takes a value of c , finds the parameter estimates for that c , calculates the Kolmogorov-Smirnov statistic, and returns everything in a dictionary. Same as GB2.

```

328 def estimate_Dagum_fit_for_c(data, c, var, wgt, ecdf, x, y):
329     shift_data = data.copy()
330     shift_data[var] = shift_data[var] - c
331     shift_data = shift_data[shift_data[var] > 0]
332     a, b, p = estimate_Dagum_unshift(shift_data, var, wgt)
333     def F(x):
334         return cdf_Dagum(x, [a, b, p, c])
335     temp = {
336         "fit": kolmogorov_smirnov(F, ecdf=ecdf, x=x, y=y),
337         "parameters": [a, b, p, c]}
338     return temp

```

Now we code the main estimation for Dagum distribution. Unlike with GB2, we can ask Python to numerically minimize the Kolmogorov-Smirnov value for us. This code uses Brent's method.

```

339 def estimate_Dagum(data, var, wgt=None, *, ecdf=None, x=None, y=None):
340     validate_var_wgt(data, var, wgt)
341     if isinstance(ecdf, type(None)):

```

```

342     ecdf = make_ecdf(data, var, wgt)
343     x = var
344     if wgt:
345         y = wgt
346     else:
347         if var == "y":
348             y = "y1"
349         else:
350             y = "y"
351     def check_c(c):
352         return estimate_Dagum_fit_for_c(data, c, var, wgt, ecdf, x, y)["fit"]
353     sol = min_s(check_c, bracket=[-20000,-5000], options={"xtol":1.4e-5})

```

I can't figure out a good/easy way to get the parameters without another function call.
Oh well, not a big deal.

```

354     return estimate_Dagum_fit_for_c(data, sol.x, var, wgt, ecdf, x, y)

```

And add the functions to the appropriate dictionaries.

```

355 distribution["Dagum"] = cdf_Dagum
356 density["Dagum"]      = density_Dagum
357 likelihood["Dagum"]   = L_Dagum
358 estimator["Dagum"]    = estimate_Dagum

```

Burr Burr (Singh-Maddala) distribution. Parameters:

- $\alpha > 0$: shape.
- $\beta > 0$: scale.
- $q > 0$: shape.
- c : shift (for shifted Burr)

The cumulative distribution is

```

359 def cdf_Burr(x, params):
360     a, b, q, c = params
361     x = x.astype(np.float64)
362     mask = x <= c
363     x[mask] = 0
364     frac = exp(a * (log(x[~mask] - c) - log(b)))
365     x[~mask] = 1 - (1 + frac) ** (-q)
366     return x

```

Density is

$$y = \frac{\alpha q (x - c)^{\alpha-1}}{\beta^{\alpha} \left(1 + \left(\frac{x-c}{\beta}\right)^{\alpha}\right)^{q+1}}$$

Code:

```

367 def density_Burr(x, params):

```

```

368 a, b, q, c = params
369 if x <= c:
370     return 0
371 else:
372     num = log(a) + log(q) + (a-1) * log(x-c)
373     denom = a * log(b) + (1+q) * log(1 + ((x-c)/b)**a)
374     return exp(num - denom)

```

Likelihood is

$$\begin{aligned}
 L = n \log \alpha + n \log q + (\alpha - 1) \sum_{i=1}^n w_i \log(x_i - c) \\
 - n \alpha \log \beta - (1 + q) \sum_{i=1}^n w_i \log \left(1 + \left(\frac{x_i - c}{\beta} \right)^\alpha \right)
 \end{aligned}$$

```

375 def L_Burr(data, params, var, wgt=None):
376     validate_var_wgt(data, var, wgt)
377     a, b, q, c = params
378     if wgt:
379         n = data[wgt].sum()
380         term1 = n * log(a)
381         term2 = n * log(q)
382         term3 = (a-1) * (data[wgt] * log(data[var] - c)).sum()
383         term4 = -n * a * log(b)
384         term5 = -(q+1) * (data[wgt] * log(1 + ((data[var] - c)/b)**a)).sum()
385     else:
386         n = len(data)
387         term1 = n * log(a)
388         term2 = n * log(q)
389         term3 = (a-1) * log(data[var] - c).sum()
390         term4 = -n * a * log(b)
391         term5 = -(q+1) * log(1 + ((data[var] - c)/b)**a).sum()
392     return term1 + term2 + term3 + term4 + term5

```

Unshifted version:

```

393 def L_Burr_unshift(data, params, var, wgt=None):
394     a, b, q = params
395     return L_Burr(data, [a, b, q, 0], var, wgt)

```

Again, use Fisk moment estimators.

```

396 def estimate_Burr_initial(data, var, wgt):
397     a, b = fisk_moments(data, var, wgt)
398     return [a, b, 1]

```

Same thing as Dagum. The next functions are almost exactly identical to Burr.

```

399 def estimate_Burr_unshift(data, var, wgt):

```

```

400 guess = estimate_Burr_initial(data, var, wgt)
401 def neg_L(params):
402     return -L_Burr_unshift(data, params, var, wgt)
403 sol = min(neg_L, x0=guess, method="Nelder-Mead")
404 return sol.x

```

Now code fit_for_c.

```

405 def estimate_Burr_fit_for_c(data, c, var, wgt, ecdf, x, y):
406     shift_data = data.copy()
407     shift_data[var] = shift_data[var] - c
408     shift_data = shift_data[shift_data[var] > 0]
409     a, b, q = estimate_Burr_unshift(shift_data, var, wgt)
410     def F(x):
411         return cdf_Burr(x, [a, b, q, c])
412     temp = {
413         "fit": kolmogorov_smirnov(F, ecdf=ecdf, x=x, y=y),
414         "parameters": [a, b, q, c]}
415     return temp

```

Main estimation function.

```

416 def estimate_Burr(data, var, wgt=None, *, ecdf=None, x=None, y=None):
417     validate_var_wgt(data, var, wgt)
418     if isinstance(ecdf, type(None)):
419         ecdf = make_ecdf(data, var, wgt)
420         x = var
421         if wgt:
422             y = wgt
423         else:
424             if var == "y":
425                 y = "y1"
426             else:
427                 y = "y"
428     def check_c(c):
429         return estimate_Burr_fit_for_c(data, c, var, wgt, ecdf, x, y)["fit"]
430     sol = min_s(check_c, bracket=[-20000, -5000], options={"xtol": 1.4e-5})
431     return estimate_Burr_fit_for_c(data, sol.x, var, wgt, ecdf, x, y)

```

And add the functions to the correct dictionaries.

```

432 distribution["Burr"] = cdf_Burr
433 density["Burr"]      = density_Burr
434 likelihood["Burr"]   = L_Burr
435 estimator["Burr"]    = estimate_Burr

```

And Burr is finished!

Fisk Fisk distribution. Parameters:

- $\alpha > 0$: shape.

- $\beta > 0$: scale.
- c : shift (for shifted Fisk)

The cumulative distribution is

$$y = \frac{1}{1 + \left(\frac{x-c}{\beta}\right)^{-\alpha}}$$

```

436 def cdf_Fisk(x, params):
437     a, b, c = params
438     x = x.astype(np.float64)
439     mask = x <= c
440     x[mask] = 0
441     frac = exp(-a * (log(x[~mask] - c) - log(b)))
442     x[~mask] = 1 / (1 + frac)
443     return x

```

The density is

$$y = \frac{\alpha(x-c)^{\alpha-1}}{\beta^{\alpha} \left(1 + \left(\frac{x-c}{\beta}\right)^{\alpha}\right)^2}$$

Code:

```

444 def density_Fisk(x, params):
445     a, b, c = params
446     if x <= c:
447         return 0
448     else:
449         num = log(a) + (a-1) * log(x-c)
450         denom = a * log(b) + 2 * log(1 + ((x-c)/b)**a)
451         return exp(num - denom)

```

Likelihood is

$$L = n \log \alpha + (\alpha - 1) \sum_{i=1}^n w_i \log(x_i - c) - n \alpha \log \beta - 2 \sum_{i=1}^n w_i \log \left(1 + \left(\frac{x_i - c}{\beta}\right)^{\alpha}\right)$$

Code:

```

452 def L_Fisk(data, params, var, wgt=None):
453     validate_var_wgt(data, var, wgt)
454     a, b, c = params
455     if wgt:
456         n = data[wgt].sum()
457         term1 = n * log(a)
458         term2 = (a-1) * (data[wgt] * log(data[var] - c)).sum()
459         term3 = -n * a * log(b)
460         term4 = -2 * (data[wgt] * log(1 + ((data[var]-c)/b)**a)).sum()

```

```

461 else:
462     n = len(data)
463     term1 = n * log(a)
464     term2 = (a-1) * log(data[var] - c).sum()
465     term3 = -n * a * log(b)
466     term4 = -2 * log(1 + ((data[var]-c)/b)**a).sum()
467     return term1 + term2 + term3 + term4
468 def L_Fisk_unshift(data, params, var, wgt=None):
469     a, b = params
470     return L_Fisk(data, [a, b, 0], var, wgt)

```

Use Fisk moment estimators to make initial guess for Fisk distribution. :)

```

471 estimate_Fisk_initial = fisk_moments

```

Functions are the same as previously.

```

472 def estimate_Fisk_unshift(data, var, wgt):
473     guess = estimate_Fisk_initial(data, var, wgt)
474     def neg_L(params):
475         return -L_Fisk_unshift(data, params, var, wgt)
476     sol = min(neg_L, x0=guess, method="Nelder-Mead")
477     return sol.x

```

The fit_for_c function.

```

478 def estimate_Fisk_fit_for_c(data, c, var, wgt, ecdf, x, y):
479     shift_data = data.copy()
480     shift_data[var] = shift_data[var] - c
481     shift_data = shift_data[shift_data[var] > 0]
482     a, b = estimate_Fisk_unshift(shift_data, var, wgt)
483     def F(x):
484         return cdf_Fisk(x, [a, b, c])
485     temp = {
486         "fit": kolmogorov_smirnov(F, ecdf=ecdf, x=x, y=y),
487         "parameters": [a, b, c]}
488     return temp

```

And main estimation function.

```

489 def estimate_Fisk(data, var, wgt=None, *, ecdf=None, x=None, y=None):
490     validate_var_wgt(data, var, wgt)
491     if isinstance(ecdf, type(None)):
492         ecdf = make_ecdf(data, var, wgt)
493     x = var
494     if wgt:
495         y = wgt
496     else:
497         if var == "y":
498             y = "y1"

```

```

499     else:
500         y = "y"
501     def check_c(c):
502         return estimate_Fisk_fit_for_c(data, c, var, wgt, ecdf, x, y)["fit"]
503     sol = min_s(check_c, bracket=[data[var].min() - 2000, data[var].min()],
504         options={"xtol":1.4e-5})
505     return estimate_Fisk_fit_for_c(data, sol.x, var, wgt, ecdf, x, y)

```

Add functions to the dictionaries.

```

506 distribution["Fisk"] = cdf_Fisk
507 density["Fisk"]      = density_Fisk
508 likelihood["Fisk"]   = L_Fisk
509 estimator["Fisk"]    = estimate_Fisk

```

Done with the GB2 family of models!

InvG Inverse-Gamma distribution. Parameters:

- $\alpha > 0$: shape.
- $\beta > 0$: scale.
- c : shift. (For shifted inverse gamma)

The cumulative distribution is

$$y = Q\left(\alpha, \frac{\beta}{x - c}\right),$$

where Q is the upper incomplete gamma function.

```

510 def cdf_InvG(x, params):
511     a, b, c = params
512     x = x.astype(np.float64)
513     mask = x <= c
514     x[mask] = 0
515     frac = exp(log(b) - log(x[~mask] - c))
516     x[~mask] = Q(a, frac)
517     return x

```

The density is

$$y = \frac{\beta^\alpha}{\Gamma(\alpha)} \frac{e^{-\frac{\beta}{x-c}}}{(x-c)^{1+\alpha}}$$

This model is much simpler than previous distributions.

```

518 def density_InvG(x, params):
519     a, b, c = params

```

Similar to GB2, we manually check if we need to approximate $\log \Gamma(\alpha)$.

```

520 if a < 10:
521     log_G_a = log_G(a)

```

```

522 else:
523     log_G_a = log_G_approx(a)
524 if x <= c:
525     return 0
526 else:
527     num = a * log(b) - b / (x-c)
528     denom = log_G_a + (1+a) * log(x-c)
529     return exp(num - denom)

```

Likelihood:

$$L = n\alpha \log \beta - n \log \Gamma(\alpha) - \beta \sum_{i=1}^n \frac{w_i}{x_i - c} - (1 + \alpha) \sum_{i=1}^n w_i \log(x_i - c)$$

```

530 def L_InvG(data, params, var, wgt=None):
531     validate_var_wgt(data, var, wgt)
532     a, b, c = params

```

Again manually check if we need to approximate $\log \Gamma(\alpha)$.

```

533 if a < 10:
534     log_G_a = log_G(a)
535 else:
536     log_G_a = log_G_approx(a)
537 if wgt:
538     n = data[wgt].sum()
539     term1 = n * a * log(b)
540     term2 = -n * log_G_a
541     term3 = -b * (data[wgt] / (data[var] - c)).sum()
542     term4 = -(1+a) * (data[wgt] * log(data[var] - c)).sum()
543 else:
544     n = len(data)
545     term1 = n * a * log(b)
546     term2 = -n * log_G_a
547     term3 = -b * (1 / (data.income - c)).sum()
548     term4 = -(1+a) * log(data.income - c).sum()
549 return term1 + term2 + term3 + term4
550 def L_InvG_unshift(data, params, var, wgt=None):
551     validate_var_wgt(data, var, wgt)
552     a, b = params
553     return L_InvG(data, [a, b, 0], var, wgt)

```

We can find the MLE in more-or-less closed form. Sort of.

$$\frac{dL}{d\alpha} = n \log \beta - n\psi(\alpha) - \sum_{i=1}^n w_i \log(x_i - c) = 0$$

$$\frac{dL}{d\beta} = n \frac{\alpha}{\beta} - \sum_{i=1}^n \frac{w_i}{x_i - c} = 0$$

$$\beta = \frac{\alpha}{\frac{1}{n} \sum_{i=1}^n \frac{w_i}{x_i - c}}$$

$$\log \alpha - \psi(\alpha) - \log \left(\frac{1}{n} \sum_{i=1}^n \frac{w_i}{x_i - c} \right) - \frac{1}{n} \sum_{i=1}^n w_i \log(x_i - c) = 0$$

If we know c , the last equation gives us α , and the second-to-last equation gives us β in terms of α . So we have to numerically find the positive real root of $\log x - \psi(x) = C$. We define a dummy function called `_numeric` inside the `unshift` estimator, and we solve it numerically. Then we find β as above.

```

554 def estimate_InvG_unshift(data, var, wgt):
555     if wgt:
556         n = data[wgt].sum()
557         recip_data = (1/n) * (data[wgt] / data[var]).sum()
558         log_data = (1/n) * (data[wgt] * log(data[var])).sum()
559     else:
560         n = len(data)
561         recip_data = (1/n) * (1 / data[var]).sum()
562         log_data = (1/n) * log(data[var]).sum()
563     def estimate_InvG_numeric(x):
564         return log(x) - psi(x) - log(recip_data) - log_data

```

Getting the bounds on α right for the numerical solver can be a little tricky. In most cases and certainly for all the final estimates, α will be at most 5 or 10, but occasionally, we want to find α and β where α ends up being very large. So we dynamically find bounds for a good `bracket` for α . We store the current value of the numerical expression in `temp` because then we can handle the case when (if) we (miraculously) stumble on the root when we're searching for good bounds.

```

565     temp = estimate_InvG_numeric(10)
566     if temp < 0:
567         temp1 = estimate_InvG_numeric(0.01)
568         a = root(estimate_InvG_numeric, bracket=[0.01, 10],
569                 fprime = lambda x: (1/x) - psi1(x)).root
570     elif temp == 0:
571         a = 10
572     else:
573         prev_bound = 10
574         curr_bound = 30
575         temp = estimate_InvG_numeric(curr_bound)
576         while temp > 0:
577             prev_bound = curr_bound
578             curr_bound = curr_bound + 20
579             temp = estimate_InvG_numeric(curr_bound)
580         if temp == 0:
581             a = curr_bound

```

```

582     else:
583         a = root(estimate_InvG_numeric, bracket=[prev_bound,curr_bound],
584                 fprime = lambda x: (1/x) - psi1(x)).root
585     b = a / recip_data
586     return [a, b]

```

Now the `fit_for_c` function.

```

587 def estimate_InvG_fit_for_c(data, c, var, wgt, ecdf, x, y):
588     shift_data = data.copy()
589     shift_data[var] = shift_data[var] - c
590     shift_data = shift_data[shift_data[var] > 0]
591     a, b = estimate_InvG_unshift(shift_data, var, wgt)
592     def F(x):
593         return cdf_InvG(x, [a, b, c])
594     temp = {
595         "fit": kolmogorov_smirnov(F, ecdf=ecdf, x=x, y=y),
596         "parameters": [a, b, c]}
597     return temp

```

And finally the unshifted estimator

```

598 def estimate_InvG(data, var, wgt=None, *, ecdf=None, x=None, y=None):
599     validate_var_wgt(data, var, wgt)
600     if isinstance(ecdf, type(None)):
601         ecdf = make_ecdf(data, var, wgt)
602         x = var
603         if wgt:
604             y = wgt
605         else:
606             if var == "y":
607                 y = "y1"
608             else:
609                 y = "y"
610     def check_c(c):
611         return estimate_InvG_fit_for_c(data, c, var, wgt, ecdf, x, y)["fit"]
612     sol = min_s(check_c, method="bounded", bounds=[-23000,-1000])
613     return estimate_InvG_fit_for_c(data, sol.x, var, wgt, ecdf, x, y)

```

And add to the dictionaries.

```

614 distribution["InvG"] = cdf_InvG
615 density["InvG"]      = density_InvG
616 likelihood["InvG"]   = L_InvG
617 estimator["InvG"]    = estimate_InvG

```

Davis Davis distribution. Parameters:

- $\alpha > 0$: shape
- $\beta > 0$: scale

- c : shift

This probability distribution hasn't gotten as much attention recently as others. It has similar appearance to the inverse-gamma distribution except that it puts more mass in the tail. The probability density is

$$f(x) = \frac{\beta^\alpha}{\Gamma(\alpha) \zeta(\alpha)} \left[\frac{1}{\left(e^{\frac{\beta}{x-c}} - 1\right) (x-c)^{1+\alpha}} \right]$$

We can't write the cumulative distribution neatly, but we can express it as an infinite sum in terms of the upper incomplete gamma function.

$$\begin{aligned} F(x) &= \frac{\beta^\alpha}{\Gamma(\alpha) \zeta(\alpha)} \int_c^x \frac{dt}{\left(e^{\frac{\beta}{t-c}} - 1\right) (t-c)^{1+\alpha}} \\ &= \frac{\beta^\alpha}{\Gamma(\alpha) \zeta(\alpha)} \int_{\frac{\beta}{x-c}}^{\frac{\beta}{x-c}} \frac{u^{1+\alpha}}{\beta^{1+\alpha} (e^u - 1)} \left(-\frac{\beta}{u^2}\right) du & u = \frac{\beta}{s-c} \\ &= \frac{1}{\Gamma(\alpha) \zeta(\alpha)} \int_{\frac{\beta}{x-c}}^{\infty} \frac{u^{\alpha-1}}{e^u - 1} du \\ &= \frac{1}{\Gamma(\alpha) \zeta(\alpha)} \int_{\frac{\beta}{x-c}}^{\infty} \frac{u^{\alpha-1} e^{-u}}{1 - e^{-u}} du. \end{aligned}$$

At this point, we can rewrite the denominator as a geometric series with ratio e^{-u} (since $0 < e^{-u} < 1$). If we write the lower limit as p and consider just the integral, we have

$$\begin{aligned} \int_p^\infty \frac{u^{\alpha-1} e^{-u}}{1 - e^{-u}} du &= \int_p^\infty u^{\alpha-1} e^{-u} \sum_{k=0}^{\infty} e^{-ku} du \\ &= \int_p^\infty u^{\alpha-1} \sum_{k=1}^{\infty} e^{-ku} du \\ &= \sum_{k=1}^{\infty} \int_p^\infty u^{\alpha-1} e^{-ku} du \\ &= \sum_{k=1}^{\infty} \int_{kp}^\infty \left(\frac{v}{k}\right)^{\alpha-1} e^{-v} \frac{dv}{k} & v = ku \\ &= \sum_{k=1}^{\infty} \frac{1}{k^\alpha} \int_{kp}^\infty v^{\alpha-1} e^{-v} dv. \end{aligned}$$

It follows that

$$F(x) = \frac{1}{\Gamma(\alpha) \zeta(\alpha)} \sum_{k=1}^{\infty} \frac{1}{k^\alpha} \Gamma\left(\alpha, \frac{k\beta}{x-c}\right) = \frac{1}{\zeta(\alpha)} \sum_{k=1}^{\infty} \frac{1}{k^\alpha} Q\left(\alpha, \frac{k\beta}{x-c}\right),$$

where Q is the regularized upper incomplete gamma function. What a fun distribution.

To code it, we start by writing a function to calculate the sum. We expect that the input x will have the form $\beta/(x_i - c)$, so we only need to provide α as an argument. (Note: we use `np.max` because x could be an array.)

```

618     def cdf_Davis_sum(x, alpha):
619         k = 1
620         temp = Q(alpha, x)
621         var = temp
622         while np.max(np.abs(temp)) > 0.000001:
623             k = k + 1
624             temp = exp(-alpha * log(k) + Q(alpha, exp(log(k) + log(x))))
625             var = var + temp
626         return var

```

And calculate the actual cumulative distribution value.

```

627     def cdf_Davis(x, params):
628         a, b, c = params
629         x = x.astype(float)
630         mask = (x <= c)
631         x[mask] = 0
632         x[~mask] = cdf_Davis_sum(b / (np.array(x[~mask]) - c), a)
633         return x / zeta(a)

```

Code for density.

```

634     def density_Davis(x, params):
635         a, b, c = params
636         if x <= c:
637             return 0
638         else:
639             num = a * log(b)
640             denom = log(G(a)) + log(zeta(a)) + \
641                 log(exp(exp(log(b) - log(x - c))) - 1) + \
642                 (1 + a) * log(x - c)
643             return exp(num - denom)

```

The likelihood is

$$L = n\alpha \log \beta - n \log \Gamma(\alpha) - n \log \zeta(\alpha) - \sum_{i=1}^n w_i \log \left(e^{\frac{\beta}{x_i - c}} - 1 \right) - (1 + \alpha) \sum_{i=1}^n w_i \log(x_i - c)$$

Code.

```

644     def L_Davis(data, params, var, wgt=None):
645         validate_var_wgt(data, var, wgt)
646         a, b, c = params
647         if isinstance(wgt, type(None)): # if no weight provided
648             n = len(data)
649             term4 = -log(exp(exp(log(b) - log(data[var] - c))) - 1).sum()

```

```

650     term5 = -(1 + a) * log(data[var] - c).sum()
651 else:                                     # if weight provided
652     n = data[wgt].sum()
653     term4 = -(data[wgt] *
654               log(exp(exp(log(b) - log(data[var] - c))) - 1)).sum()
655     term5 = -(1 + a) * (data[wgt] * log(data[var] - c)).sum()

```

The first three terms don't depend on the weights, only on the (weight or length) total n , so we can code them after the conditional block, which makes the code cleaner. (Even if the terms are out of order now.)

```

656     term1 = n * a * log(b)
657     term2 = -n * log(G(a))
658     term3 = -n * log(zeta(a))
659     return term1 + term2 + term3 + term4 + term5

```

Unshifted version is

```

660     def L_Davis_unshift(data, params, var, wgt=None):
661         return L_Davis(data, [*params, 0], var, wgt)

```

Differentiating the likelihood function and setting equal to 0 gives us

$$\frac{dL}{d\alpha} = n \log \beta - n\psi(\alpha) - n \frac{\zeta'(\alpha)}{\zeta(\alpha)} - \sum_{i=1}^n w_i \log(x_i - c) = 0$$

$$\frac{dL}{d\beta} = \frac{n\alpha}{\beta} - \sum_{i=1}^n \left(\frac{w_i}{x_i - c} \right) \frac{e^{\frac{\beta}{x_i - c}}}{e^{\frac{\beta}{x_i - c}} - 1} = 0$$

We would like to solve these equations numerically. First we need to code ζ' . We have

$$\zeta'(z) = - \sum_{k=2}^{\infty} \frac{\log k}{k^z}$$

Code:

```

662     def zeta_prime(z):
663         k = 2
664         temp = 1
665         val = -log(k) / k ** z
666         while np.abs(temp) > 0.000001:
667             k = k + 1
668             temp = -exp(log(log(k)) - z * log(k))
669             val = val + temp
670         return val

```

Simultaneously finding roots α and β for $dL/d\alpha$ and $dL/d\beta$ doesn't work because the numerics don't converge. However, we can still maximize the likelihood function. Notice that it is straightforward to find α that maximizes likelihood conditional on β and

c , and we write the `_a_from_b` function to do this.

```

671     def estimate_Davis_a_from_b(data, b, var, wgt=None):
672         if isinstance(wgt, type(None)):
673             n = len(data)
674             wgt = pd.Series(1, index=data.index)
675         else:
676             n = data[wgt].sum()
677             wgt = data[wgt]
```

The function `dL_da` is what we want to solve. Because it doesn't depend on α , we store the $\log \beta$ and $\log(x_i - c)$ terms in `temp` so we don't have to calculate their value every time the numerical solver evaluates `dL_da`. We've also divided everything by `n` to make things simpler.

```

678         temp = log(b) - (1/n) * (wgt * log(data[var])).sum()
679         def dL_da(a):
680             return temp - psi(a) - zeta_prime(a) / zeta(a)
```

The expression

$$-\psi(\alpha) - \frac{\zeta'(\alpha)}{\zeta(\alpha)}$$

is monotonic decreasing in α , and that makes it very easy to find a bracket for `dL_da`. This is the same way we found a bracket for α when we estimated the `InvG` model. We start with `lower_bound = 2`, and we repeatedly scale `lower_bound` up or down until we find `lower_bound` and `upper_bound` where `dL_da` has different signs.

```

681     lower_bound = 2
682     if dL_da(lower_bound) == 0:           # if lower_bound is a root
683         return lower_bound
684     elif dL_da(lower_bound) > 0:         # if lower_bound is too small
685         upper_bound = 1.2 * lower_bound
686         while dL_da(upper_bound) > 0:
687             lower_bound = upper_bound
688             upper_bound = 1.2 * upper_bound
```

If `lower_bound` is too big, we decrease it. We take roots instead of multiplying or dividing because we don't want to cross the asymptote of ζ at $\alpha = 1$.

```

689         elif dL_da(lower_bound) < 0:     # if lower_bound is too big
690             while dL_da(lower_bound) < 0:
691                 upper_bound = lower_bound
692                 lower_bound = lower_bound ** 0.8
693         else:
694             raise RuntimeError("Something weird happened")
695     return opt.root_scalar(dL_da, x0=2,
696                           bracket=[lower_bound, upper_bound]).root
```

Now we do the estimation. Same thing as usual. We have the `unshift` function to calculate α and β given c , the `fit_for_c` function to find the Kolmogorov-Smirnov statistic, and the final estimation function to do the overall estimation. As initial guesses, we use α and β from the inverse-gamma model.

```
697     def estimate_Davis_unshift(data, var, wgt):
698         guess = estimate_InvG_unshift(data, var, wgt)
```

Because we know how to find α in terms of β , we can numerically maximize the likelihood as a function of a single argument β . This function is nicely behaved.

```
699     def neg_L_with_a_in_terms_of_b(b):
700         a = estimate_Davis_a_from_b(data, b, var, wgt)
701         return -L_Davis_unshift(data, [a,b], var, wgt)
702     sol = opt.minimize_scalar(neg_L_with_a_in_terms_of_b,
703                             bracket=[guess[1] - 20000, guess[1]])
704     return [estimate_Davis_a_from_b(data, sol.x, var, wgt), sol.x]
```

The `fit_for_c` function.

```
705     def estimate_Davis_fit_for_c(data, c, var, wgt, ecdf, x, y):
706         shift_data = data.copy()
707         shift_data[var] = data[var] - c
708         shift_data = shift_data[shift_data[var] > 0]
709         a, b = estimate_Davis_unshift(shift_data, var, wgt)
710         def F(x):
711             return cdf_Davis(x, [a, b, c])
712         temp = {
713             "fit": kolmogorov_smirnov(F, ecdf=ecdf, x=x, y=y),
714             "parameters": [a, b, c]}
715         return temp
```

And the main estimation function.

```
716     def estimate_Davis(data, var, wgt=None, *, ecdf=None, x=None, y=None):
717         validate_var_wgt(data, var, wgt)
718         if isinstance(ecdf, type(None)):
719             ecdf = make_ecdf(data, var, wgt)
720             x = var
721             if wgt:
722                 y = wgt
723             else:
724                 if var == "y":
725                     y = "y1"
726                 else:
727                     y = "y"
728         def check_c(c):
729             return estimate_Davis_fit_for_c(data, c, var, wgt, ecdf, x, y)["fit"]
730         sol = opt.minimize_scalar(check_c,
731                                 bracket=[-20000, -7000], options={"xtol": 1.4e-5})
732         return estimate_Davis_fit_for_c(data, sol.x, var, wgt, ecdf, x, y)
```

And add to the dictionaries.

```

733     distribution["Davis"] = cdf_Davis
734     density["Davis"]      = density_Davis
735     likelihood["Davis"]   = L_Davis
736     estimator["Davis"]    = estimate_Davis

```

CS_InvG We're not using this model, but I've kept it for historical reasons. Constant-shift inverse-gamma distribution. Parameters:

- $\alpha > 0$: shape parameter.
- β : scale parameter.

The distribution also has a constant ϕ . This is the same as shifted inverse-gamma distribution except that we set $c = \phi\beta$, so the shift is proportional to the scale. For our purposes, $\phi \approx -0.13$. Cumulative distribution is

$$y = Q\left(\alpha, \frac{\beta}{x - \phi\beta}\right).$$

```

737 def cdf_CS_InvG(x, phi, params):
738     a, b = params
739     shift = phi * b
740     x = x.astype(np.float64)
741     mask = x <= shift
742     x[mask] = 0
743     frac = exp(log(b) - log(x[~mask] - shift))
744     x[~mask] = Q(a, frac)
745     return x

```

Density:

$$y = \frac{\beta^\alpha}{\Gamma(\alpha)} \frac{e^{-\frac{\beta}{x - \phi\beta}}}{(x - \phi\beta)^{1+\alpha}}$$

```

746 def density_CS_InvG(x, phi, params):
747     a, b = params
748     shift = phi * b
749     if x <= shift:
750         return 0
751     else:
752         num = a * log(b) - b / (x - shift)
753         denom = log_G(a) + (1 + alpha) * log(x - shift)
754         return exp(num - denom)

```

Likelihood:

$$L = n\alpha \log \beta - n \log \Gamma(\alpha) - \beta \sum_{i=1}^n \frac{w_i}{x_i - \phi\beta} - (1 + \alpha) \sum_{i=1}^n w_i \log(x_i - \phi\beta)$$

```

755 def L_CS_InvG(data, phi, params, var, wgt=None):

```



```

756 validate_var_wgt(data, var, wgt)
757 a, b = params
758 shift = phi * b
759 if wgt:
760     n = data[wgt].sum()
761     term1 = n * a * log(b)
762     term2 = -n * log_G(a)
763     term3 = -b * (data[wgt] / (data[var] - shift)).sum()
764     term4 = -(1+a) * (data[wgt] * log(data[var] - shift)).sum()
765 else:
766     n = len(data)
767     term1 = n * a * log(b)
768     term2 = -n * log_G(a)
769     term3 = -b * (1 / (data[var] - shift)).sum()
770     term4 = -(1+a) * log(data[var] - shift).sum()
771 return term1 + term2 + term3 + term4

```

We can't solve the MLE for the constant-shift inverse-gamma distribution in closed form, but we can get close. Unfortunately, it's a bit more complicated than with the separate c parameter.

$$\begin{aligned}
\frac{dL}{d\alpha} &= n \log \beta - n\psi(\alpha) - \sum_{i=1}^n w_i \log(x_i - \phi\beta) \\
&= -n\psi(\alpha) + \sum_{i=1}^n w_i \log\left(\frac{\beta}{x_i - \phi\beta}\right) = 0 \\
\frac{dL}{d\beta} &= n\frac{\alpha}{\beta} - \sum_{i=1}^n w_i \left(\frac{1}{x_i - \phi\beta} + \frac{\phi\beta}{(x_i - \phi\beta)^2} \right) - (1 + \alpha) \sum_{i=1}^n \frac{-\phi w_i}{x_i - \phi\beta} \\
&= n\frac{\alpha}{\beta} - \sum_{i=1}^n \frac{w_i x_i}{(x_i - \phi\beta)^2} + \phi(1 + \alpha) \sum_{i=1}^n \frac{w_i}{x_i - \phi\beta} = 0
\end{aligned}$$

Looking at the $d\beta$ equation, we see that

$$\begin{aligned}
n\frac{\alpha}{\beta} + \phi(1 + \alpha) \sum_{i=1}^n \frac{w_i}{x_i - \phi\beta} &= \sum_{i=1}^n \frac{w_i x_i}{(x_i - \phi\beta)^2} \\
n\frac{\alpha}{\beta} + \phi\alpha \sum_{i=1}^n \frac{w_i}{x_i - \phi\beta} + \phi \sum_{i=1}^n \frac{w_i}{x_i - \phi\beta} &= \sum_{i=1}^n \frac{w_i x_i}{(x_i - \phi\beta)^2} \\
\frac{\alpha}{\beta} + \frac{\phi\alpha}{n} \sum_{i=1}^n \frac{w_i}{x_i - \phi\beta} &= \frac{1}{n} \sum_{i=1}^n \frac{w_i x_i}{(x_i - \phi\beta)^2} - \frac{\phi}{n} \sum_{i=1}^n \frac{w_i}{x_i - \phi\beta} \\
\alpha \left[\frac{1}{\beta} + \frac{\phi}{n} \sum_{i=1}^n \frac{w_i}{x_i - \phi\beta} \right] &= \frac{1 - \phi}{n} \sum_{i=1}^n \frac{w_i x_i}{(x_i - \phi\beta)^2} + \frac{\phi^2 \beta}{n} \sum_{i=1}^n \frac{w_i}{(x_i - \phi\beta)^2}
\end{aligned}$$

Altogether we get

$$\psi(\alpha) = \frac{1}{n} \sum_{i=1}^n w_i \log \left(\frac{\beta}{x_i - \phi\beta} \right)$$

$$\alpha = \frac{\frac{1-\phi}{n} \sum_{i=1}^n \frac{w_i x_i}{(x_i - \phi\beta)^2} + \frac{\phi^2 \beta}{n} \sum_{i=1}^n \frac{w_i}{(x_i - \phi\beta)^2}}{\frac{1}{\beta} + \frac{\phi}{n} \sum_{i=1}^n \frac{w_i}{x_i - \phi\beta}}$$

To do the estimation, we can solve this system of equations numerically for β . We begin with the two expressions of interest, which we call `alpha1` and `alpha2`. We assume that the data has not been modified at all, so we should translate it and remove all negative (after shifting) incomes. We store the new dataset in `shift_data`.

```

772 def estimate_CS_InvG_alpha1(data, phi, beta, var, wgt):
773     shift = phi * beta
774     shift_data = data.copy()
775     shift_data = shift_data[shift_data[var] > shift]
776     shift_data[var] = shift_data[var] - shift
777     if wgt:
778         n = data[wgt].sum()
779         sum = (shift_data[wgt] * log(beta / shift_data[var])).sum()
780     else:
781         n = len(data)
782         sum = log(beta / shift_data[var]).sum()
783     return (1/n) * sum

```

Same thing for the second α expression.

```

784 def estimate_CS_InvG_alpha2(data, phi, beta, var, wgt):
785     shift = phi * beta
786     shift_data = data.copy()
787     shift_data = shift_data[shift_data[var] > shift]
788     shift_data[var] = shift_data[var] - shift
789     if wgt:
790         n = data[wgt].sum()
791         num1 = (1-phi) / n * (shift_data[wgt] * (shift_data[var] + shift) /
792             shift_data[var] ** 2).sum()
793         num2 = (phi**2) * (beta/n) * (shift_data[wgt] /
794             shift_data[var] ** 2).sum()
795         denom = (1/beta) + (phi/n) * (shift_data[wgt] / shift_data[var]).sum()
796     else:
797         n = len(data)
798         num1 = (1-phi) / n * ((shift_data[var] + shift) /
799             shift_data[var] ** 2).sum()
800         num2 = (phi**2) * (beta/n) * (1 / shift_data[var] ** 2).sum()
801         denom = (1/beta) + (phi/n) * (1 / shift_data[var]).sum()
802     return (num1 + num2) / denom

```

Because we don't have a separate shift parameter, we can perform the estimation in a single step rather than two (or more). The `temp_numeric` function will be the expression we want to numerically find the root of with respect to β . Consider the asymptotics of both expressions for α . (Assuming that total weight has been normalized to n .) As β gets large, the sum approaches

$$\begin{aligned}\frac{1}{n} \sum_{i=1}^n w_i \log \left(\frac{\beta}{x_i - \phi\beta} \right) &= \frac{1}{n} \sum_{i=1}^n w_i \log \left(\frac{1}{\frac{x_i}{\beta} - \phi} \right) \\ &\rightarrow \frac{1}{n} \sum_{i=1}^n w_i \log \left(\frac{1}{\phi} \right) = \log \left(\frac{1}{\beta} \right).\end{aligned}$$

For the big fraction, we have

$$\begin{aligned}&\frac{\frac{1-\phi}{n} \sum_{i=1}^n \frac{w_i x_i}{(x_i - \phi\beta)^2} + \frac{\phi^2 \beta}{n} \sum_{i=1}^n \frac{w_i}{(x_i - \phi\beta)^2}}{\frac{1}{\beta} + \frac{\phi}{n} \sum_{i=1}^n \frac{w_i}{x_i - \phi\beta}} \\ &= \frac{\frac{1-\phi}{n} \sum_{i=1}^n \frac{w_i x_i}{\left(\frac{x_i}{\sqrt{\beta}} - \phi\sqrt{\beta}\right)^2} + \frac{\phi^2}{n} \sum_{i=1}^n \frac{w_i}{\left(\frac{x_i}{\beta} - \phi\right)^2}}{1 + \frac{\phi}{n} \sum_{i=1}^n \frac{w_i}{\frac{x_i}{\beta} - \phi}}\end{aligned}$$

The first term in the numerator approaches 0, and the second approaches 1. The denominator approaches 0, so the quotient blows up. It follows that the `temp_numeric` function in the estimator gets big negative for large β .

```

803 def estimate_CS_InvG(data, phi, var, wgt=None, *, ecdf=None, x=None, y=None):
804     validate_var_wgt(data, var, wgt)
805     if isinstance(ecdf, type(None)):
806         ecdf = make_ecdf(data, var, wgt)
807     x = var
808     if wgt:
809         y = wgt
810     else:
811         if var == "y":
812             y = "y1"
813         else:
814             y = "y"
815     def temp_numeric(b):
816         return estimate_CS_InvG_alpha1(data, phi, b, var, wgt) - \
817             psi(estimate_CS_InvG_alpha2(data, phi, b, var, wgt))

```

Around the root, the second term in `temp_numeric` grows faster than the first, but

this doesn't hold for small β . There the function jumps around a lot and may give bad results, i.e. roots that we do not care about. Because the β value jumps around a lot, we calculate bounds dynamically. This is very similar to what we did to find α in the estimation for InvG. In the limit as β gets large, `temp_numeric` gets big negative, so our approach will be to start with a large value of β and successively decrease it until we get bounds that work. Then we find α accordingly.

```

818 right_bound = 200000
819 left_bound = 100000
820 temp = temp_numeric(left_bound)
821 if temp > 0:
822     rtemp = temp_numeric(right_bound)
823     while rtemp > 0:
824         left_bound = right_bound
825         right_bound = right_bound * 1.05
826         rtemp = temp_numeric(right_bound)
827     if rtemp == 0:
828         beta = right_bound
829     else:
830         beta = root(temp_numeric, bracket=[left_bound, right_bound]).root
831 elif temp == 0:
832     beta = left_bound
833 else:
834     while temp < 0:
835         right_bound = left_bound
836         left_bound = left_bound * 0.98
837         temp = temp_numeric(left_bound)
838     if temp == 0:
839         beta = left_bound
840     else:
841         beta = root(temp_numeric, bracket=[left_bound, right_bound]).root
842 alpha = estimate_CS_InvG_alpha2(data, phi, beta, var, wgt)

```

Create and return the dictionary. Unlike dictionaries from the other estimators, this one contains a key called "phi".

```

843 def F(x):
844     return cdf_CS_InvG(x, phi, [alpha, beta])
845 temp = {
846     "fit": kolmogorov_smirnov(F, ecdf=ecdf, x=x, y=y),
847     "phi": phi,
848     "parameters": [alpha, beta]}
849 return temp

```

And add to dictionaries.

```

850 distribution["CS_InvG"] = cdf_CS_InvG
851 density["CS_InvG"]      = density_CS_InvG
852 likelihood["CS_InvG"]  = L_CS_InvG
853 estimator["CS_InvG"]   = estimate_CS_InvG

```

CSS_InvG Constant-shift-scale inverse-gamma distribution. Parameters:

- $\alpha > 0$: shape parameter

This distribution is the same as in inverse-gamma distribution except that we take values for β and c based on linear functions of α . Specifically, $c = \psi_0 + \psi_1 t + \psi_2 \alpha$, and $c = \phi \beta$. From `check_constants.py`, we approximate values of the constants as follows:

- ϕ : -0.134
- ψ_0 : 721,038
- ψ_1 : -363
- ψ_2 : -2252

See Table 7 for more information. For this set of functions, we will have several more arguments for the functions than usual. The cumulative distribution function is

$$y = Q\left(\alpha, \frac{1}{\phi} \frac{\psi_0 + \psi_1 t + \psi_2 \alpha}{x - (\psi_0 + \psi_1 t + \psi_2 \alpha)}\right)$$

Code:

```

854 def cdf_CSS_InvG(x, t, phi, psi, a):
855     psi0, psi1, psi2 = psi
856     b = (psi0 + psi1 * t + psi2 * a) / phi
857     c = psi0 + psi1 * t + psi2 * a
858     x = x.astype(np.float64)
859     mask = x <= c
860     x[mask] = 0
861     frac = exp(log(b) - log(x[~mask] - c))
862     x[~mask] = Q(a, frac)
863     return x

```

The density is

$$y = \frac{(\psi_0 + \psi_1 t + \psi_2 \alpha)^\alpha}{\Gamma(\alpha)} \frac{e^{-\frac{\psi_0 + \psi_1 t + \psi_2 \alpha}{x - \phi(\psi_0 + \psi_1 t + \psi_2 \alpha)}}}{(x - \phi(\psi_0 + \psi_1 t + \psi_2 \alpha))^{1+\alpha}}$$

```

864 def density_CSS_InvG(x, t, phi, psi, a):
865     psi0, psi1, psi2 = psi
866     beta = (psi0 + psi1 * t + psi2 * a) / phi
867     c = psi0 + psi1 * t + psi2 * a
868     if x <= c:
869         return 0
870     else:
871         return density_InvG(x, [a, beta, c])

```

Likelihood:

$$L = n\alpha \log(\psi_0 + \psi_1 t + \psi_2 \alpha) - n\Gamma(\alpha) - (\psi_0 + \psi_1 t + \psi_2 \alpha) \sum_{i=1}^n \frac{w_i}{x_i - \phi(\psi_0 + \psi_1 t + \psi_2 \alpha)}$$

$$-(1 + \alpha) \sum_{i=1}^n w_i \log[x_i - \phi(\psi_0 + \psi_1 t + \psi_2 \alpha)]$$

```

872 def L_CSS_InvG(data, t, phi, psi, a, var, wgt=None):
873     validate_var_wgt(data, var, wgt)
874     psi0, psi1, psi2 = psi
875     beta = (psi0 + psi1 * t + psi2 * a) / phi
876     c = psi0 + psi1 * t + psi2 * a
877     return L_InvG(data, [a, beta, c], var, wgt)

```

Unfortunately, maximum likelihood does not work properly with the constant-shift-scale inverse-gamma distribution. So we minimize the Kolmogorov-Smirnov statistic instead. The `fit_for_a` function will calculate the Kolmogorov-Smirnov statistic.

```

878 def estimate_CSS_InvG_fit_for_a(t, phi, psi, a, ecdf, x, y):
879     psi0, psi1, psi2 = psi
880     beta = (psi0 + psi1 * t + psi2 * a) / phi
881     c = psi0 + psi1 * t + psi2 * a
882     def F(x):
883         return cdf_InvG(x, [a, beta, c])
884     temp = {
885         "fit": kolmogorov_smirnov(F, ecdf=ecdf, x=x, y=y),
886         "phi": phi,
887         "psi": psi,
888         "parameters": [a]}
889     return temp

```

And the main estimation function calculates the minimum of the `fit_for_a` function. The `a0` argument is the starting value for α in the minimization. We will take `a0` to be the value of α from year `t` under the 3-parameter model.

```

890 def estimate_CSS_InvG(data, t, phi, psi, a0, var, wgt=None, *, \
891                       ecdf=None, x=None, y=None):
892     validate_var_wgt(data, var, wgt)
893     if isinstance(ecdf, type(None)):
894         ecdf = make_ecdf(data, var, wgt)
895     x = var
896     if wgt:
897         y = wgt
898     else:
899         if var == "y":
900             y = "y1"
901         else:
902             y = "y"
903     def check_a(a):
904         return estimate_CSS_InvG_fit_for_a(t, phi, psi, a, ecdf, x, y)["fit"]
905     sol = min(check_a, x0=a0).x[0]
906     return estimate_CSS_InvG_fit_for_a(t, phi, psi, sol, ecdf, x, y)

```

And add entries in the dictionaries.

```

907     distribution["CSS_InvG"] = cdf_CSS_InvG
908     density["CSS_InvG"]      = density_CSS_InvG
909     likelihood["CSS_InvG"]   = L_CSS_InvG
910     estimator["CSS_InvG"]    = estimate_CSS_InvG

```

CSS_InvG_prop Constant-shift-scale inverse-gamma distribution where c is proportional to α and time instead of a linear function of α and time. Parameters:

- $\alpha > 0$: shape parameter

This distribution is the same as in inverse-gamma distribution except that we fix β and c to be proportional to α . Specifically, $c_t = \alpha_t(\psi_0 + \psi_1 t) = \phi\beta_t$ (and equivalently $\beta = c/\phi$). From `check_constants.py`, we approximate values of the constants as follows:

- ϕ : -0.1337
- ψ_0 : \$206,824
- ψ_1 : -\$105.33/year

See also Table 7. For this set of estimation functions, we will have several more arguments than for the other models. The cumulative distribution function is

$$y = Q\left(\alpha, \frac{1}{\phi} \frac{\alpha(\psi_0 + \psi_1 t)}{x - \alpha(\psi_0 + \psi_1 t)}\right)$$

Code:

```

911 def cdf_CSS_InvG_prop(x, t, phi, psi, a):
912     psi0, psi1 = psi
913     b = a * (psi0 + psi1 * t) / phi
914     c = a * (psi0 + psi1 * t)
915     x = x.astype(np.float64)
916     mask = (x <= c)
917     x[mask] = 0
918     frac = exp(log(b) - log(x[~mask] - c))
919     x[~mask] = Q(a, frac)
920     return x

```

The density is

$$y = \frac{\alpha^\alpha (\psi_0 + \psi_1 t)^\alpha}{\phi^\alpha \Gamma(\alpha)} \frac{e^{-\frac{1}{\phi} \frac{\alpha(\psi_0 + \psi_1 t)}{x - \alpha(\psi_0 + \psi_1 t)}}}{(x - \alpha(\psi_0 + \psi_1 t))^{1+\alpha}}$$

```

921 def density_CSS_InvG_prop(x, t, phi, psi, a):
922     psi0, psi1 = psi
923     beta = a * (psi0 + psi1 * t) / phi
924     c = a * (psi0 + psi1 * t)
925     if x <= c:
926         return 0
927     else:
928         return density_InvG(x, [a, beta, c])

```

Likelihood:

$$L = n\alpha \log \left(\frac{\alpha(\psi_0 + \psi_1 t)}{\phi} \right) - n \log \Gamma(\alpha) - \left(\frac{\alpha(\psi_0 + \psi_1 t)}{\phi} \right) \sum_{i=1}^n \frac{w_i}{x - \alpha(\psi_0 + \psi_1 t)} \\ - (1 + \alpha) \sum_{i=1}^n w_i \log(x_i - \alpha(\psi_0 + \psi_1 t))$$

```

929 def L_CSS_InvG_prop(data, t, phi, psi, a, var, wgt=None):
930     validate_var_wgt(data, var, wgt)
931     psi0, psi1 = psi
932     beta = a * (psi0 + psi1 * t) / phi
933     c = a * (psi0 + psi1 * t)
934     return L_InvG(data, [a, beta, c], var, wgt)

```

Unfortunately, maximum likelihood does not work properly with the constant-shift-scale inverse-gamma distribution. So we minimize the Kolmogorov-Smirnov statistic instead. The `fit_for_a` function will calculate the Kolmogorov-Smirnov statistic.

```

935 def estimate_CSS_InvG_prop_fit_for_a(t, phi, psi, a, ecdf, x, y):
936     if a <= 0:
937         return {"fit": 1, "parameters": a}
938     else:
939         psi0, psi1 = psi
940         beta = a * (psi0 + psi1 * t) / phi
941         c = a * (psi0 + psi1 * t)
942         def F(x):
943             return cdf_InvG(x, [a, beta, c])
944         temp = {
945             "fit": kolmogorov_smirnov(F, ecdf=ecdf, x=x, y=y),
946             "parameters": [a]}
947         return temp

```

And the main estimation function calculates the minimum of the `fit_for_a` function. The `a0` argument is the starting value for α in the minimization. We will take `a0` to be the value of α from year `t` under the 3-parameter model.

```

948 def estimate_CSS_InvG_prop(data, t, phi, psi, a0, var, wgt=None, *, \
949                             ecdf=None, x=None, y=None):
950     validate_var_wgt(data, var, wgt)
951     if isinstance(ecdf, type(None)):
952         ecdf = make_ecdf(data, var, wgt)
953     x = var
954     if wgt:
955         y = wgt
956     else:
957         if var == "y":
958             y = "y1"
959         else:

```



```

960     y = "y"
961     def check_a(a):
962         return estimate_CSS_InvG_prop_fit_for_a(t, phi, psi, a, ecdf, x, y)["fit"]
963     sol = min(check_a, x0=a0).x[0]
964     return estimate_CSS_InvG_prop_fit_for_a(t, phi, psi, sol, ecdf, x, y)

```

And add entries in the dictionaries.

```

965     distribution["CSS_InvG_prop"] = cdf_CSS_InvG_prop
966     density["CSS_InvG_prop"]      = density_CSS_InvG_prop
967     likelihood["CSS_InvG_prop"]   = L_CSS_InvG_prop
968     estimator["CSS_InvG_prop"]    = estimate_CSS_InvG_prop

```

All inverse-gamma distributions done!

LogN_P_cut Pareto and log-normal distributions with cutoff. Parameters:

- μ : log mean
- $\sigma^2 > 0$: log variance
- $k > x_m + c$: cutoff (we need the cutoff to exceed the lower bound on the support of shifted Pareto)
- $x_m > 0$: Pareto scale parameter
- $\alpha > 0$ Pareto shape parameter
- c : shift parameter (for shifted distribution)

Cumulative distribution function:

$$y = \begin{cases} \Phi\left(\frac{\log(x - c) - \mu}{\sigma}\right) & \text{if } x < k \\ 1 - \left(\frac{x_m}{x - c}\right)^\alpha & \text{if } x \geq k \end{cases}$$

```

969 def cdf_LogN_P_cut(x, params):
970     mu, sigma_sq, k, x_m, a, c = params
971     x = x.astype(np.float64)
972     mask_c = x <= c
973     mask_k = x <= k
974     mask_xm = x < x_m
975     if k < c:
976         raise RuntimeError("k was less than c; problem?")

```

The LogN_P_cut distribution is a bit different from the other models in that its domain gets partitioned into two separate pieces. So we will end up transforming two separate subarrays, one for data less than k and one for data greater than k .

```

977     x[mask_c] = 0
978     x[(~mask_c) & mask_k] = Phi(
979         (log(x[(~mask_c) & mask_k] - c) - mu) / sqrt(sigma_sq))
980     x[(~mask_k) & mask_xm] = 0
981     x[(~mask_k) & (~mask_xm)] = \
982         1 - (x_m / (x[(~mask_k) & (~mask_xm)] - c)) ** a
983     return x

```

Density:

$$y = \begin{cases} \frac{1}{(x-c)\sigma\sqrt{2\pi}} e^{-(\log(x-c)-\mu)^2/2\sigma^2} & \text{if } x < k \\ \frac{\alpha x_m^\alpha}{(x-c)^{1+\alpha}} & \text{if } x \geq k \end{cases}$$

Code:

```

984 def density_LogN_P_cut(x, params):
985     mu, sigma_sq, k, x_m, a, c = params
986     if k < c:
987         print("k was less than c; setting k=c")
988         k = c
989     if x <= c:
990         return 0
991     elif x > c and x < k:
992         frac = 1/((x-c) * sqrt(2 * pi * sigma_sq))
993         exponent = -(log(x-c) - mu)**2 / (2 * sigma_sq)
994         return frac * exp(exponent)
995     elif x >= k:
996         if x <= x_m + c:
997             return 0
998         else:
999             return exp(log(a) + a * log(x_m) - (1+a) * log(x-c))
1000 else:
1001     raise RuntimeError("Something weird happened")

```

Likelihood:

$$L = \sum_{i=1}^n w_i \psi_i,$$

where

$$\psi_i = \begin{cases} -\log(x_i - c) - \log(\sigma) - \frac{1}{2}\log(2\pi) - \frac{(\log(x_i - c) - \mu)^2}{2\sigma^2} & \text{if } x_i < k \\ \log \alpha + \alpha \log(x_m) - (1 + \alpha) \log(x_i - c) & \text{if } x_i \geq k \end{cases}$$

Code:

```

1002 def L_LogN_P_cut(data, params, var, wgt=None):
1003     validate_var_wgt(data, var, wgt)
1004     mu, sigma_sq, k, x_m, a, c = params
1005     data_down = data[data[var] < k]
1006     data_up = data[data[var] >= k]
1007     if wgt:
1008         n_down = data_below_k[wgt].sum()
1009         n_up = data_above_k[wgt].sum()
1010         down_terms = -(data_down[wgt] * log(data_down[var] - c)).sum() - \

```

```

1011     (n_down/2) * log(2 * sigma_sq * pi) - \
1012     (data_down[wgt] * (log(data_down[var] - c) - mu)**2).sum() * \
1013     (1 / (2*sigma_sq))
1014     up_terms = n_up * (log(a) + a * log(x_m)) + \
1015     (1+a) * (data_up[wgt] * log(data_up[var] - c)).sum()
1016 else:
1017     n_down = len(data_below_k)
1018     n_up = len(data_above_k)
1019     down_terms = -log(data_down[var] - c).sum() - \
1020     (n_down/2) * log(2 * sigma_sq * pi) - \
1021     ((log(data_down.income() - c) - mu)**2).sum() / (2*sigma_sq)
1022     up_terms = n_up * (log(a) + a * log(x_m)) + \
1023     (1+a) * log(data_up.income - c).sum()
1024 return down_terms + up_terms

```

Unshifted version:

```

1025 def L_LogN_P_cut_unshift(data, params, var, wgt=None):
1026     mu, sigma_sq, k, x_m, a = params
1027     return L_LogN_P_cut(data, [mu, sigma_sq, k, x_m, a, 0], var, wgt)

```

We can solve for the MLE symbolically. (Or at least solve most of the way.) The algebra is a bit ugly though. The one restriction is that the density should integrate to 1. In other words,

$$\Phi\left(\frac{\log k - \mu}{\sigma}\right) + \left(\frac{x_m}{k}\right)^\alpha = 1,$$

where Φ is the cumulative distribution of a standard normal random variable. We can maximize the likelihood subject to this restriction. The restricted likelihood function/Lagrangian is

$$L = \sum_{i=1}^n w_i \psi_i + \lambda \left(1 - \Phi\left(\frac{\log k - \mu}{\sigma}\right) - \left(\frac{x_m}{k}\right)^\alpha \right)$$

Assume that k and c are constant, and let j be the index that divides the incomes between [less than k] and [greater than k]. If we let n_\downarrow and n^\uparrow be the sums of weights for incomes less than k and greater than k respectively, then differentiating gives us

$$\begin{aligned}
\frac{dL}{d\mu} &= \sum_{i=1}^j w_i \left(-\frac{1}{2\sigma^2} \right) (-1)(\log(x_i - c) - \mu) - \lambda \left(-\frac{1}{\sigma} \right) \Phi' \left(\frac{\log k - \mu}{\sigma} \right) \\
&= \frac{1}{\sigma^2} \sum_{i=1}^j w_i (\log(x_i - c) - \mu) + \lambda \left(\frac{1}{\sigma} \right) \Phi' \left(\frac{\log k - \mu}{\sigma} \right) = 0 \\
\frac{dL}{d\sigma} &= -n_\downarrow \left(\frac{1}{\sigma} \right) - \left(-\frac{1}{\sigma^3} \right) \sum_{i=1}^j w_i (\log(x_i - c) - \mu)^2 \\
&\quad - \lambda \left(-\frac{\log k - \mu}{\sigma^2} \right) \Phi' \left(\frac{\log k - \mu}{\sigma} \right)
\end{aligned}$$

$$\begin{aligned}
&= -\frac{n_{\downarrow}}{\sigma} + \frac{1}{\sigma^3} \sum_{i=1}^j w_i (\log(x_i - c) - \mu)^2 \\
&\quad + \lambda \left(\frac{\log k - \mu}{\sigma^2} \right) \Phi' \left(\frac{\log k - \mu}{\sigma} \right) = 0 \\
\frac{dL}{d\alpha} &= n^{\uparrow} \left[\frac{1}{\alpha} + \log(x_m) \right] - \sum_{i=j+1}^n w_i \log(x_i - c) - \lambda \log \left(\frac{x_m}{k} \right) \left(\frac{x_m}{k} \right)^{\alpha} = 0 \\
\frac{dL}{dx_m} &= n^{\uparrow} \frac{\alpha}{x_m} - \lambda \left(\frac{\alpha}{k} \right) \left(\frac{x_m}{k} \right)^{\alpha-1} = 0
\end{aligned}$$

From the x_m equation, we have

$$\begin{aligned}
n^{\uparrow} \frac{\alpha}{x_m} - \lambda \left(\frac{\alpha}{k} \right) \left(\frac{x_m}{k} \right)^{\alpha-1} &= 0 \\
n^{\uparrow} - \lambda \left(\frac{x_m}{\alpha} \right) \left(\frac{\alpha}{k} \right) \left(\frac{x_m}{k} \right)^{\alpha-1} &= 0 \\
n^{\uparrow} - \lambda \left(\frac{x_m}{k} \right)^{\alpha} &= 0
\end{aligned}$$

Combining with the derivative for α gives us

$$\begin{aligned}
n^{\uparrow} \left[\frac{1}{\alpha} + \log(x_m) \right] - \sum_{i=j+1}^n w_i \log(x_i - c) - n^{\uparrow} \log \left(\frac{x_m}{k} \right) &= 0 \\
\frac{1}{\alpha} + \log(x_m) - \frac{1}{n^{\uparrow}} \sum_{i=j+1}^n w_i \log(x_i - c) - \log(x_m) + \log k &= 0 \\
\frac{1}{\alpha} - \frac{1}{n^{\uparrow}} \sum_{i=j+1}^n w_i \log \left(\frac{x_i - c}{k} \right) &= 0 \\
\alpha &= \frac{1}{\frac{1}{n^{\uparrow}} \sum_{i=j+1}^n w_i \log \left(\frac{x_i - c}{k} \right)}
\end{aligned}$$

And that gives us α in terms of data and constants! (Again assuming that c and k are constant.) Consider the derivatives with respect to μ and σ . Rearranging and dividing both equations gives us

$$\begin{aligned}
\frac{-\frac{n_{\downarrow}}{\sigma} + \frac{1}{\sigma^3} \sum_{i=1}^j w_i (\log(x_i - c) - \mu)^2}{\frac{1}{\sigma^2} \sum_{i=1}^j w_i (\log(x_i - c) - \mu)} &= \frac{-\lambda \left(\frac{\log k - \mu}{\sigma^2} \right) \Phi' \left(\frac{\log k - \mu}{\sigma} \right)}{-\lambda \left(\frac{1}{\sigma} \right) \Phi' \left(\frac{\log k - \mu}{\sigma} \right)} \\
-\frac{n_{\downarrow}}{\sigma} + \frac{1}{\sigma^3} \sum_{i=1}^j w_i (\log(x_i - c) - \mu)^2 &= \left(\frac{\log k - \mu}{\sigma} \right) \frac{1}{\sigma^2} \sum_{i=1}^j w_i (\log(x_i - c) - \mu)
\end{aligned}$$

Altogether this gives us

$$\sigma^2 = \frac{1}{n_{\downarrow}} \sum_{i=1}^j w_i (\log(x_i - c) - \mu)^2 - \frac{\log k - \mu}{n_{\downarrow}} \sum_{i=1}^j w_i (\log(x_i - c) - \mu)$$

for σ^2 in terms of μ . Substituting back into (either) one of the derivatives gives us

$$\begin{aligned} -\frac{n_{\downarrow}}{\sigma} + \frac{1}{\sigma^3} \sum_{i=1}^j w_i (\log(x_i - c) - \mu)^2 + \lambda \left(\frac{\log k - \mu}{\sigma^2} \right) \Phi' \left(\frac{\log k - \mu}{\sigma} \right) &= 0 \\ -n_{\downarrow} \sigma^2 + \sum_{i=1}^j w_i (\log(x_i - c) - \mu)^2 + \lambda \sigma (\log k - \mu) \Phi' \left(\frac{\log k - \mu}{\sigma} \right) &= 0 \\ -n_{\downarrow} \left[\frac{1}{n_{\downarrow}} \sum_{i=1}^j w_i (\log(x_i - c) - \mu)^2 - \frac{\log k - \mu}{n_{\downarrow}} \sum_{i=1}^j w_i (\log(x_i - c) - \mu) \right] \\ + \sum_{i=1}^j w_i (\log(x_i - c) - \mu)^2 + \lambda \sigma (\log k - \mu) \Phi' \left(\frac{\log k - \mu}{\sigma} \right) &= 0 \\ (\log k - \mu) \sum_{i=1}^j w_i (\log(x_i - c) - \mu) + \lambda \sigma (\log k - \mu) \Phi' \left(\frac{\log k - \mu}{\sigma} \right) &= 0 \\ \sum_{i=1}^j w_i (\log(x_i - c) - \mu) + \lambda \sigma \Phi' \left(\frac{\log k - \mu}{\sigma} \right) &= 0 \\ \sum_{i=1}^j w_i \log(x_i - c) - n_{\downarrow} \mu + \lambda \sigma \Phi' \left(\frac{\log k - \mu}{\sigma} \right) &= 0 \end{aligned}$$

If we examine the formula for σ^2 , we see that we can simplify it.

$$\begin{aligned} \sigma^2 &= \frac{1}{n_{\downarrow}} \sum_{i=1}^j w_i (\log(x_i - c) - \mu)^2 - \frac{\log k - \mu}{n_{\downarrow}} \sum_{i=1}^j w_i (\log(x_i - c) - \mu) \\ &= \frac{1}{n_{\downarrow}} \sum_{i=1}^j w_i \log(x_i - c)^2 - 2\mu \frac{1}{n_{\downarrow}} \sum_{i=1}^j w_i \log(x_i - c) + \mu^2 \\ &\quad - \frac{\log k}{n_{\downarrow}} \sum_{i=1}^j w_i \log(x_i - c) + \frac{\mu}{n_{\downarrow}} \sum_{i=1}^j w_i \log(x_i - c) + (\log k - \mu) \mu \\ &= \frac{1}{n_{\downarrow}} \sum_{i=1}^j w_i \log(x_i - c)^2 - \frac{\log k}{n_{\downarrow}} \sum_{i=1}^j w_i \log(x_i - c) \\ &\quad + \mu \left[\log k - \frac{1}{n_{\downarrow}} \sum_{i=1}^j w_i \log(x_i - c) \right] \\ &= p_0 + p_1 \mu \end{aligned}$$

The result is that σ^2 is a linear function of μ . It follows that $\sigma = \sqrt{p_0 + p_1\mu}$. Looking back to our mass condition and the condition for x_m , we see that

$$\begin{aligned} n^\uparrow &= \lambda \left(\frac{x_m}{k} \right)^\alpha \\ \lambda &= \frac{n^\uparrow}{\left(\frac{x_m}{k} \right)^\alpha} \\ &= \frac{n^\uparrow}{1 - \Phi\left(\frac{\log k - \mu}{\sigma}\right)} \end{aligned}$$

It follows that μ must solve

$$\sum_{i=1}^j w_i \log(x_i - c) - n_\downarrow \mu + n^\uparrow \sqrt{p_0 + p_1 \mu} \frac{\Phi'\left(\frac{\log k - \mu}{\sqrt{p_0 + p_1 \mu}}\right)}{1 - \Phi\left(\frac{\log k - \mu}{\sqrt{p_0 + p_1 \mu}}\right)} = 0$$

where

$$\begin{aligned} p_0 &= \frac{1}{n_\downarrow} \sum_{i=1}^j w_i \log(x_i - c)^2 - \frac{\log k}{n_\downarrow} \sum_{i=1}^j w_i \log(x_i - c) \\ p_1 &= \log k - \frac{1}{n_\downarrow} \sum_{i=1}^j w_i \log(x_i - c) \end{aligned}$$

This is a single equation with a single unknown and nothing that is obviously a problem for a numerical solver. So we will solve the equation numerically for μ . Then we find σ^2 easily because $\sigma^2 = p_0 + p_1\mu$. We already have α , and we can find x_m by solving the mass condition. This gives

$$x_m = k \left[1 - \Phi\left(\frac{\log k - \mu}{\sigma}\right) \right]^{1/\alpha}$$

for x_m .

We code this solution in several steps:

1. Write a function that returns p_0 and p_1 in terms of k and (unshifted) data. For convenience, we will say

$$p_2 = \frac{1}{n_\downarrow} \sum_{i=1}^j w_i \log(x_i - c) = \log k - p_1$$

and return that value too. The `constants` will be the return values of this function, specifically p_0 , p_1 , p_2 , n_\downarrow , and n^\uparrow .

2. Write a function that finds μ given the `constants`.
3. Write a function that finds α in terms of k and (unshifted) data.
4. Once we know μ and α it is easy to find the remaining parameters.
5. Then check whether this c - k pair results in a good-fitting model to the data.

We are expecting to feed pre-shifted data to the function, so the function does not have a shift parameter in its argument. (And the k value will also have been shifted by c . Shifting things with a cutoff parameter can get complicated quickly.) We use `_down` and `_up` to distinguish between DataFrames with incomes less than or greater than k respectively, similarly to how we notated n_{\downarrow} and n^{\uparrow} .

```

1028 def estimate_LogN_P_cut_cons(data, k, var, wgt):
1029     data_down = data[data[var] < k]
1030     data_up = data[data[var] >= k]
1031     if wgt:
1032         n_down = data_down[wgt].sum()
1033         n_up = data_up[wgt].sum()
1034         p2 = (1/n_down) * (data_down[wgt] * log(data_down[var])).sum()
1035         p0 = (1/n_down) * (data_down[wgt] * log(data_down[var])**2).sum() - \
1036             log(k) * p2
1037         p1 = log(k) - p2
1038     else:
1039         n_down = len(data_down)
1040         n_up = len(data_up)
1041         p2 = (1/n_down) * log(data_down[var]).sum()
1042         p0 = (1/n_down) * (log(data_down[var])**2).sum() - log(k) * p2
1043         p1 = log(k) - p2
1044     return [p0, p1, p2, n_down, n_up, k]

```

Now we define a function that accepts data as its argument and returns the value of μ . The function will also accept k as its argument and a list of constants (which come from the `_const` estimation function). The `eq_to_solve` returns the value of the expression that μ should solve.

```

1045 def estimate_LogN_P_cut_get_mu(cons):
1046     p0, p1, p2, n_down, n_up, k = cons
1047     def eq_to_solve(mu):
1048         ratio = (log(k) - mu) / sqrt(p0 + p1 * mu)
1049         if Phi(ratio) == 1:
1050             hazard = np.inf
1051         else:
1052             hazard = Phi_prime(ratio) / (1 - Phi(ratio))
1053     return n_down * p2 - n_down * mu + n_up * sqrt(p0 + p1 * mu) * hazard

```

For the bounds on μ for the numerical solver, we set `min_val` to be the value of μ where $\sigma = 0$, i.e. $-p_0/p_1$. An application of L'Hospital's rule shows that

$$\lim_{x \rightarrow \infty} \frac{1}{\sqrt{x}} \frac{\Phi'(x)}{1 - \Phi(x)} = \infty,$$

so close to `min_val`, the expression will be positive. Then we set `max_val` equal to `min_val` and steadily increment it until we reach a candidate for μ where the expression is negative. This is guaranteed to happen eventually because as μ increases, the ratio $\log k - \mu/\sigma$ decreases and eventually becomes negative. It follows that the quotient $\Phi'/1 - \Phi$ is eventually bounded above by 1, and because $p_1 > 0$, the $-n_{\downarrow}\mu$ term dominates. Then we feed everything into the numerical solver to recover μ .

```

1054 min_val = -p0 / p1 + 1.0e-6
1055 max_val = min_val
1056 while eq_to_solve(max_val) > 0:
1057     max_val = max_val + 2
1058 sol = root(eq_to_solve, bracket=[min_val, max_val])
1059 return sol.root

```

And we need a function to find α .

```

1060 def estimate_LogN_P_cut_get_alpha(data, cons, var, wgt):
1061     p0, p1, p2, n_down, n_up, k = cons
1062     data_up = data[data[var] >= k]
1063     if wgt:
1064         temp = (data_up[wgt] * log(data_up[var] / k)).sum()
1065     else:
1066         temp = log(data_up[var] / k).sum()
1067     return n_up / temp

```

If we know μ , we can recover the other parameters.

```

1068 def estimate_LogN_P_cut_get_params(mu, alpha, cons):
1069     p0, p1, p2, n_down, n_up, k = cons
1070     sigma_sq = p0 + p1 * mu
1071     x_m = k * (1 - Phi((log(k) - mu) / sqrt(sigma_sq))) ** (1/alpha)
1072     return [sigma_sq, x_m]

```

And we are finally ready to proceed with the estimation. Same as before, we define the `_unshift` estimation function that does the heavy lifting, and then the `fit_for_c_k` and main estimator functions will find c and k to minimize the chi-squared statistic. In this function, we have already shifted k . (Because everything is unshifted here.)

```

1073 def estimate_LogN_P_cut_unshift(data, k, var, wgt):
1074     constants = estimate_LogN_P_cut_cons(data, k, var, wgt)
1075     mu = estimate_LogN_P_cut_get_mu(constants)
1076     alpha = estimate_LogN_P_cut_get_alpha(data, constants, var, wgt)
1077     sigma_sq, x_m = estimate_LogN_P_cut_get_params(mu, alpha, constants)
1078     return [mu, sigma_sq, x_m, alpha]

```

And now the `_fit_for_c_k` function.

```

1079 def estimate_LogN_P_cut_fit_for_c(data, c, k, var, wgt, ecdf, x, y):
1080     shift_data = data.copy()
1081     shift_data[var] = shift_data[var] - c

```



```

1082 shift_data = shift_data[shift_data[var] > 0]
1083 shift_k = k - c
1084 mu, sigma_sq, x_m, alpha = estimate_LogN_P_cut_unshift(shift_data,
1085     shift_k, var, wgt)
1086 def F(x):
1087     return cdf_LogN_P_cut(x, [mu, sigma_sq, k, x_m, alpha, c])
1088 temp = {
1089     "fit": kolmogorov_smirnov(F, ecdf=ecdf, x=x, y=y),
1090     "parameters": [mu, sigma_sq, k, x_m, alpha, c]}
1091 return temp

```

The main estimator function will minimize the Kolmogorov-Smirnov statistic like in the previous other cases. Now we have an initial guess rather than bounds because the minimization is two-dimensional. We will use Nelder-Mead algorithm. Examination of the data suggests that the minimum chi-squared value occurs somewhere around $c = -\$9000$ and $k = \$110,000$, at least for 2019.

```

1092 def estimate_LogN_P_cut(data, var, wgt=None, *, ecdf=None, x=None, y=None):
1093     validate_var_wgt(data, var, wgt)
1094     if isinstance(ecdf, type(None)):
1095         ecdf = make_ecdf(data, var, wgt)
1096         x = var
1097         if wgt:
1098             y = wgt
1099         else:
1100             if var == "y":
1101                 y = "y1"
1102             else:
1103                 y = "y"
1104     def check_c_k(p):
1105         c, k = p
1106         f = estimate_LogN_P_cut_fit_for_c(data, c, k, var, wgt, ecdf, x, y)["fit"]
1107         return f
1108     sol = min(check_c_k, method="Nelder-Mead", x0=[-9000, 110000],
1109         options={"xatol": 0.1})
1110     c, k = sol.x
1111     return estimate_LogN_P_cut_fit_for_c(data, c, k, var, wgt, ecdf, x, y)

```

And add the functions to the appropriate dictionaries

```

1112 distribution["LogN_P_cut"] = cdf_LogN_P_cut
1113 density["LogN_P_cut"]     = density_LogN_P_cut
1114 likelihood["LogN_P_cut"]  = L_LogN_P_cut
1115 estimator["LogN_P_cut"]   = estimate_LogN_P_cut

```

LogN_P_mix Log-normal and Pareto mixture. Parameters:

- μ : log mean
- $\sigma^2 > 0$: log variance

- $x_m > 0$: Pareto scale parameter
- $\alpha > 0$ Pareto shape parameter
- γ : mixture parameter
- c : shift parameter (for shifted distribution)

Cumulative distribution function:

$$y = \gamma \Phi\left(\frac{\log(x - c) - \mu}{\sigma}\right) + (1 - \gamma) \chi_{x \geq x_m + c} \left[1 - \left(\frac{x_m}{x - c}\right)^\alpha\right]$$

```

1116 def cdf_LogN_P_mix(x, params):
1117     mu, sigma_sq, gamma, x_m, alpha, c = params
1118     if gamma < 0 or gamma > 1:
1119         raise ValueError("gamma is outside unit interval")
1120     mask_c = x <= c
1121     mask_c_xm = x < x_m + c
1122     term1 = x.astype(np.float64)
1123     term1[mask_c] = 0
1124     term1[~mask_c] = Phi((log(term1[~mask_c] - c) - mu) / sqrt(sigma_sq))
1125     term2 = x.astype(np.float64)
1126     term2[mask_c_xm] = 0
1127     term2[~mask_c_xm] = \
1128         1 - exp(alpha * (log(x_m) - log(term2[~mask_c_xm] - c)))
1129     return gamma * term1 + (1 - gamma) * term2

```

where $\chi_{x \geq x_m + c}$ is an indicator function. Density is

$$y = \frac{\gamma}{(x - c)\sigma\sqrt{2\pi}} e^{-(\log(x - c) - \mu)^2 / 2\sigma^2} + (1 - \gamma) \chi_{x \geq x_m + c} \frac{\alpha x_m^\alpha}{(x - c)^{1 + \alpha}},$$

where $\chi_{x \geq x_m + c}$ is an indicator function. Code:

```

1130 def density_LogN_P_mix(x, params):
1131     mu, sigma_sq, gamma, x_m, alpha, c = params
1132     if gamma < 0 or gamma > 1:
1133         raise ValueError("gamma is outside unit interval")
1134     if x <= c:
1135         return 0
1136     else:
1137         log_n_term = exp(log(gamma) - log(x - c) -
1138             0.5 * log(2 * pi * sigma_sq) -
1139             (log(x - c) - mu) ** 2 / (2 * sigma_sq))
1140         if x <= x_m + c:
1141             pareto_term = 0
1142         else:
1143             pareto_term = exp(log(1 - gamma) + log(alpha) + alpha * log(x_m) -
1144                 (1 + alpha) * log(x - c))
1145         return log_n_term + pareto_term

```

Likelihood:

$$L = \sum_{i=1}^n w_i \log \left[\frac{\gamma}{(x_i - c) \sigma \sqrt{2\pi}} e^{-(\log(x_i - c) - \mu)^2 / 2\sigma^2} + (1 - \gamma) \chi_{x_i \geq x_m + c} \frac{\alpha x_m^\alpha}{(x_i - c)^{1+\alpha}} \right]$$

Again, we use the `data_up/data_down` naming convention for the portion of the data that is above versus below the Pareto cutoff.

```

1146 def L_LogN_P_mix(data, params, var, wgt=None):
1147     validate_var_wgt(data, var, wgt)
1148     mu, sigma_sq, gamma, x_m, alpha, c = params
1149     if gamma < 0 or gamma > 1:
1150         raise ValueError("gamma is outside unit interval")
1151     data_down = data[data[var] <= x_m + c]
1152     data_up = data[data[var] > x_m + c]
1153     if wgt:
1154         down_terms = (data_down[wgt] *
1155             (log(gamma) - log(data_down[var] - c) - 0.5 * log(2 * pi * sigma_sq) -
1156             (log(data_down[var] - c) - mu) ** 2 / (2 * sigma_sq))).sum()
1157         up_terms = (data_up[wgt] * log(
1158             exp(log(gamma) - log(data_up[var] - c) -
1159             0.5 * log(2 * pi * sigma_sq) -
1160             (log(data_up[var] - c) - mu) ** 2 / (2 * sigma_sq)) +
1161             exp(log(1 - gamma) + log(alpha) + alpha * log(x_m) -
1162             (1 + alpha) * log(data_up[var] - c)))).sum()
1163     else:
1164         down_terms = (log(gamma) -
1165             log(data_down[var] - c) - 0.5 * log(2 * pi * sigma_sq) -
1166             (log(data_down[var] - c) - mu) ** 2 / (2 * sigma_sq)).sum()
1167         up_terms = log(
1168             exp(log(gamma) - log(data_up[var] - c) -
1169             0.5 * log(2 * pi * sigma_sq) -
1170             (log(data_up[var] - c) - mu) ** 2 / (2 * sigma_sq)) +
1171             exp(log(1 - gamma) + log(alpha) + alpha * log(x_m) -
1172             (1 + alpha) * log(data_up[var] - c)))).sum()
1173     return down_terms + up_terms

```

Unshifted version:

```

1174 def L_LogN_P_mix_unshift(data, params, wgt, var=None):
1175     mu, sigma_sq, gamma, x_m, alpha = params
1176     return L_LogN_P_mix(data, [mu, sigma_sq, gamma, x_m, alpha, 0], var, wgt)

```

Estimating the mixture distribution is tricky. Direct numerical maximization of the likelihood function tends to converge but not to the same results every time. Accordingly, we will estimate in three steps:

1. The first (innermost) step is to maximize the likelihood with respect to μ , σ^2 , and α conditional on the remaining parameters. We find μ and σ^2 using a truncated maximum-likelihood estimator on the portion of the data below $x_m + c$. Then we find α using a one-dimensional maximization.

2. The second step is choosing γ according to a first-order condition from the maximum-likelihood estimation. Potential values of γ must be large enough that the log-normal component of the mixture places enough mass on the portion of its support below $x_m + c$.
3. The third (outermost) step is a Kolmogorov-Smirnov statistic minimization with regard to c and x_m . Exploratory analysis suggests this minimum occurs for values of c and x_m relatively similar to their estimates under the cutoff model.

The three-step estimation appears to fit the model correctly.

For the log-normal parameters, we use a truncated maximum-likelihood estimator. However, simply maximizing the weighted log-sum of

$$\begin{aligned}
\tilde{f}(x) &= \frac{\frac{\gamma}{(x-c)\sigma\sqrt{2\pi}}e^{-(\log(x-c)-\mu)^2/2\sigma^2} + (1-\gamma)\chi_{x \geq x_m+c} \frac{\alpha x_m^\alpha}{(x-c)^{1+\alpha}}}{\gamma\Phi\left(\frac{\log(x-c)-\mu}{\sigma}\right) + (1-\gamma)\chi_{x \geq x_m+c} \left[1 - \left(\frac{x_m}{x-c}\right)^\alpha\right] \Big|_{x=x_m+c}} \\
&= \frac{\gamma}{(x-c)\sigma\sqrt{2\pi}}e^{-(\log(x-c)-\mu)^2/2\sigma^2} \Big/ \gamma\Phi\left(\frac{\log(x_m)-\mu}{\sigma}\right) \\
&= \frac{1}{(x-c)\sigma\sqrt{2\pi}}e^{-(\log(x-c)-\mu)^2/2\sigma^2} \Big/ \Phi\left(\frac{\log(x_m)-\mu}{\sigma}\right)
\end{aligned}$$

over data points less than $x_m + c$ doesn't work because it reduces to the standard maximum-likelihood estimator for the truncated log-normal distribution! In this case, the resulting probability mass allocated by the model to $[c, x_m + c]$ will be too small by a factor of λ . Instead, we maximize the likelihood subject to a mass constraint on $[c, x_m + c]$ given by the empirical mass fraction in this region. Let η be the fraction of data points less than $x_m + c$. Then the mass condition becomes

$$\Phi\left(\frac{\log(x_m)-\mu}{\sigma}\right) = \frac{\eta}{\gamma}$$

For algebraic convenience, we take logs of both sides. Then the Lagrangian is

$$\begin{aligned}
L &= -n \log \Phi\left(\frac{\log(x_m)-\mu}{\sigma}\right) - n \log \sigma - \frac{n}{2} \log(2\pi) - \sum_{i=1}^n w_i \log x_i \\
&\quad - \sum_{i=1}^n \frac{(\log x_i - \mu)^2}{2\sigma^2} + \lambda \left[\log\left(\frac{\eta}{\gamma}\right) - \log \Phi\left(\frac{\log(x_m)-\mu}{\sigma}\right) \right]
\end{aligned}$$

where λ is the Lagrange multiplier. (Throughout we assume for notational purposes that total weight has been normalized to n . But the exact total weight doesn't matter.) Notice also that we dropped the cs . To keep the notation cleaner, we perform the calculation for data that begins at 0. (So to apply to a DataFrame of income data, we

should first subtract c from all incomes.) Taking derivatives gives us

$$\begin{aligned}\frac{dL}{d\mu} &= -n \left(-\frac{1}{\sigma} \right) \frac{\Phi' \left(\frac{\log(x_m) - \mu}{\sigma} \right)}{\Phi \left(\frac{\log(x_m) - \mu}{\sigma} \right)} - \frac{1}{\sigma^2} \sum_{i=1}^n (\mu - \log x_i) - \lambda \left(-\frac{1}{\sigma} \right) \frac{\Phi' \left(\frac{\log(x_m) - \mu}{\sigma} \right)}{\Phi \left(\frac{\log(x_m) - \mu}{\sigma} \right)} \\ &= -\frac{1}{\sigma^2} \left[n\mu - \sum_{i=1}^n w_i \log x_i \right] + (n + \lambda) \left(\frac{1}{\sigma} \right) \frac{\Phi' \left(\frac{\log(x_m) - \mu}{\sigma} \right)}{\Phi \left(\frac{\log(x_m) - \mu}{\sigma} \right)} = 0\end{aligned}$$

and

$$\begin{aligned}\frac{dL}{d\sigma} &= -n \left(-\frac{\log(x_m) - \mu}{\sigma^2} \right) \frac{\Phi' \left(\frac{\log(x_m) - \mu}{\sigma} \right)}{\Phi \left(\frac{\log(x_m) - \mu}{\sigma} \right)} - \frac{n}{\sigma} - \left(-\frac{1}{\sigma^3} \right) \sum_{i=1}^n w_i (\log x_i - \mu)^2 \\ &\quad - \lambda \left(-\frac{\log(x_m) - \mu}{\sigma^2} \right) \frac{\Phi' \left(\frac{\log(x_m) - \mu}{\sigma} \right)}{\Phi \left(\frac{\log(x_m) - \mu}{\sigma} \right)} \\ &= -\frac{n}{\sigma} + \frac{1}{\sigma^3} \sum_{i=1}^n (\log x_i - \mu^2) + (\lambda + n) \left(\frac{\log(x_m) - \mu}{\sigma^2} \right) \frac{\Phi' \left(\frac{\log(x_m) - \mu}{\sigma} \right)}{\Phi \left(\frac{\log(x_m) - \mu}{\sigma} \right)} = 0\end{aligned}$$

We bring the (rightmost) terms containing λ to the other side and divide both equations. This gives us

$$\begin{aligned}\frac{\log(x_m) - \mu}{\sigma} &= \frac{\frac{n}{\sigma} - \frac{1}{\sigma^3} \sum_{i=1}^n (\log x_i - \mu)^2}{\frac{n\mu}{\sigma^2} - \frac{1}{\sigma^2} \sum_{i=1}^n \log x_i} \\ \Phi^{-1} \left(\frac{\eta}{\gamma} \right) \left[\frac{n\mu}{\sigma^2} - \frac{1}{\sigma^2} \sum_{i=1}^n \log x_i \right] &= \frac{n}{\sigma} - \frac{1}{\sigma^3} \sum_{i=1}^n (\log x_i - \mu)^2 \\ \Phi^{-1} \left(\frac{\eta}{\gamma} \right) \left[\mu\sigma - \frac{\sigma}{n} \sum_{i=1}^n \log x_i \right] &= \sigma^2 - \frac{1}{n} \sum_{i=1}^n (\log x_i - \mu)^2 \\ \Phi^{-1} \left(\frac{\eta}{\gamma} \right) (\mu\sigma - \sigma \overline{\log x_i}) &= \sigma^2 - (\overline{\log x_i^2} - 2\mu \overline{\log x_i} + \mu^2)\end{aligned}$$

$$\Phi^{-1}\left(\frac{\eta}{\gamma}\right)\sigma(\mu - \overline{\log x_i}) = \sigma^2 - \overline{\log x_i^2} + 2\mu \overline{\log x_i} - \mu^2$$

From the mass condition, we know that

$$\begin{aligned}\Phi^{-1}\left(\frac{\eta}{\gamma}\right) &= \frac{\log(x_m) - \mu}{\sigma} \\ \sigma &= \frac{1}{\Phi^{-1}\left(\frac{\eta}{\gamma}\right)}(\log(x_m) - \mu)\end{aligned}$$

Substituting back gives us

$$\begin{aligned}(\log(x_m) - \mu)(\mu - \overline{\log x_i}) &= \left[\frac{1}{\Phi^{-1}\left(\frac{\eta}{\gamma}\right)}(\log(x_m) - \mu) \right]^2 - \overline{\log x_i^2} + 2\mu \overline{\log x_i} - \mu^2 \\ \Phi^{-1}\left(\frac{\eta}{\gamma}\right)^2 \left[-\mu^2 + \mu(\log(x_m) + \overline{\log x_i}) - \log(x_m) \overline{\log x_i} \right] \\ &= (\log(x_m) - \mu)^2 + \Phi^{-1}\left(\frac{\eta}{\gamma}\right)^2 \left[-\overline{\log x_i^2} + 2\mu \overline{\log x_i} - \mu^2 \right] \\ \Phi^{-1}\left(\frac{\eta}{\gamma}\right)^2 \left[\mu(\log(x_m) + \overline{\log x_i}) - \log(x_m) \overline{\log x_i} \right] \\ &= \mu^2 - 2\mu \log(x_m) + \log(x_m)^2 + \Phi^{-1}\left(\frac{\eta}{\gamma}\right)^2 \left[-\overline{\log x_i^2} + 2\mu \overline{\log x_i} \right] \\ \mu^2 + \mu \left[\Phi^{-1}\left(\frac{\eta}{\gamma}\right)^2 (\overline{\log x_i} - \log(x_m)) - 2\log(x_m) \right] \\ &\quad + \left[\Phi^{-1}\left(\frac{\eta}{\gamma}\right)^2 (\log(x_m) \overline{\log x_i} - \overline{\log x_i^2}) + \log(x_m)^2 \right] = 0\end{aligned}$$

So μ solves a quadratic equation $x^2 + bx + c = 0$ with

$$\begin{aligned}b &= \Phi^{-1}\left(\frac{\eta}{\gamma}\right)^2 \left[\overline{\log x_i} - \log(x_m) \right] - 2\log(x_m) \\ c &= \Phi^{-1}\left(\frac{\eta}{\gamma}\right)^2 \left[\log(x_m) \overline{\log x_i} - \overline{\log x_i^2} \right] + \log(x_m)^2\end{aligned}$$

```

1177 def estimate_LogN_P_mix_logn_trunc(data, gamma, k, var, wgt):
1178     data_down = data[data[var] <= k]
1179     if wgt:
1180         n_down = data_down[wgt].sum()
1181         eta = n_down / data[wgt].sum()
1182         sum_log = (data_down[wgt] * log(data_down[var])).sum() / n_down

```

```

1183     sum_log_sq = (data_down[wgt] * log(data_down[var]) ** 2).sum() / n_down
1184 else:
1185     n_down = len(data_down)
1186     eta = n_down / len(data)
1187     sum_log = log(data_down[var]).sum() / n_down
1188     sum_log_sq = (log(data_down[var]) ** 2).sum() / n_down
1189 phi_coef = Phinv(eta / gamma) # eta is empirical mass
1190 phi_coef_sq = phi_coef ** 2
1191 b = phi_coef_sq * (sum_log - log(k)) - 2 * log(k)
1192 c = phi_coef_sq * (log(k) * sum_log - sum_log_sq) + log(k) ** 2

```

Given the estimate for μ , we recover σ from the mass condition. The two solutions μ_1 and μ_2 should yield σ_1 and σ_2 with opposite signs, so we take the greater σ_i as our estimate for σ (because we require $\sigma > 0$). Because σ is monotonic decreasing in μ , this also means we want the smaller μ_i . We have

$$\mu = \frac{-b \pm \sqrt{b^2 - 4c}}{2},$$

so we use the value of μ with the negative square root.

```

1193 mu = (-b - sqrt(b ** 2 - 4 * c)) / 2
1194 sigma = (log(k) - mu) / phi_coef
1195 return [mu, sigma ** 2]

```

The `_inner_params` function finds the values of μ , σ^2 , and α conditional on remaining parameters. Here the `data` should be unmodified (start from $-c$), and we shift it ourselves. We find α by directly maximizing the likelihood function.

```

1196 def estimate_LogN_P_mix_inner_params(data, gamma, x_m, c, var, wgt):
1197     data_above_c = data[data[var] > c]
1198     shift_data = data_above_c.copy()
1199     shift_data[var] = shift_data[var] - c
1200     mu, sigma_sq = estimate_LogN_P_mix_logn_trunc(shift_data,
1201         gamma, x_m, var, wgt)
1202     def neg_L(a):
1203         return -L_LogN_P_mix(data_above_c,
1204             [mu, sigma_sq, gamma, x_m, a, c], var, wgt)
1205     sol = min(neg_L, x0=1, bounds=[(0.5, 5)], method="Nelder-Mead")
1206     return [mu, sigma_sq, sol.x[0]]

```

Now we need to find γ . Let j be the index dividing points less than $x_m + c$ from points greater. Taking the derivative of the likelihood function with respect to γ gives us

$$\begin{aligned} \frac{dL}{d\gamma} &= \sum_{i=1}^n w_i \frac{\frac{1}{(x_i - c)\sigma\sqrt{2\pi}} e^{-(\log(x_i - c) - \mu)^2 / 2\sigma^2} - \chi_{x_i \geq x_m + c} \frac{\alpha x_m^\alpha}{(x_i - c)^{1+\alpha}}}{f(x_i)} \\ &= \sum_{i=1}^j \frac{w_i}{\gamma} + \sum_{i=j+1}^n w_i \frac{\frac{1}{(x_i - c)\sigma\sqrt{2\pi}} e^{-(\log(x_i - c) - \mu)^2 / 2\sigma^2}}{f(x_i)} \end{aligned}$$

$$\begin{aligned}
& - \sum_{i=j+1}^n \frac{\chi_{x_i \geq x_m + c} \frac{\alpha x_m^\alpha}{(x_i - c)^{1+\alpha}}}{f(x_i)} \\
& = \frac{n_{\downarrow}}{\gamma} + \sum_{i=j+1}^n \frac{w_i}{\gamma + (1 - \gamma)l_i} - \sum_{i=j+1}^n \frac{w_i l_i}{\gamma + (1 - \gamma)l_i} \\
& = \frac{n_{\downarrow}}{\gamma} + \sum_{i=j+1}^n \frac{w_i(1 - l_i)}{\gamma + l_i},
\end{aligned}$$

where f is the mixture density and l_i is the likelihood ratio

$$l_i = \frac{\alpha x_m^\alpha}{(x_i - c)^{1+\alpha}} \sigma \sqrt{2\pi} e^{-(\log(x_i - c) - \mu)^2 / 2\sigma^2},$$

i.e. the ratio of Pareto to log-normal densities. We want to choose γ such that the derivative is 0. Explatory analysis suggests that the derivative is always negative for γ sufficiently close to 1, and we also require that $\eta < \gamma$ for the amount of probability mass below $x_m + c$ to be well-defined. Accordingly, we first search for a maximum on $(\eta, 1)$, and then we look for a root between the maximum and 1. In the `find_gamma` function, we will assume unmodified data (i.e. incomes start at c).

```

1207 def estimate_LogN_P_mix_find_gamma(data, c, x_m, var, wgt):
1208     data_up = data[data[var] > x_m + c]
1209     if wgt:
1210         n_down = data[data[var] <= x_m + c][wgt].sum()
1211         n_up = data_up[wgt].sum()
1212         eta = n_down / data[wgt].sum()
1213     else:
1214         n_down = len(data[data[var] <= x_m + c])
1215         n_up = len(data_up)
1216         eta = n_down / len(data)

```

The function `dL` returns the derivative of the likelihood function. On $[\eta, 1]$, the maximum occurs at `max_g` and has value `max_dL`.

```

1217 def dL(g):
1218     if g <= eta or g >= 1:
1219         print("bad g for dL")
1220         return float("nan")
1221     mu, sigma_sq, alpha = estimate_LogN_P_mix_inner_params(data,
1222     g, x_m, c, var, wgt)
1223     l = exp(log(alpha) + alpha * log(x_m) -
1224     alpha * log(data_up[var] - c) + 0.5 * log(2 * pi * sigma_sq) +
1225     (log(data_up[var] - c) - mu) ** 2 / (2 * sigma_sq))
1226     term1 = n_down / g
1227     term2 = (data_up.weight * (1 - l) / (g + (1 - g) * l)).sum()
1228     return {"val": term1 + term2, "parameters": [mu, sigma_sq, g, alpha]}
1229 max_g = min_s(lambda x: -dL(x)["val"], method="bounded",
1230     bounds=[0.99 * eta + 0.01, 0.01 * eta + 0.99]).x
1231 max_dL = dL(max_g)["val"]

```


Now we want to find the root of dL , which will be our estimate of γ conditional on x_m and c . We know that dL is maximized at max_g , and the function should be negative near 1. If max_dL is negative, then the function is negative everywhere, and we have no estimate for γ , which means the original c and x_m values were bad. If max_dL is positive, then dL has a root somewhere between max_g and 1, which we find with a bounded root finder.

```

1232 if max_dL == 0:
1233     return dL(max_g)["parameters"]
1234 elif max_dL > 0:
1235     delta = 0.01
1236     while dL((delta) * max_g + (1 - delta))["val"] > 0:
1237         delta = delta / 10
1238     bound = (delta) * max_g + (1 - delta)
1239     g = root(lambda x: dL(x)["val"], bracket=[max_g, bound], xtol=1e-5).root
1240     return dL(g)["parameters"]
1241 elif max_dL < 0:
1242     return False

```

Now we have the `fit_for_c` function. This function calls the `find_gamma` function to calculate optimal parameter values conditional on x_m and c . Then it returns the value of the objective function and parameters. If the return value of `find_gamma` is a dictionary, we carry out the calculation normally. Otherwise, we have bad values, so we set the fit value to 1.

```

1243 def estimate_LogN_P_mix_fit_for_c(data, c, x_m, var, wgt, ecdf, x, y):
1244     shift_data = data.copy()
1245     shift_data[var] = shift_data[var] - c
1246     shift_data = shift_data[shift_data[var] > 0]
1247     temp_params = estimate_LogN_P_mix_find_gamma(data, c, x_m, var, wgt)
1248     if temp_params:
1249         mu, sigma_sq, gamma, alpha = temp_params
1250         def F(x):
1251             return cdf_LogN_P_mix(x, [mu, sigma_sq, gamma, x_m, alpha, c])
1252         temp = {
1253             "fit": kolmogorov_smirnov(F, ecdf=ecdf, x=x, y=y),
1254             "parameters": [mu, sigma_sq, gamma, x_m, alpha, c]}
1255         return temp
1256     else:
1257         print("bad gamma, xm = {0}, c = {1}".format(x_m, c))
1258         return {"fit": 1, "parameters": []}

```

Unfortunately, numerically minimizing the objective function appears to work poorly for the mixture model. We appear to be running into some practical identifiability issues here. So we do a brute-force minimization similar to how we handled GB2.

```

1259 def estimate_LogN_P_mix_get_cxm(dict):
1260     return [dict["parameters"][-1], dict["parameters"][3]]

```

Now we code the main estimation function. The setup is the same as for GB2 with a brute-force search using three iterations with the same mesh sizes. In the first iteration, we use a partition mesh size of \$500, and in the second iteration, we use a partition mesh size of \$100 around the ten best-fitting pairs from the first iteration. In the third iteration, we use a partition mesh size of \$10 around the ten best-fitting pairs from the second iteration. Unlike with GB2, we don't have to check the value of the α parameter.

```

1261 def estimate_LogN_P_mix(data, var, wgt=None, *, ecdf=None, x=None, y=None):
1262     validate_var_wgt(data, var, wgt)
1263     if isinstance(ecdf, type(None)):
1264         ecdf = make_ecdf(data, var, wgt)
1265         x = var
1266         if wgt:
1267             y = wgt
1268         else:
1269             if var == "y":
1270                 y = "y1"
1271             else:
1272                 y = "y"

```

First iteration. The $\langle x_m \text{ or } c \rangle_vals$ lists contain the test values of the variables we are using, and the `best_fits` list stores corresponding values of the objective function.

```

1273 x_m_vals = [35000 + 500*i for i in range(41)] # xm in [35000, 55000]
1274 c_vals = [-12000 + 500*i for i in range(21)] # c in [-12000,-2000]
1275 best_fits = []
1276 the_time()
1277 print("First iteration of brute-force search")
1278 for x_m in x_m_vals:
1279     for c in c_vals:
1280         temp = estimate_LogN_P_mix_fit_for_c(data, c, x_m, var, wgt, ecdf, x, y)

```

Again, same as with GB2, we modify `best_fits` dynamically. If the list has fewer than ten entries, we add the current results. Otherwise, we check whether the current result fits better than the last entry in `best_fits`. If no, we ignore. If yes, we insert it into `best_fits` at the first position where the current result fits worse than previous items in `best_fits`.

```

1281     if len(best_fits) < 10:
1282         best_fits.append(temp)
1283     else:
1284         if temp["fit"] < best_fits[-1]["fit"]:
1285             i = 0
1286             while temp["fit"] > best_fits[i]["fit"]:
1287                 i = i + 1
1288             best_fits.insert(i, temp)
1289             best_fits.pop()

```

For the second iteration, we store the (c, x_m) values in `pairs`. First extract the values

from `best_fits`, and then we populate the grid of parameter values to test. We turn `pairs` into a set to avoid duplicate parameter pairs.

```

1290 pairs = [estimate_LogN_P_mix_get_cxm(i) for i in best_fits]
1291 pairs = [(p[0] - 500 + 100*i, p[1]) for p in pairs for i in range(11)]
1292 pairs = [(p[0], p[1] - 500 + 100*i) for p in pairs for i in range(11)]
1293 pairs = set(pairs)
1294 best_fits = []
1295 the_time()
1296 print("Second iteration of brute-force search")
1297 for p in pairs:
1298     temp = estimate_LogN_P_mix_fit_for_c(data, *p, var, wgt, ecdf, x, y)
1299     if len(best_fits) < 10:
1300         best_fits.append(temp)
1301     else:
1302         if temp["fit"] < best_fits[-1]["fit"]:
1303             i = 0
1304             while temp["fit"] > best_fits[i]["fit"]:
1305                 i = i + 1
1306             best_fits.insert(i, temp)
1307             best_fits.pop()

```

Now do this again for the third iteration. This time, we save just the best-fitting entry, not the whole list.

```

1308 pairs = [estimate_LogN_P_mix_get_cxm(i) for i in best_fits]
1309 pairs = [(p[0] - 100 + 10*i, p[1]) for p in pairs for i in range(21)]
1310 pairs = [(p[0], p[1] - 100 + 10*i) for p in pairs for i in range(21)]
1311 pairs = set(pairs)
1312 solution = {"fit":2}
1313 the_time()
1314 print("Third iteration of brute-force search")
1315 for p in pairs:
1316     temp = estimate_LogN_P_mix_fit_for_c(data, *p, var, wgt, ecdf, x, y)
1317     if temp["fit"] < solution["fit"]:
1318         solution = temp
1319 return solution

```

And add functions to dictionaries.

```

1320 distribution["LogN_P_mix"] = cdf_LogN_P_mix
1321 density["LogN_P_mix"]      = density_LogN_P_mix
1322 likelihood["LogN_P_mix"]   = L_LogN_P_mix
1323 estimator["LogN_P_mix"]    = estimate_LogN_P_mix

```

File V

check_constants.py

This file contains code to investigate possible values for the constants in the inverse-gamma variants. We know that β and c are proportional, and we observe a time-dependent relationship between these parameters and α . Accordingly, we will model our coefficients as

$$c_t = \phi\beta_t$$
$$c_t = \psi_0 + \psi_1 t + \psi_2 \alpha_t \quad \text{or} \quad c_t = \alpha_t(\psi_0 + \psi_1 t),$$

and we want to estimate ϕ , ψ_0 , ψ_1 , and possibly ψ_2 for this relationship. It doesn't make sense to estimate ϕ in a way that depends on α since ϕ has nothing to do with α . So we'll estimate these constants in two steps: (1) find ϕ and then (2) minimize an objective function for ψ . If β and c are proportional, then we only need to determine the ratio of their magnitudes. It makes most sense to use the L^1 norm for magnitude in this context, so

$$\phi = \frac{-||c||_1}{||\beta||_1} = \sum_t c_t / \sum_t \beta_t,$$

Then we find ψ . We want to balance making $\alpha_t(\psi_0 + \psi_1 t)$ close to both $\phi\beta_t$ and c_t , so our objective function will be

$$L = \sum_t (\phi\beta_t - f(\alpha_t, t))^2 + \sum_t (c_t - f(\alpha_t, t))^2,$$

where ϕ is the conversion factor from above and f is the relationship between α and the other two parameters. It will be helpful to write our objective function in terms of vectors and matrices because the analysis is simpler with linear algebra. Let

$$A = \begin{bmatrix} 1 & \alpha_{t_0} & t_0 \\ 1 & \alpha_{t_1} & t_1 \\ \vdots & \vdots & \vdots \\ 1 & \alpha_{t_{n-1}} & t_{n-1} \end{bmatrix} \quad \text{or} \quad \begin{bmatrix} \alpha_{t_0} & \alpha_{t_0} t_0 \\ \alpha_{t_1} & \alpha_{t_1} t_1 \\ \vdots & \vdots \\ \alpha_{t_{n-1}} & \alpha_{t_{n-1}} t_{n-1} \end{bmatrix}$$
$$B = \begin{bmatrix} \beta_{t_0} \\ \beta_{t_1} \\ \vdots \\ \beta_{t_{n-1}} \end{bmatrix} \quad C = \begin{bmatrix} c_{t_0} \\ c_{t_1} \\ \vdots \\ c_{t_{n-1}} \end{bmatrix}$$

be matrices containing the parameter data and (in the case of A) years. If $\psi = [\psi_0 \ \psi_1 \ \psi_2]^T$ or $[\psi_0 \ \psi_1]^T$, then our modeling equations become

$$C = A\psi \quad \phi B = A\psi,$$

so we can write our objective function as

$$L = ||C - A\psi||^2 + ||\phi B - A\psi||^2$$

$$\begin{aligned}
&= (C - A\psi)^T(C - A\psi) + (\phi B - A\psi)^T(\phi B - A\psi) \\
&= C^T C - \psi^T A^T C - C^T A\psi + \phi^2 B^T B - \phi \psi^T A^T B \\
&\quad - \phi B^T A\psi + 2\psi^T A^T A\psi.
\end{aligned}$$

Each term in this sum is a scalar, so we can freely take the transpose of $C^T A\psi$ and $B^T A\psi$. This gives us

$$L = \phi^2 B^T B + C^T C - 2\psi^T A^T (\phi B + C) + 2\psi^T A^T A\psi.$$

Taking the derivative with respect to ψ and setting equal to 0 gives us

$$\begin{aligned}
\frac{dL}{d\psi} &= -2(\phi B^T + C^T) A\psi + 4\psi^T A^T A = 0 \\
4\psi^T A^T A &= 2(\phi B^T + C^T) A \\
\psi^T &= \frac{1}{2}(\phi B^T + C^T) A (A^T A)^{-1} \\
\psi &= \frac{1}{2}(A^T A)^{-1} A^T (\phi B + C)
\end{aligned}$$

Interestingly, this is the same estimate for ψ that we get from minimizing

$$\left\| \frac{\phi B + C}{2} - A\psi \right\|.$$

Now we implement. We first import modules and data.

```

1 import pandas as pd
2 import numpy as np
3 import numpy.linalg as la

```

We will import this file within `main.py`. To avoid running code on import, we run the analysis inside a function. First we define a couple helper functions. The first helper function raises a `TypeError` if `x` is not a `DataFrame`.

```

4 def check_df(x):
5     if not isinstance(x, pd.DataFrame):
6         msg = """
7 The second arguments of main() should be a Pandas
8 DataFrame. Right now one it is {0}
9 instead.\n""".format(type(x))
10    raise TypeError(msg)

```

Second helper function. It checks if `col` is a column in `x`.

```

11 def check_col(x, col):
12     if col not in x.columns:
13         msg = """
14 The third, fourth, and fifth arguments of main() should
15 be columns in data (second argument). However, it looks
16 like {0} is not a column in data.\n""".format(col)
17     raise KeyError(msg)

```

Table 7: (Proportional) Constants for Inverse Gamma

Variable	Value
ϕ	-0.1337
ψ_0	\$206,824
ψ_1	-\$105.33/year

Main function. We begin by checking arguments and disabling the chained assignment warning from Pandas.

```

18 def main(years, data, a_col, b_col, c_col):
19     check_df(data)
20     check_col(data, a_col)
21     check_col(data, b_col)
22     check_col(data, c_col)
23     pd.options.mode.chained_assignment = None

```

It will be easier to code everything if we save everything as arrays now.

```

24     alpha = data[a_col].to_numpy()
25     temp = data[a_col].to_numpy() * years.to_numpy()
26     A_lin = np.concatenate([np.ones([len(years), 1]),
27                             np.transpose([years]), np.transpose([alpha])], axis=1)
28     A_prop = np.concatenate([np.transpose([alpha]),
29                             np.transpose([temp])], axis=1)
30     B = data[b_col].to_numpy()
31     C = data[c_col].to_numpy()

```

Now estimate ϕ .

```

32     phi = np.sum(C) / np.sum(B)

```

And ψ .

```

33     psi_lin = 0.5 * la.inv(A_lin.transpose() @ A_lin) @ \
34         A_lin.transpose() @ (phi * B + C)
35     psi_prop = 0.5 * la.inv(A_prop.transpose() @ A_prop) @ \
36         A_prop.transpose() @ (phi * B + C)

```

Finally, we return the values.

```

37     return {
38         "linear":
39             {"phi": phi, "psi0": psi_lin[0],
40              "psi1": psi_lin[1], "psi2": psi_lin[2]},
41         "proportional":
42             {"phi": phi, "psi0": psi_prop[0], "psi1": psi_prop[1]}

```

Table 7 contains the values of constants from this estimation on the public CPS data.

File VI

bootstrap.py

This file implements the bootstrapping to find standard errors on the constant-shift-scale inverse-gamma distributions. The original approach was to implement a synthetic clustering and stratification design based on the synthetic survey design described in Jolliffe (2003). However, that approach appears to be problematic because it produces biased estimates of the parameter under the synthetic datasets. So this file will contain three functions: one for a naive bootstrap, one for a bootstrap with a synthetic survey design based on Jolliffe (2003), and one with a synthetic survey design based on Jolliffe (2003) except with a different set of clusters that results in unbiased parameter estimates. Eventually we will bound the standard error in the parameter between standard deviations from the first and third bootstrap functions above.

In this file, we will define the following functions:

- `bootstrap_naive(<estim>, <data>, <var>, <wgt>, <n>=100, **kwargs)`—function that implements a naive bootstrap using the variable `var` from `data` with `wgt` as the survey weights. The number of replicates `n` defaults to 200. The `kwargs` get fed to the estimation function.
- `bootstrap_Jol(<estim>, <data>, <var>, <wgt>, <strat>, <clust>, <n>=100, **kwargs)`—function that implements a bootstrap according to the synthetic survey design outlined in Jolliffe (2003). The `strat` and `clust` arguments should be names of columns in the data that we use for stratification and clustering. For income data, `strat` will be region or state, and `clust` will be household id.
- `bootstrap_Jol_sim(<estim>, <data>, <var>, <wgt>, <strat>, <cluster_size>=10, <n>=100, **kwargs)`—function that implements a simplified version of the Jolliffe (2003) synthetic survey design. In this case, the clusters will be collections of successive observations when we sort the dataset by `var`.

A minor technical note is that changes in the weights in the head and the tail can bias individual parameter estimates, i.e. create extreme outliers. To avoid issues from this problem, we use a normalized interquartile range to estimate the standard deviation of the parameter estimates.

We begin by importing modules and defining a few bookkeeping functions.

```
1 import numpy as np
2 import pandas as pd
3 def check_column(data, col):
4     if not isinstance(data, pd.DataFrame):
5         raise TypeError("Data should be a DataFrame.")
6     if col not in data.columns:
7         raise KeyError("{0} is not a column in the data.".format(col))
8 generator = np.random.default_rng(12345)
```

The bootstrap function is relatively straightforward.

```
9 def bootstrap_naive(F, df, var, wgt=None, n=100, **kwargs):
10     data = df.copy()
```

```

11  check_column(data, var)
12  #data.sort_values(var, inplace=True)

```

The `with_wgt` boolean will store whether we called the `bootstrap` function with a `wgt` column specified in the function argument.

```

13  if isinstance(wgt, type(None)):
14      with_wgt = False
15  else:
16      check_column(data, wgt)
17      with_wgt = True

```

Now add the synthetic weights to the `DataFrame`, and estimate parameters. If the user called `bootstrap` with `wgt` specified, then we set `lam` to the `wgt` column of `data`. Otherwise, we use all 1's for the weights. We create `new_weight` by simulating `n` Poisson distribution values for each row of the data, where the parameter of the Poisson distribution is `lam` (either 1 for `wgt` unspecified or the corresponding entry from `wgt` column otherwise).

```

18  if with_wgt:
19      data[wgt] = data[wgt] / data[wgt].sum() * len(data)
20      lam = data[wgt]
21  else:
22      lam = 1

```

The `generator.poisson` creates a matrix of Poisson realizations using a vector `lam` of parameters. Each element of `lam` determines the Poisson rate parameter used for the realizations in each column, and the elements of a single column are identically distributed. So to use this array in the definition of our `DataFrame`, we first transpose it.

```

23  new_weights = pd.DataFrame(
24      generator.poisson(lam=lam, size=[n,len(data)]).transpose(),
25      index=data.index,
26      columns=["_brw" + str(i) for i in range(n)])
27  data = pd.concat([data, new_weights], axis=1)
28  temp = [F(data, var=var, wgt="_brw"+str(i), **kwargs) for i in range(n)]
29  return temp

```

The function returns a list of dictionaries.

The `bootstrap_Jol()` function will calculate the bootstrap using a synthetic stratification and clustering scheme from Jolliffe (2003). The process contains several steps:

1. Find the strata by using `.unique()` method on the stratification column.
2. Within each stratum, create a list of total income for each household (distinct entries in the `clust` column) using `.groupby()`, and sort the household id's by income.
3. Use `np.arrange()` to add a column of integers to the `DataFrame` and integer-divide by 4 to get the synthetic cluster id's. These will go in the column `id`.
4. Loop through household identifiers, and for each household, add the corresponding cluster id to the `id` column in `data`.
5. Then we make the bootstrap replicate weights. Each replicate weight is a Poisson random variable. Choosing the Poisson parameter is a bit tricky. Let n_i denote the number of clusters in stratum i , and let w_i denote the total weight attached to stratum i . Let c_j be the total weight on cluster j in stratum i . Then c_j/w_i is the fraction of

stratum weight that corresponds to cluster j . From Wolter (2007), we know that we should think in terms of picking clusters for the bootstrap rather than individual observations, so we want to multiply this fraction by n_i to get the correct Poisson parameter for cluster j . Because

$$c_j = \sum_{k \text{ in cluster } j} p_k,$$

where p_k is the weight attached to observation k , we should use $n_i p_k / w_i$ for the Poisson parameters for each individual.

6. Replace each individual weight with the total weight of individuals in the same cluster.
7. Normalize the weight in each stratum to be the same as the original weight. The result will be the bootstrap replicate weights.

Once we have the new weights, we call the estimation function on the data with each set of synthetic weights. The function returns the results from each estimate in a list. We begin the function definition with some error checking and bookkeeping variables.

```
30 def bootstrap_Jol(F, df, var, wgt, strat, clust, n=100, **kwargs):
31     data = df.copy()
32     check_column(data, var)
33     check_column(data, strat)
34     check_column(data, clust)
35     check_column(data, wgt)
36     #data.sort_values(var, inplace=True)
37     household_incomes = {}
38     clusters_per_strat = {}
39     strata = data[strat].unique()
40     weight_per_strat = data.groupby(strat).sum()[wgt]
```

We add an `id` column containing household id, a column `_weight_per_strat` containing the total weight of the stratum corresponding to the current row, and a column `_clusters_per_strat` containing the total number of clusters in the stratum corresponding to the current row. We don't do checking of column names since this function is mostly for the current project, where we already know all the column names.

```
41     data["_id"] = 0
42     data["_weight_per_strat"] = 0
43     data["_clusters_per_strat"] = 0
```

Loop through the strata. Put the Series of household income levels in `household_incomes` dictionary. We add an `id` column that contains the cluster identifier.

```
44     for s in strata:
45         household_incomes[s] = \
46             data[data[strat] == s].groupby(clust).sum()[[var]].sort_values(var)
47         household_incomes[s]["_id"] = np.arange(len(household_incomes[s])) // 4
48         clusters_per_strat[s] = household_incomes[s]["_id"].max() + 1
49         data.loc[data[strat] == s, "_weight_per_strat"] = weight_per_strat[s]
50         data.loc[data[strat] == s, "_clusters_per_strat"] = clusters_per_strat[s]
```

Now loop through the list of households in the stratum, and add each household's cluster identifier to `data`. It would be simpler and faster to use `np.repeat()` here (as well as in a

number of other places in this function and the `Jol_sim` function). Unfortunately, I didn't know about `np.repeat()` when I wrote this file, and it's not worth the effort to redo it.

```
51     for h in household_incomes[s].index:
52         data.loc[data[clust] == h, "_id"] = household_incomes[s].loc[h, "_id"]
```

We are ready to create the new synthetic survey weights. Again `new_weights` is a DataFrame of Poisson distribution realizations, where the weights are the total weight on each synthetic cluster.

```
53     new_weights = pd.DataFrame(
54         generator.poisson(lam=(data[wgt] * data["_clusters_per_strat"] /
55             data["_weight_per_strat"])),
56         size=[n, len(data)]).transpose().astype(float),
57     index=data.index,
58     columns=["_brw" + str(i) for i in range(n)])
```

Now sum the Poisson weights across people in each synthetic cluster.

```
59     for s in strata:
60         for c in range(clusters_per_strat[s]):
61             new_weights.loc[(data[strat] == s) & (data["_id"] == c)] = \
62                 np.array(new_weights.loc[(data[strat] == s) &
63                     (data["_id"] == c)].sum(axis=0))
```

Normalize the total weight in each stratum to be the original weight in that stratum.

```
64     for s in strata:
65         new_weights[data[strat] == s] = (
66             new_weights[data[strat] == s]
67             * weight_per_strat[s]
68             / np.array(new_weights[data[strat] == s].sum()))
```

Finally, we concatenate the weights with the original data and call the estimation function on each set of bootstrap replicate weights.

```
69     data = pd.concat([data, new_weights], axis=1)
70     temp = [F(data, var=var, wgt="_brw"+str(i), **kwargs) for i in range(n)]
71     return temp
```

We return a list of estimates.

Using synthetic clusters adapted directly as it's written in the paper produces biased parameter estimates, so we use a simplified approach to get (what we hope to be?) an upper bound on the standard error in the parameter estimate. The `bootstrap_Jol_sim` function groups observations by stratum, then forms clusters by ranking observations by income and grouping 10 consecutive observations, since clusters in the CPS have about four households on average, and the average household size in the U.S. is about 2.5 people. The `clust_size` argument stores the cluster size.

```
72 def bootstrap_Jol_sim(F, df, var, wgt, strat, clust_size=10, n=100, **kwargs):
73     data = df.copy()
74     check_column(data, var)
75     check_column(data, strat)
76     check_column(data, wgt)
77     #data.sort_values(var, inplace=True)
```

```

78 clusters_per_strat = {}
79 strata = data[strat].unique()
80 weight_per_strat = data.groupby(strat).sum()[wgt]

```

Add a `_weight_per_strat` and `_clusters_per_strat` columns containing weight and number of clusters per stratum. We don't need household id's this time since we aren't using household information, but we will still add an `_id` column to keep track of the synthetic clusters.

```

81 data["_id"] = 0
82 data["_weight_per_strat"] = 0
83 data["_clusters_per_strat"] = 0

```

Loop through the strata. Put the Series of household income levels in `household_incomes` dictionary. We set the `id` column to contain the cluster identifier.

```

84 for s in strata:
85     temp = data[data[strat] == s].sort_values(var)
86     temp["_id"] = np.arange(len(temp)) // clust_size
87     clusters_per_strat[s] = temp["_id"].max() + 1
88     data.loc[data[strat] == s, "_id"] = temp["_id"]
89     data.loc[data[strat] == s, "_weight_per_strat"] = weight_per_strat[s]
90     data.loc[data[strat] == s, "_clusters_per_strat"] = clusters_per_strat[s]

```

We are ready to create the new synthetic survey weights. For each individual, we replace that individual's weight by the sum of all weights of individuals in the same cluster. Again `new_weights` is a DataFrame of Poisson distribution realizations, where the weights are the total weight on each synthetic cluster.

```

91 new_weights = pd.DataFrame(
92     generator.poisson(lam=(data[wgt] * data["_clusters_per_strat"] /
93         data["_weight_per_strat"])),
94     size=[n,len(data)]).transpose().astype(float),
95     index=data.index,
96     columns=["_brw" + str(i) for i in range(n)])

```

Now sum the Poisson weights across people in each synthetic cluster.

```

97 for s in strata:
98     for c in range(clusters_per_strat[s]):
99         new_weights.loc[(data[strat] == s) & (data["_id"] == c)] = \
100             np.array(new_weights.loc[(data[strat] == s) &
101                 (data["_id"] == c)].sum(axis=0))

```

Normalize the total weight in each stratum to be the original weight in that stratum.

```

102 for s in strata:
103     new_weights[data[strat] == s] = (
104         new_weights[data[strat] == s]
105         * weight_per_strat[s]
106         / np.array(new_weights[data[strat] == s].sum()))

```

As in the other two functions, we concatenate the weights and data and call the estimation function using each set of bootstrap replicate weights.

```

107 data = pd.concat([data, new_weights], axis=1)

```

```
108     temp = [F(data, var=var, wgt="_brw"+str(i), **kwargs) for i in range(n)]
109     return temp
```

Again the result will be a list of estimates (dictionaries in this context).

File VII

make_figures.py

This file contains the code to make figures. All handling of Matplotlib and Pyplot commands happens here, and `main.py` calls the functions and provides the data. We define several graphics-producing functions:

- `single_graph`—makes a single graph
- `triple_graph`—makes three panels where the first two are next to each other, and the third is below them and centered
- `lin_loglog_graphs`—makes three rows of two graphs each, where the first column is linear and the second column is loglog scaling
- `lin_graphs`—makes six linearly scaled graphs
- `loglog_graphs`—makes six loglog scaled graphs
- `four_graphs`—makes four graphs in 2x2 layout; hook for extra code
- `single_graph_ext`—makes a single graph but with more control over the contents; similar to `four_graphs` in its implementation/interface

All of these functions use the `plt.savefig()` function to save the figure directly as a pdf at the correct size. The first thing to do is import modules.

```
1 import bin
2 import matplotlib as mpl
3 import matplotlib.pyplot as plt
4 import numpy as np
5 import pandas as pd
```

Some rcParams to make the figures look nice. You will need a working TeX distribution on your machine to produce figures, or you may have to rewrite some figure titles and labels.

```
6 mpl.rcParams["legend.fontsize"] = "small"
7 mpl.rcParams["font.family"] = "serif"
8 mpl.rcParams["text.usetex"] = True      # tells Pyplot to use TeX
9 mpl.rcParams["figure.constrained_layout.use"] = True
10 mpl.rcParams["savefig.dpi"] = 300
11 mpl.rcParams["savefig.format"] = "pdf"
```

Bookkeeping functions.

```
12 def check_var(data, var):
13     if var not in data:
14         raise KeyError("{0} is not a column in the data".format(var))
15 def check_list(x):
16     is_list_like = hasattr(x, "__len__") and \
17                     hasattr(x, "__getitem__") and \
18                     hasattr(x, "__iter__")
19     if not is_list_like or isinstance(x, str):
20         raise TypeError("Please use list instead of {0} for {1}".format(type(x), x))
```

We begin by creating a function to show a single line graph. We will use this to create the

figure of constant-shift-scale inverse-gamma parameters.

```
21 def single_graph(data, parameter, *, filename, title=None):
22     plt.close(plt.gcf())
23     if isinstance(title, type(None)):
24         called_with_title = False
25     else:
26         called_with_title = True
27     check_var(data, parameter)
28     plt.plot(data.index, data[parameter], c="black", lw=0.5)
29     if called_with_title:
30         plt.title(title)
31     #plt.show()
32     plt.gcf().set_size_inches(3.25, 2.5)
33     plt.savefig(filename)
34     plt.close(plt.gcf())
```

The next function creates three panels (for showing inverse-gamma parameter estimates). We accept one DataFrame (of inverse-gamma parameters), parameter names, and titles for the subplots.

```
35 def triple_graph(data, shape, scale, shift, *, filename, titles=[]):
36     plt.close(plt.gcf())
37     check_list(titles)
38     if len(titles) > 0 and len(titles) < 3:
39         raise ValueError("Please specify zero or all titles")
40     elif len(titles) == 0:
41         called_with_titles = False
42     else:
43         called_with_titles = True
44     check_var(data, shape)
45     check_var(data, scale)
46     check_var(data, shift)
```

Now fill the first three subplots. We use a GridSpec object to get the positioning right.

```
47     plt.figure()
48     grid = mpl.gridspec.GridSpec(4, 4, figure=plt.gcf())
49     grid_boxes = [grid[0:2, 0:2], grid[0:2, 2:4], grid[2:4, 1:3]]
50     parameter_names = [shape, scale, shift]
51     for i in range(3):
52         plt.subplot(grid_boxes[i])
53         plt.plot(data.index, data[parameter_names[i]], lw=0.5, c="black")
54         if called_with_titles:
55             plt.title(titles[i])
```

Plot the figure.

```
56     #plt.tight_layout()
57     #plt.show()
58     plt.gcf().set_size_inches(6.5, 5)
59     plt.savefig(filename)
60     plt.close(plt.gcf())
```

The `lin_loglog_graphs` function is more complicated. It makes a 3x2 plot of three rows of sample and model densities in linear and log-log scales. Arguments of the function are

- Three DataFrames
- `var` (and `wgt=None`) columns
- Several lists of things, where we plot everything from each list on one row of the subplot array
 - `list_F1`, `list_F2`, and `list_F3`—density functions, which take a single argument and return a (nonnegative) real number
 - `list_c1`, `list_c2`, and `list_c3`—minimum values of the support
 - `list_opt1`, `list_opt2`, and `list_opt3`—list of dictionaries containing keyword-options for plotting each density function
- Bounds on the graphs
- A filename for saving the figure
- An optional list of titles and boolean to determine if using a legend

For the lists of functions, constants, and options, we will plot each list on one row of the subplot array. The left column is a linear scaling, and the right column is a loglog scaling. We need to use lists of functions and constants so that we can make a figure where each row plots multiple functions. We begin with error checking and setting the `called_with_titles` boolean.

```

61 def lin_loglog_graphs(data1, data2, data3, var,
62     list_F1, list_c1, list_opt1,
63     list_F2, list_c2, list_opt2,
64     list_F3, list_c3, list_opt3,
65     x_bounds_lin, y_bounds_lin, x_bounds_log, y_bounds_log, *,
66     filename, wgt=None, titles=[], with_legend=False):
67     plt.close(plt.gcf())
68     check_list(titles)
69     if len(titles) > 0 and len(titles) < 6:
70         raise ValueError("Please specify zero or all titles")
71     elif len(titles) == 0:
72         called_with_titles = False
73     else:
74         called_with_titles = True

```

We loop through the rows of the subplot. For each row, we store the lists of density functions, constants, and option dictionaries in new variables. Then we check that all of them are in fact a list. We also make a new variable for the data, check that `var` is a column, and feed the data to `bin.bin_data`.

```

75     for i in range(1, 4):
76         list_F = eval("list_F" + str(i))
77         list_c = eval("list_c" + str(i))
78         list_opt = eval("list_opt" + str(i))
79         check_list(list_F)
80         check_list(list_c)
81         check_list(list_opt)

```

```

82     d = eval("data" + str(i))
83     check_var(d, var)
84     binned_data = bin.bin_data(d, var, wgt)

```

Start with the linear plot. We add a title, a curve for the density, and sample density points. The `x_vals` list starts with `c_i` so that it is clear on the figure where the support of the model ends. We loop through the lists of functions, constants, and options concurrently and plot each one separately on the current plot.

```

85     plt.subplot(3, 2, 2 * i - 1)
86     if called_with_titles:
87         plt.title(titles[2 * i - 2])
88     for F, c, opt in zip(list_F, list_c, list_opt):
89         x_vals = np.linspace(c, x_bounds_lin[1], 200)
90         y_vals = list(map(F, x_vals))
91         plt.plot(x_vals, y_vals, **opt)
92     plt.scatter(binned_data["mid"], binned_data["dens"], s=2, c="blue")
93     plt.xlim(x_bounds_lin)
94     plt.ylim(y_bounds_lin)
95     if with_legend:
96         plt.legend(loc="upper right")

```

Now do the loglog plot. The code is similar except that this time, the `x_vals` list runs for the entire length of the horizontal axis.

```

97     plt.subplot(3, 2, 2 * i)
98     if called_with_titles:
99         plt.title(titles[2 * i - 1])
100    x_vals = np.geomspace(*x_bounds_log, 200)
101    for F, c, opt in zip(list_F, list_c, list_opt):
102        y_vals = list(map(F, x_vals))
103        plt.plot(x_vals, y_vals, **opt)
104    plt.scatter(binned_data["mid"], binned_data["dens"], s=2, c="blue")
105    plt.xlim(x_bounds_log)
106    plt.ylim(y_bounds_log)
107    plt.loglog()
108    if with_legend:
109        plt.legend(loc="lower left")

```

After the for-loop, plot the figure.

```

110    #plt.show()
111    plt.gcf().set_size_inches(6.5, 7.5)
112    plt.savefig(filename)
113    plt.close(plt.gcf())

```

The `lin_graphs` and `loglog_graphs` functions will be the same as `lin_loglog_graphs` except that they take more data and scale their plots all the same way.

```

114 def lin_graphs(data1, data2, data3,
115                data4, data5, data6, var,
116                list_F1, list_c1, list_opt1,
117                list_F2, list_c2, list_opt2,

```



```

118     list_F3, list_c3, list_opt3,
119     list_F4, list_c4, list_opt4,
120     list_F5, list_c5, list_opt5,
121     list_F6, list_c6, list_opt6,
122     x_bounds, y_bounds, *,
123     filename, wgt=None, titles=[], with_legend=True):
124 plt.close(plt.gcf())
125 check_list(titles)
126 if len(titles) > 0 and len(titles) < 6:
127     raise ValueError("Please specify zero or all titles")
128 elif len(titles) == 0:
129     called_with_titles = False
130 else:
131     called_with_titles = True

```

Now loop through the panels, create pointers, and check list properties.

```

132 for i in range(1, 7):
133     list_F = eval("list_F" + str(i))
134     list_c = eval("list_c" + str(i))
135     list_opt = eval("list_opt" + str(i))
136     check_list(list_F)
137     check_list(list_c)
138     check_list(list_opt)
139     d = eval("data" + str(i))
140     check_var(d, var)
141     binned_data = bin.bin_data(d, var, wgt)

```

Plot the figures.

```

142 plt.subplot(3, 2, i)
143 if called_with_titles:
144     plt.title(titles[i - 1])
145 for F, c, opt in zip(list_F, list_c, list_opt):
146     x_vals = np.linspace(c, x_bounds[1], 200)
147     y_vals = list(map(F, x_vals))
148     plt.plot(x_vals, y_vals, **opt)
149 plt.scatter(binned_data["mid"], binned_data["dens"], s=2, c="blue")
150 plt.xlim(x_bounds)
151 plt.ylim(y_bounds)
152 if with_legend:
153     plt.legend(loc="upper right")

```

After the for-loop, plot the figure.

```

154 #plt.show()
155 plt.gcf().set_size_inches(6.5, 7.5)
156 plt.savefig(filename)
157 plt.close(plt.gcf())

```

Same but with loglog plots.

```

158 def loglog_graphs(data1, data2, data3,
159     data4, data5, data6, var,

```

```

160     list_F1, list_opt1,
161     list_F2, list_opt2,
162     list_F3, list_opt3,
163     list_F4, list_opt4,
164     list_F5, list_opt5,
165     list_F6, list_opt6,
166     x_bounds, y_bounds, *,
167     filename, wgt=None, titles=[], with_legend=True):
168 plt.close(plt.gcf())
169 check_list(titles)
170 if len(titles) > 0 and len(titles) < 6:
171     raise ValueError("Please specify zero or all titles")
172 elif len(titles) == 0:
173     called_with_titles = False
174 else:
175     called_with_titles = True

```

Now loop through the panels, create pointers, and check list properties.

```

176 for i in range(1, 7):
177     list_F = eval("list_F" + str(i))
178     list_opt = eval("list_opt" + str(i))
179     check_list(list_F)
180     check_list(list_opt)
181     d = eval("data" + str(i))
182     check_var(d, var)
183     binned_data = bin.bin_data(d, var, wgt)

```

Create the figures.

```

184     plt.subplot(3, 2, i)
185     if called_with_titles:
186         plt.title(titles[i - 1])
187     for F, opt in zip(list_F, list_opt):
188         x_vals = np.geomspace(*x_bounds, 200)
189         y_vals = list(map(F, x_vals))
190         plt.plot(x_vals, y_vals, **opt)
191     plt.scatter(binned_data["mid"], binned_data["dens"], s=2, c="blue")
192     plt.xlim(x_bounds)
193     plt.ylim(y_bounds)
194     plt.loglog()
195     if with_legend:
196         plt.legend(loc="lower left")

```

After the for-loop, plot the figure.

```

197     #plt.show()
198     plt.gcf().set_size_inches(6.5, 7.5)
199     plt.savefig(filename)
200     plt.close(plt.gcf())

```

The `four_graphs` is similar to the `lin_loglog_graphs`, etc. in that it accepts lists of information for creating multiple figures on the subplot. It is slightly lower-level than previous

plotting functions in that we feed it actual data rather than a DataFrame and column. The arguments are

- `x_data` for the horizontal axes
- Several lists of things, one for each panel of the figure
 1. `list_y1`, `list_y2`, `list_y3`, `list_y4`—lists of data for the vertical axis
 2. `list_plot1`, `list_plot2`, `list_plot3`, `list_plot4`—lists of plotting functions. Most of these will be `plt.plot()`
 3. `list_opt1`, `list_opt2`, `list_opt3`, `list_opt4`—lists of options for plotting functions
- A filename for the figure
- An optional list of titles for the subgraphs

We begin with error checking and setting the `called_with_titles` boolean.

```
201 def four_graphs(  
202     list_x1, list_y1, list_plot1, list_opt1,  
203     list_x2, list_y2, list_plot2, list_opt2,  
204     list_x3, list_y3, list_plot3, list_opt3,  
205     list_x4, list_y4, list_plot4, list_opt4, *,  
206     filename, titles=[], extra_code=""):  
207     plt.close(plt.gcf())  
208     check_list(titles)  
209     if len(titles) > 0 and len(titles) < 4:  
210         raise ValueError("Please specify zero or all titles")  
211     elif len(titles) == 0:  
212         called_with_titles = False  
213     else:  
214         called_with_titles = True
```

We loop through the four panels. On each iteration, we first create new pointers to the lists for that iteration and check that they are actually lists. Then we loop through the three lists of *y*-axis data, plotting functions, and options. We call the plotting function with the *x*-axis data, corresponding *y*-axis data, and corresponding plotting options.

```
215     for i in range(1,5):  
216         list_x = eval("list_x" + str(i))  
217         list_y = eval("list_y" + str(i))  
218         list_plot = eval("list_plot" + str(i))  
219         list_opt = eval("list_opt" + str(i))  
220         check_list(list_x)  
221         check_list(list_y)  
222         check_list(list_plot)  
223         check_list(list_opt)
```

Now plot the data for this subplot.

```
224     plt.subplot(2, 2, i)  
225     for x, y, plot, opt in zip(list_x, list_y, list_plot, list_opt):  
226         plot(x, y, **opt)  
227     if called_with_titles:
```

```
228     plt.title(titles[i-1])
```

Execute any extra code. Used to provide a hook into the function.

```
229     exec(extra_code)
```

Then save the figure.

```
230     #plt.show()
```

```
231     plt.gcf().set_size_inches(6.5, 5)
```

```
232     plt.savefig(filename)
```

```
233     plt.close(plt.gcf())
```

A single-panel with more control. This function uses the same interface as the `four_graphs` function.

```
234 def single_graph_ext(
```

```
235     list_x, list_y, list_plot, list_opt, *,
```

```
236     filename, title=None, extra_code=""):
```

```
237     plt.close(plt.gcf())
```

```
238     if isinstance(title, type(None)):
```

```
239         called_with_title = False
```

```
240     else:
```

```
241         called_with_title = True
```

Now actually make the graph and execute the extra code.

```
242     for x, y, plot, opt in zip(list_x, list_y, list_plot, list_opt):
```

```
243         plot(x, y, **opt)
```

```
244     if called_with_title:
```

```
245         plt.title(title)
```

```
246     exec(extra_code)
```

And save the figure.

```
247     #plt.show()
```

```
248     plt.gcf().set_size_inches(3.25, 2.5)
```

```
249     plt.savefig(filename)
```

```
250     plt.close(plt.gcf())
```

File VIII

main.py

This file contains the main code to carry out the estimations and write parameters to files. We have several tasks:

1. Load the data files that we made through `gen_files.py`. Make sure to drop the incomes of \$0!
2. Estimate the parameters of all models for 2023, store the parameter estimates, and write them to files.
3. Estimate parameters, Fisk, inverse-gamma, and constant-shift-scale inverse-gamma distributions in all years.
4. Bootstrap the standard errors for the constant-shift-scale inverse-gamma parameter estimates.
5. Make figures.

- List `dist` of distribution strings. See table 6 for the distribution strings.
- List `years` of years. Both `dist` and `years` are for reference.
- Dictionary `data_min_max` of minimum and maximum cutoff values for data in each year. The keys are years, and the values are 2-tuples of boundary values.
- Dictionary `data` with years as keys and DataFrames as values.
- Dictionary `ecdfs` with years as keys and ecdfs as values.
- DataFrames of `parameters` containing the parameter estimates for the (shifted, constant-shift, and constant-shift-scale) inverse-gamma distributions. The columns are parameters, and the rows are years.
- Several `CSS_bootstrap` DataFrames containing synthetic parameter estimates and estimates of standard errors.
- `trim_data(<data>, <var>, <year>)`—trims the data according to the information in `data_min_max` and returns the trimmed DataFrame.

Several switches control which portion of this file runs. Make sure to set the appropriate switches to `True`.

```
1 ##### Switches! #####
2 do_load_data = False
3 do_2023_short = False
4 do_2023_long = False
5 do_Fisk = False
6 do_InvG = False
7 do_CS_InvG = False    # <--- present for historical reasons
8 do_CSS_InvG = False
9 do_bootstrap = False
10 do_figures = False
11 do_test = False
```

```

12 #####
13 if not (do_load_data
14         or do_2023_short
15         or do_2023_long
16         or do_InvG
17         or do_Fisk
18         or do_CS_InvG
19         or do_CSS_InvG
20         or do_bootstrap
21         or do_figures
22         or do_test):
23     print()
24     print("Warning. This run of estimate_parameters will not do anything.")
25     print("To enable data analysis, set one or more do_ booleans to True.")
26     print()
27     quit()

```

We begin by importing modules, creating the list of distribution strings, and creating the list of years. We need Pandas and the file `estimate_params.py`. See also Table 6. We have data for all years between 1967 and 2023 with the exception of 1970.

```

28 import bin
29 import bootstrap
30 import check_constants as cc
31 import estimate_parameters as est
32 import matplotlib as mpl
33 import matplotlib.pyplot as plt
34 import make_figures
35 import numpy as np
36 import pandas as pd
37 import scipy.special as spec
38 import time
39 dist = ["GB2", "Dagum", "Burr", "Fisk", "InvG", "CS_InvG", "CSS_InvG",
40         "Davis", "LogN_P_cut", "LogN_P_mix"]
41 dist_names = {
42     "GB2": "Generalized Beta, type II",
43     "Dagum": "Dagum",
44     "Burr": "Burr (Singh-Maddala)",
45     "Fisk": "Fisk",
46     "InvG": "Inverse Gamma",
47     "CS_InvG": "Constant-Shift Inverse Gamma",
48     "CSS_InvG": "Constant-Shift-Scale Inverse Gamma",
49     "Davis": "Davis",
50     "LogN_P_cut": "Log-Normal, Pareto Cutoff",
51     "LogN_P_mix": "Log-Normal, Pareto Mixture"}
52 years = [i for i in range(1967, 2024) if (i != 1970)]
53 data = {}
54 ecdfs = {}

```

A function to print the time.

```
55 def the_time():  
56     print("The time is", time.asctime())
```

The maximum and minimum values are by inspection.

```
57 dat_max_min = {  
58     1967:(-4000,60000),  
59     1968:(-4000,90000),  
60     1969:(-4000,90000),  
61     1971:(-9000,90000),  
62     1972:(-9000,90000),  
63     1973:(-9000,90000),  
64     1974:(-9000,90000),  
65     1975:(-9000,90000),  
66     1976:(-9000,90000),  
67     1977:(-9000,100000),  
68     1978:(-9000,100000),  
69     1979:(-9000,100000),  
70     1980:(-9000,100000),  
71     1981:(-9000,100000),  
72     1982:(-9000,120000),  
73     1983:(-9000,120000),  
74     1984:(-9000,120000),  
75     1985:(-9000,200000),  
76     1986:(-9000,150000),  
77     1987:(-9000,200000),  
78     1988:(-9000,200000),  
79     1989:(-9000,200000),  
80     1990:(-9000,200000),  
81     1991:(-9000,200000),  
82     1992:(-9000,200000),  
83     1993:(-9000,200000),  
84     1994:(-9000,250000),  
85     1995:(-9000,250000),  
86     1996:(-8000,450000),  
87     1997:(-4000,450000),  
88     1998:(-9000,450000),  
89     1999:(-9000,450000),  
90     2000:(-9000,350000),  
91     2001:(-9000,450000),  
92     2002:(-11000,450000),  
93     2003:(-11000,550000),  
94     2004:(-11000,550000),  
95     2005:(-11000,650000),  
96     2006:(-11000,550000),  
97     2007:(-9000,650000),  
98     2008:(-9000,650000),  
99     2009:(-9000,550000),
```

```

100 2010:(-9000,550000),
101 2011:(-9000,1000000),
102 2012:(-9000,1000000),
103 2013:(-9000,1000000),
104 2014:(-9000,1000000),
105 2015:(-9000,1000000),
106 2016:(-9000,1000000),
107 2017:(-9000,1000000),
108 2018:(-9000,1000000),
109 2019:(-9000,1000000),
110 2020:(-9000,1000000),
111 2021:(-9000,1000000),
112 2022:(-9000,1000000),
113 2023:(-9000,1000000)}

```

The `parameters` dictionaries will hold a `DataFrame` to store the (shifted, constant-shift, and constant-shift-scale) inverse-gamma parameter estimates. We initialize them to be empty `DataFrames`. By specifying the column information at initialization, we can more easily concatenate parameter estimates later. See Tables 1 and 6 for discussion of the notation for different parameters. We store constants in the `CS_InvG` and `CSS_InvG` parameter tables for reference. The `CSS_bootstrap` `DataFrames` will store the synthetic parameter estimates from each bootstrapping function, and `CSS_bootstrap_se` stores the standard error estimates for each year and type of bootstrap.

```

114 Fisk_parameters = pd.DataFrame(index=years, columns=["alpha", "beta", "c"])
115 InvG_parameters = pd.DataFrame(index=years, columns=["alpha", "beta", "c"])
116 CS_InvG_parameters = pd.DataFrame({"alpha": [], "beta": [], "phi": []})
117 CSS_InvG_parameters = pd.DataFrame(index=years, columns=["alpha"])
118 CSS_bootstrap_naive = pd.DataFrame()
119 CSS_bootstrap_Jol = pd.DataFrame()
120 CSS_bootstrap_Jol_sim = pd.DataFrame()
121 CSS_bootstrap_se = pd.DataFrame(index=years)

```

The function `trim_data` accepts a year as its argument. It selects the corresponding data from `raw_data`, slices it, and returns the trimmed version.

```

122 def trim_data(data, var, year):
123     if not isinstance(data, pd.DataFrame):
124         raise TypeError("The data should be a DataFrame")
125     if var not in data.columns:
126         raise KeyError("Variable of interest is not a column in the data")
127     bounds = dat_max_min[year]
128     return data[(data[var] >= bounds[0]) &
129                (data[var] <= bounds[1])].reset_index(drop=True).dropna()

```

Results

First we load all the data files, and then we calculate the parameter estimates for 2023. See Table 8 for approximate lengths of time each estimator takes. The time depends on the computer's processing speed and may vary between machines.

Table 8: Approximate Time to Perform Each Estimation

Model	Approximate Time to Estimate
GB2	Several hours
Dagum	Two minutes
Burr	Two minutes
Fisk	One minute
InvG	Less than one minute
Davis	One minute
CS_InvG	Less than one minute
CSS_InvG	Two minutes
LogN_P_cut	Two minutes
LogN_P_mix	One day

Load all the data files. We trim the ends according to `dat_max_min`. (Don't forget to remove the 0's!)

```

130 #####
131 ## Load Data ##
132 #####
133 print()
134 the_time()
135 if do_load_data:
136     print("Loading data:")
137     for i in years:
138         print("Year {0}".format(i))
139         temp_data = pd.read_csv("data_{0}.txt".format(i), header=0) # read file
140         temp_data = trim_data(temp_data, "income", i) # trim ends
141         temp_data = temp_data[temp_data["income"] != 0] # remove 0s
142         temp_data = temp_data[temp_data["weight"] >= 0] # weight > 0
143         data[i] = temp_data
144         ecdfs[i] = est.make_ecdf(data[i], "income", "weight")
145 else:
146     print("Skipping loading data")
147 print()

```

Now find parameters for 2023 data. We handle the constant-scale and constant-shift-scale inverse-gamma distributions separately. If we loaded the 2023 data previously, we reference it from `data`. Otherwise, we load it now.

```

148 #####
149 ## 2023 Data ##
150 #####
151 the_time()
152 if do_2023_short or do_2023_long:
153     if 2023 in data: # if loaded data_2023.txt
154         data_2023 = data[2023]

```

```

155 else:                # otherwise, load it now
156     data_2023 = pd.read_csv("data_2023.txt", header=0)
157     data_2023 = trim_data(data_2023, "income", 2023)
158     data_2023 = data_2023[data_2023["income"] != 0]
159     data_2023 = data_2023[data_2023["weight"] >= 0]
160     print("Estimating parameters for 2023 data:")

```

First the case of the “short” distributions. (Where the time to run the estimation algorithm is short.) We open the file `2023_parameters_short.txt` for writing, loop through the models, and write in the file the parameter estimate for 2023 data under each model.

```

161 if do_2023_short:
162     f = open("2023_parameters_short.txt", "w")
163     for model in dist:
164         if model != "CS_InvG" and model != "CSS_InvG" and \
165             model != "LogN_P_mix" and model != "GB2":
166             print(model)
167             f.write("---{0}---\n".format(model))
168             temp = est.estimate(data_2023, "income", "weight")
169             for i in temp:
170                 f.write("{0}: {1}\n".format(i, temp[i]))
171     f.close()
172 else:
173     print("Skipping Dagum, Burr, Fisk, InvG, Davis, and LogN_P_cut")

```

The “long” distributions are those whose estimation algorithms take a long time, specifically GB2 and LogN_P_mix.

```

174 if do_2023_long:
175     f = open("2023_parameters_long.txt", "w")
176     for model in ["GB2", "LogN_P_mix"]:
177         print(model)
178         f.write("---{0}---\n".format(model))
179         temp = est.estimate(data_2023, "income", "weight")
180         for i in temp:
181             f.write("{0}: {1}\n".format(i, temp[i]))
182     f.close()
183 else:
184     print("Skipping GB2 and LogN_P_mix")
185 else:
186     print("Skipping estimating 2023 data")
187 print()

```

Time for the Fisk distribution. We estimate the parameters for all years of data. First we check that `do_load_data` is true. Then we loop through `years` and estimate parameters for that year. We store them in `Fisk_parameters` DataFrame, and after the loop, we write the results to `Fisk_parameters.txt`.

```

188 #####
189 ## Fisk ##
190 #####
191 the_time()

```

```

192 if do_Fisk:
193     if not do_load_data:
194         raise RuntimeError("Please load data before estimating Fisk")
195     print("Estimating Fisk parameters:")
196     for i in years:
197         print("Year {0}".format(i))
198         p = est.estimate["Fisk"](data[i], "income", "weight",
199             ecdf=ecdfs[i], x="income", y="weight")["parameters"]
200         Fisk_parameters.loc[i] = p
201     Fisk_parameters.to_csv("Fisk_parameters.txt")
202 else:
203     print("Skipping Fisk estimation")
204 print()

```

Now calculate the parameter estimates for all years of data using inverse gamma. Same approach as with Fisk distribution. First we check that `do_load_data` is true. Then we loop through years. For each year, we estimate inverse-gamma parameters for that year and store them in `InvG_parameters` DataFrame. After the loop, we write the results to `InvG_parameters.txt`.

```

205 #####
206 ## Inverse Gamma ##
207 #####
208 the_time()
209 if do_InvG:
210     if not do_load_data:
211         raise RuntimeError("Please load data before estimating InvG")
212     print("Estimating inverse-gamma parameters:")
213     for i in years:
214         print("Year {0}".format(i))
215         p = est.estimate["InvG"](data[i], "income", "weight",
216             ecdf=ecdfs[i], x="income", y="weight")["parameters"]
217         InvG_parameters.loc[i] = p
218     InvG_parameters.to_csv("InvG_parameters.txt")
219 else:
220     print("Skipping inverse-gamma estimation")
221 print()

```

For historical reasons, we include code to estimate parameters for the constant-shift inverse-gamma model, where we have $c = \phi\beta$. We don't use this code for anything. It involves the formula

$$\phi = \sum_{i=1}^n \beta_i c_i / \sum_{i=1}^n \beta_i^2,$$

which minimizes the total squared distance between β and ϕc . This comes out to be approximately $\phi = -0.13$. The procedure is the same as for inverse-gamma distribution: loop through years, estimate the parameters for that year's data, and store the results in `CS_InvG_parameters` DataFrame.

```

222 #####

```

```

223 ## Constant-shift Inverse Gamma ##
224 #####
225 the_time()
226 if do_CS_InvG:
227     if not do_load_data:
228         raise RuntimeError("Please load data before estimating CS_InvG")
229     print("Estimating constant-scale inverse-gamma parameters:")
230     phi = (InvG_parameters["beta"] * InvG_parameters["c"]).sum() / \
231         (InvG_parameters["beta"] ** 2).sum()
232     for i in years:
233         print("Year {0}".format(i))
234         p = est.estimate["CS_InvG"](data[i], phi, "income", "weight",
235             ecdf=ecdfs[i], x="income", y="weight")["parameters"]
236         CS_InvG_parameters = pd.concat([CS_InvG_parameters,
237             pd.DataFrame({"alpha": p[0], "beta": p[1], "phi": phi}, index=[i])])
238     CS_InvG_parameters.to_csv("CS_InvG_parameters.txt")
239 else:
240     print("Skipping constant-shift inverse-gamma estimation")
241 print()

```

For the constant-shift-scale case, we need to calculate constants before we can estimate it. So we import `check_constants` and call the `main()` function from that file. We store the inverse-gamma parameter estimates for each year in `InvG_parameters` DataFrame.

```

242 #####
243 ## Constant-shift-scale Inverse Gamma ##
244 #####
245 the_time()
246 if do_CSS_InvG:
247     if not do_load_data:
248         raise RuntimeError("Please load data before estimating CSS_InvG")
249     print("Checking constants")
250     InvG_parameters = pd.read_csv("InvG_parameters.txt", header=0, index_col=0)
251     temp = cc.main(InvG_parameters.index, InvG_parameters,
252         "alpha", "beta", "c")["linear"]
253     phi = temp["phi"]
254     psi = [temp["psi0"], temp["psi1"], temp["psi2"]]

```

Then we write the values in a file.

```

255     f = open("CSS_InvG_constants.txt", "w")
256     for i in temp:
257         f.write("{0},{1}\n".format(i, temp[i]))
258     f.close()

```

Now we estimate parameters for the constant-shift-scale inverse-gamma distribution. Similar to Fisk and InvG, we store the end results in `CSS_InvG_parameters.txt`.

```

259     print("Estimating constant-shift-scale inverse-gamma parameters:")
260     for i in years:
261         print("Year {0}".format(i))

```

```

262     p = est.estimated["CSS_InvG"](data[i], i, phi, psi,
263         InvG_parameters.loc[i, "alpha"], "income", "weight", ecdf=ecdfs[i],
264         x="income", y="weight")["parameters"]
265     CSS_InvG_parameters.loc[i] = p
266     CSS_InvG_parameters.to_csv("CSS_InvG_prop_parameters.txt")
267 else:
268     print("Skipping constant-shift-scale inverse-gamma distribution")
269 print()

```

Finally, it is bootstrapping time. First, we read in the inverse-gamma parameter estimates and the variables from `check_constants.py`.

```

270 #####
271 ## bootstraps ##
272 #####
273 the_time()
274 if do_bootstrap:
275     if not do_load_data:
276         raise RuntimeError("Please load data before bootstrapping")
277     print("Bootstrapping standard errors:")
278     InvG_parameters = pd.read_csv("InvG_parameters.txt", header=0, index_col=0)
279     f = open("CSS_InvG_constants.txt")
280     for line in f:
281         temp = line[:-1].split(",")
282         exec("{0} = {1}".format(temp[0], temp[1]))
283     f.close()

```

For each year, we take the data for that year and feed it into one of the three bootstrapping functions. We will store the `parameter` column values in a DataFrame. The `CSS_pos_args` list and `CSS_named_args` dictionary pack common arguments of the bootstrap functions for convenience. We manually set the index of each DataFrame to be a list of integers if we want to be able to easily add the synthetic parameter estimates as columns. If we change the number of estimates in the bootstrap, we will need to change the index to a different range or list.

```

284     CSS_bootstrap_naive.index = range(100)
285     CSS_bootstrap_Jol.index = range(100)
286     CSS_bootstrap_Jol_sim.index = range(100)
287     for i in years:
288         CSS_pos_args = [est.estimated["CSS_InvG"], data[i], "income", "weight"]
289         CSS_named_args = {"t": i, "phi": phi, "psi": [psi0, psi1],
290             "a0": InvG_parameters.loc[i, "alpha"]}
291         print("Year {0}".format(i))

```

Naive bootstrap is first.

```

292     temp = pd.DataFrame(bootstrap.bootstrap_naive(
293         *CSS_pos_args,
294         **CSS_named_args)
295         )["parameters"]
296     temp = temp.explode().astype(np.float64)
297     CSS_bootstrap_naive[i] = temp

```

Jolliffe (2003) stratification and clustering is next.

```
298     temp = pd.DataFrame(bootstrap.bootstrap_Jol(  
299         *CSS_pos_args,  
300         "region", "household",  
301         **CSS_named_args)  
302     )["parameters"]  
303     temp = temp.explode().astype(np.float64)  
304     CSS_bootstrap_Jol[i] = temp
```

Stratification and simplified clustering based on Jolliffe (2003) is third.

```
305     temp = pd.DataFrame(bootstrap.bootstrap_Jol_sim(  
306         *CSS_pos_args,  
307         "region",  
308         **CSS_named_args)  
309     )["parameters"]  
310     temp = temp.explode().astype(np.float64)  
311     CSS_bootstrap_Jol_sim[i] = temp
```

After creating the DataFrames, we write them to files, find the normalized interquartile range, and store that in `CSS_bootstrap_se`. After the loop, we write `CSS_bootstrap_se` to a file.

```
312     for i in ["naive", "Jol", "Jol_sim"]:  
313         temp = eval("CSS_bootstrap_" + i)  
314         temp.to_csv("CSS_bootstrap_" + i + ".txt", index=False)  
315         iqrs = 0.7413 * (temp.quantile(0.75, numeric_only=True) -  
316                        temp.quantile(0.25, numeric_only=True))  
317         CSS_bootstrap_se[i] = iqrs  
318     CSS_bootstrap_se.to_csv("CSS_bootstrap_se.txt")  
319     print("se's are:")  
320     print(CSS_bootstrap_se)  
321 else:  
322     print("Skipping bootstrapping")  
323 print()
```

Now figures. We will make several:

1. `CSS_lin_densities.pdf`—CSS_InvG densities in 1967, 1995, and 2023, plotted on linear and log-log plots with linear relationship between parameters
2. `CSS_prop_densities.pdf`—CSS_InvG densities in 1967, 1995, and 2023, plotted on linear and log-log plots with proportional relationship between parameters
3. `InvG_densities.pdf`—InvG densities in 1967, 1995, and 2023, plotted on linear and log-log plots
4. `Fisk_densities.pdf`—Fisk densities in 1967, 1995, and 2023, plotted on linear and log-log plots
5. `InvG_parameter_graphs.pdf`—InvG parameter estimates in three panels and a fourth panel for normalized parameter estimates
6. `InvG_parameter_regression.pdf`—One panel with differences of normalized parameters, one panel with quotients of normalized parameters, and two panels with β and c predictions.

7. Fisk_parameters_normalized.pdf—normalized Fisk parameter estimates
8. CSS_InvG_parameters_graph.pdf—single panel; CSS_InvG α estimates
9. comparison_linear_graphs.pdf—linear graphs for GBII, Burr, Dagum, Davis, LogN-P mix/cut using 2023 data
10. comparison_loglog_graphs.pdf—loglog graphs for GBII, Burr, Dagum, Davis, LogN-P mix/cut using 2023 data
11. gini_compare_graphs.pdf

First, we load the InvG, Fisk, and CSS_InvG parameters.

```

324 #####
325 ## figures ##
326 #####
327 the_time()
328 if do_figures:
329     if not do_load_data:
330         raise RuntimeError("Please load data before making figures")
331     print("Making figures")
332     InvG_parameters = pd.read_csv("InvG_parameters.txt", header=0, index_col=0)
333     CSS_InvG_lin_parameters = \
334         pd.read_csv("CSS_InvG_lin_parameters.txt", header=0, index_col=0)
335     CSS_InvG_prop_parameters = \
336         pd.read_csv("CSS_InvG_prop_parameters.txt", header=0, index_col=0)
337     Fisk_parameters = pd.read_csv("Fisk_parameters.txt", header=0, index_col=0)

```

Again, we load the CSS_InvG constants

```

338 f = open("CSS_InvG_lin_constants.txt")
339 for line in f:
340     temp = line[:-1].split(",")
341     exec("{0}_lin = {1}".format(temp[0], temp[1]))
342 f.close()
343 f = open("CSS_InvG_prop_constants.txt")
344 for line in f:
345     temp = line[:-1].split(",")
346     exec("{0}_prop = {1}".format(temp[0], temp[1]))
347 f.close()

```

Figure 1: CSS_lin_densities.pdf. We begin with the plots of empirical and model density for past years of data. To keep things simple, we use data from 1967, 1995, and 2023. We define density functions for each year. To make it simple to change years later, we save the years in variables year1, year2, and year3.

```

348 print("Making CSS_lin_densities.pdf")
349 year1, year2, year3 = [1967, 1995, 2023]
350 for i in [year1, year2, year3]:
351     exec("""def F_{0}(x):
352         return est.density["CSS_InvG"](x, {0}, phi_lin,
353             [psi0_lin, psi1_lin, psi2_lin],
354             CSS_InvG_lin_parameters.loc[{0}, "alpha"])""".format(i))
355     exec("""c_{0} = psi0_lin + psi1_lin * {0} + \

```

```

356     psi2_lin * CSS_InvG_lin_parameters.loc[{0}, "alpha"]"".format(i))
357 make_figures.lin_loglog_graphs(
358     data[year1], data[year2], data[year3], "income",
359     [eval("F_{0}".format(year1))], [eval("c_{0}".format(year1))],
360     [{"c": "black", "lw": 0.7}],
361     [eval("F_{0}".format(year2))], [eval("c_{0}".format(year2))],
362     [{"c": "black", "lw": 0.7}],
363     [eval("F_{0}".format(year3))], [eval("c_{0}".format(year3))],
364     [{"c": "black", "lw": 0.7}],
365     [-20000,100000], [-0.05e-4,2.05e-4], [1000,1100000], [1e-10,2e-4],
366     filename="CSS_lin_densities", wgt="weight",
367     titles=["{0} Data (Linear Scale)".format(year1),
368            "{0} Data (Log Scale)".format(year1),
369            "{0} Data (Linear Scale)".format(year2),
370            "{0} Data (Log Scale)".format(year2),
371            "{0} Data (Linear Scale)".format(year3),
372            "{0} Data (Log Scale)".format(year3)])

```

Figure 2: CSS_prop_densities.pdf. Same thing except with proportional relationship imposed on parameters.

```

373 print("Making CSS_prop_densities.pdf")
374 year1, year2, year3 = [1967, 1995, 2023]
375 for i in [year1, year2, year3]:
376     exec("""def F_{0}(x):
377         return est.density["CSS_InvG_prop"](x, {0}, phi_prop,
378         [psi0_prop, psi1_prop],
379         CSS_InvG_prop_parameters.loc[{0}, "alpha"])"".format(i))
380     exec("""c_{0} = CSS_InvG_prop_parameters.loc[{0}, "alpha"] * \
381         (psi0_prop + psi1_prop * {0})"".format(i))
382 make_figures.lin_loglog_graphs(
383     data[year1], data[year2], data[year3], "income",
384     [eval("F_{0}".format(year1))], [eval("c_{0}".format(year1))],
385     [{"c": "black", "lw": 0.7}],
386     [eval("F_{0}".format(year2))], [eval("c_{0}".format(year2))],
387     [{"c": "black", "lw": 0.7}],
388     [eval("F_{0}".format(year3))], [eval("c_{0}".format(year3))],
389     [{"c": "black", "lw": 0.7}],
390     [-20000,100000], [-0.05e-4,2.05e-4], [1000,1100000], [1e-10,2e-4],
391     filename="CSS_prop_densities", wgt="weight",
392     titles=["{0} Data (Linear Scale)".format(year1),
393            "{0} Data (Log Scale)".format(year1),
394            "{0} Data (Linear Scale)".format(year2),
395            "{0} Data (Log Scale)".format(year2),
396            "{0} Data (Linear Scale)".format(year3),
397            "{0} Data (Log Scale)".format(year3)])

```

Figure 3: InvG_densities.pdf. Same thing except with inverse-gamma.

```

398 print("Making InvG_densities.pdf")
399 year1, year2, year3 = [1967, 1995, 2023]

```



```

400 for i in [year1, year2, year3]:
401     exec("""def F_{0}(x):
402         return est.density["InvG"](x, InvG_parameters.loc[{0}])""".format(i))
403     exec("""c_{0} = InvG_parameters.loc[{0}, "c"]""".format(i))
404 make_figures.lin_loglog_graphs(
405     data[year1], data[year2], data[year3], "income",
406     [eval("F_{0}".format(year1))], [eval("c_{0}".format(year1))],
407     [{"c": "black", "lw": 0.7}],
408     [eval("F_{0}".format(year2))], [eval("c_{0}".format(year2))],
409     [{"c": "black", "lw": 0.7}],
410     [eval("F_{0}".format(year3))], [eval("c_{0}".format(year3))],
411     [{"c": "black", "lw": 0.7}],
412     [-20000,100000], [-0.05e-4,2.05e-4], [1000,1100000], [1e-10,2e-4],
413     filename="InvG_densities", wgt="weight",
414     titles=["{0} Data (Linear Scale)".format(year1),
415            "{0} Data (Log Scale)".format(year1),
416            "{0} Data (Linear Scale)".format(year2),
417            "{0} Data (Log Scale)".format(year2),
418            "{0} Data (Linear Scale)".format(year3),
419            "{0} Data (Log Scale)".format(year3)])

```

Figure 4: Fisk_densities.pdf. Same thing except with Fisk densities.

```

420 print("Making Fisk_densities.pdf")
421 year1, year2, year3 = [1967, 1995, 2023]
422 for i in [year1, year2, year3]:
423     exec("""def F_{0}(x):
424         return est.density["Fisk"](x, Fisk_parameters.loc[{0}])""".format(i))
425     exec("""c_{0} = Fisk_parameters.loc[{0}, "c"]""".format(i))
426 make_figures.lin_loglog_graphs(
427     data[year1], data[year2], data[year3], "income",
428     [eval("F_{0}".format(year1))], [eval("c_{0}".format(year1))],
429     [{"c": "black", "lw": 0.7}],
430     [eval("F_{0}".format(year2))], [eval("c_{0}".format(year2))],
431     [{"c": "black", "lw": 0.7}],
432     [eval("F_{0}".format(year3))], [eval("c_{0}".format(year3))],
433     [{"c": "black", "lw": 0.7}],
434     [-20000,100000], [-0.05e-4,2.05e-4], [1000,1100000], [1e-10,2e-4],
435     filename="Fisk_densities", wgt="weight",
436     titles=["{0} Data (Linear Scale)".format(year1),
437            "{0} Data (Log Scale)".format(year1),
438            "{0} Data (Linear Scale)".format(year2),
439            "{0} Data (Log Scale)".format(year2),
440            "{0} Data (Linear Scale)".format(year3),
441            "{0} Data (Log Scale)".format(year3)])

```

Figure 5: InvG_parameter_graphs.pdf—The next set of graphs shows the graphs of inverse-gamma parameters. We have four panels: one for each parameter and one with the normalized parameters. Before we create the figure, we need a few more Series. The **norm** Series

are inverse-gamma parameters normalized by the sum of that parameter across all years.

```

442 print("Making InvG_parameter_graphs.pdf")
443 norm_alpha = InvG_parameters["alpha"] / InvG_parameters["alpha"].sum()
444 norm_beta = InvG_parameters["beta"] / InvG_parameters["beta"].sum()
445 norm_c = InvG_parameters["c"] / InvG_parameters["c"].sum()
446 make_figures.four_graphs(
447     [years], [InvG_parameters["alpha"]], [plt.plot],
448     [{"c": "black", "lw": 0.5}],
449     [years], [InvG_parameters["beta"]], [plt.plot],
450     [{"c": "black", "lw": 0.5}],
451     [years], [InvG_parameters["c"]], [plt.plot],
452     [{"c": "black", "lw": 0.5}],
453     [*[years]*3],
454     [norm_alpha, norm_beta, norm_c], [plt.plot, plt.plot, plt.plot],
455     [{"c": "blue", "lw": 0.7, "ls": "--", "label": "Shape"},
456     {"c": "black", "lw": 0.5, "label": "Scale"},
457     {"c": "red", "lw": 1.0, "ls": ":", "label": "Shift"}],
458     filename="InvG_parameter_graphs", titles=["Shape Parameter",
459     "Scale Parameter", "Shift Parameter", "Normalized Parameters"],
460     extra_code = \
461     """plt.subplot(2,2,4)
462 plt.legend()""")

```

Figure 6: `InvG_parameter_regression.pdf`—The next figure is a plot illustrating the relationship between the parameters. The upper row will show the difference and quotient between normalized shape and the other two parameters. The lower row will comparing actual and predicted values for β and c . The `beta_hat` and `c_hat` Series are the predicted values from the linear regression of β and c on α and the year.

```

463 print("Making InvG_parameter_regression.pdf")
464 beta_hat = (psi0_lin + psi1_lin * InvG_parameters.index + \
465     psi2_lin * InvG_parameters["alpha"]) / phi_lin
466 c_hat = psi0_lin + psi1_lin * InvG_parameters.index + \
467     psi2_lin * InvG_parameters["alpha"]
468 norm_beta_hat = beta_hat / beta_hat.sum()
469 norm_c_hat = c_hat / c_hat.sum()

```

We make numpy arrays out of the differences and quotients of the normalized parameters.

```

470 diff_beta_alpha = (norm_beta - norm_alpha).to_numpy()
471 diff_c_alpha = (norm_c - norm_alpha).to_numpy()
472 quot_beta_alpha = (norm_beta / norm_alpha).to_numpy()
473 quot_c_alpha = (norm_c / norm_alpha).to_numpy()

```

Now two linear regressions for the combinations of normalized parameters. We are approximating the quotient and differences as linear functions of time. First the differences.

```

474 year_cons = pd.DataFrame({"cons": 1, "year": years}).to_numpy()
475 reg_diff = np.linalg.inv(np.transpose(year_cons) @ year_cons) @ \
476     np.transpose(year_cons) @ (0.5 * (diff_beta_alpha + diff_c_alpha))
477 reg_diff_vals = reg_diff[0] + reg_diff[1] * InvG_parameters.index

```

And quotient.

```
478 reg_quot = np.linalg.inv(np.transpose(year_cons) @ year_cons) @ \
479     np.transpose(year_cons) @ (0.5 * (quot_beta_alpha + quot_c_alpha))
480 reg_quot_vals = reg_quot[0] + reg_quot[1] * InvG_parameters.index
```

Now make the figure. Every list of x -values to use will be the same, so we put copies of `years` in a list for each of those function arguments.

```
481 from matplotlib.markers import MarkerStyle
482 make_figures.four_graphs(
483     [*[years]*3],
484     [diff_beta_alpha, diff_c_alpha, reg_diff_vals],
485     [plt.scatter, plt.scatter, plt.plot],
486     [{"c": "blue", "s": 2.0, "label": "$\\bar{\\beta}_t-\\bar{\\alpha}_t$"},
487      {"c": "red", "s": 10.0, "label": "$\\bar{c}_t-\\bar{\\alpha}_t$",
488       "marker": "2"}],
489     [{"c": "black", "lw": 0.7, "ls": "--", "label": "Trendline"}],
490     [*[years]*3],
491     [quot_beta_alpha, quot_c_alpha, reg_quot_vals],
492     [plt.scatter, plt.scatter, plt.plot],
493     [{"c": "blue", "s": 2.0, "label": "$\\bar{\\beta}_t/\\bar{\\alpha}_t$"},
494      {"c": "red", "s": 10.0, "label": "$\\bar{c}_t/\\bar{\\alpha}_t$",
495       "marker": "2"}],
496     [{"c": "black", "lw": 0.7, "ls": "--", "label": "Trendline"}],
497     [*[years]*2],
498     [InvG_parameters["beta"], beta_hat], [plt.plot, plt.plot],
499     [{"c": "black", "lw": 0.5, "label": "Observed"},
500      {"c": "blue", "lw": 0.7, "ls": "--", "label": "Predicted"}],
501     [*[years]*2],
502     [InvG_parameters["c"], c_hat], [plt.plot, plt.plot],
503     [{"c": "black", "lw": 0.5, "label": "Observed"},
504      {"c": "blue", "lw": 0.7, "ls": "--", "label": "Predicted"}],
505     filename="InvG_parameter_regression",
506     titles=["(Normalized) Differences",
507            "(Normalized) Quotients", "Predicted Scale",
508            "Predicted Shift"],
```

The extra code for this function call will be a loop through subplots, and on each iteration, we specify adding a legend for that subplot.

```
509     extra_code= \
510     """for i in range(1, 5):
511         plt.subplot(2, 2, i)
512         plt.legend()"""
```

Figure 7: `Fisk_parameters_normalized.pdf`—a graph of normalized Fisk parameters. We start by creating series of normalized parameters, and then we put them in a figure.

```
513 print("Making Fisk_parameters_normalized.pdf")
514 Fisk_alpha_norm = Fisk_parameters["alpha"] / Fisk_parameters["alpha"].sum()
515 Fisk_beta_norm = Fisk_parameters["beta"] / Fisk_parameters["beta"].sum()
516 Fisk_c_norm = Fisk_parameters["c"] / Fisk_parameters["c"].sum()
```

```

517 make_figures.single_graph_ext(
518     [*[years]*3],
519     [Fisk_alpha_norm, Fisk_beta_norm, Fisk_c_norm],
520     [plt.plot, plt.plot, plt.plot],
521     [{"c": "blue", "lw": 0.7, "ls": "--", "label": "Shape"},
522      {"c": "black", "lw": 0.5, "label": "Scale"},
523      {"c": "red", "lw": 1.0, "ls": ":", "label": "Shift"}],
524     filename="Fisk_parameters_normalized", title="Normalized Fisk Parameters",
525     extra_code = "plt.legend()")

```

Figure 8: CSS_InvG_parameters_graph.pdf—A graph of the constant-shift-scale inverse-gamma parameter estimates.

```

526 print("Making CSS_InvG_parameters_graph.pdf")
527 make_figures.single_graph(CSS_InvG_lin_parameters, "alpha",
528     filename="CSS_InvG_parameters_graph", title="Parameter Estimates")

```

Figure 9: comparison_linear_graphs.pdf—Now we make graphs for the different densities in 2023. First we load the files with parameter estimates saved as dictionaries. The dictionary parameters will save the parameters for each distribution.

```

529 print("comparison_linear_graphs.pdf")
530 parameters = {}
531 def add_from_parameter_file(filename, parameter_dict):
532     f = open(filename)
533     curr_model = ""
534     for line in f:
535         if line[0] == "-":
536             curr_model = line[3:-4]
537         elif line[0] == "f":
538             pass
539         elif line[0] == "p":
540             temp = line[13:-2].split(", ")
541             parameter_dict[curr_model] = list(map(float, temp))
542     f.close()
543     return parameter_dict
544 parameters = add_from_parameter_file("2023_parameters_short.txt", parameters)
545 parameters = add_from_parameter_file("2023_parameters_long.txt", parameters)

```

Now we create the density functions and constants.

```

546 for i in ["GB2", "Dagum", "Burr", "Davis", "LogN_P_cut",
547          "LogN_P_mix"]:
548     exec("""def F_{0}(x):
549         return est.density['{0}'](x, parameters['{0}'])""".format(i))
550     exec("c_{0} = parameters['{0}'][-1]".format(i))

```

Constant-shift-scale inverse-gamma density.

```

551 def F_CSS_InvG(x):
552     return est.density["CSS_InvG"](x, 2023, phi_lin,
553         [psi0_lin, psi1_lin, psi2_lin],
554         CSS_InvG_lin_parameters.loc[2023, "alpha"])
555 c_CSS_InvG = psi0_lin + psi1_lin * 2023 + \

```

```
556     psi2_lin * CSS_InvG_lin_parameters.loc[2023, "alpha"]
```

Make the graphs. We create two sets of comparison graphs where we put the inverse-gamma density on each plot and one alternative model on each plot. The first set of graphs will be linear scale, and the second one will be loglog scale.

```
557     make_figures.lin_graphs(*[data[2023]]*6, "income",
558         [F_CSS_InvG, F_GB2], [c_CSS_InvG, c_GB2],
559         [{"c": "black", "lw": 0.5, "label": "Inverse-Gamma"},
560          {"c": "blue", "lw": 0.7, "ls": "--", "label": "Gen Beta II"}],
561         [F_CSS_InvG, F_Dagum], [c_CSS_InvG, c_Dagum],
562         [{"c": "black", "lw": 0.5, "label": "Inverse-Gamma"},
563          {"c": "blue", "lw": 0.7, "ls": "--", "label": "Dagum"}],
564         [F_CSS_InvG, F_Burr], [c_CSS_InvG, c_Burr],
565         [{"c": "black", "lw": 0.5, "label": "Inverse-Gamma"},
566          {"c": "blue", "lw": 0.7, "ls": "--", "label": "Burr"}],
567         [F_CSS_InvG, F_Davis], [c_CSS_InvG, c_Davis],
568         [{"c": "black", "lw": 0.5, "label": "Inverse-Gamma"},
569          {"c": "blue", "lw": 0.7, "ls": "--", "label": "Davis"}],
570         [F_CSS_InvG, F_LogN_P_cut], [c_CSS_InvG, c_LogN_P_cut],
571         [{"c": "black", "lw": 0.5, "label": "Inverse-Gamma"},
572          {"c": "blue", "lw": 0.7, "ls": "--", "label": "Cutoff"}],
573         [F_CSS_InvG, F_LogN_P_mix], [c_CSS_InvG, c_LogN_P_mix],
574         [{"c": "black", "lw": 0.5, "label": "Inverse-Gamma"},
575          {"c": "blue", "lw": 0.7, "ls": "--", "label": "Mixture"}],
576         [-20000, 100000], [-0.05e-5, 2.05e-5],
577         filename="comparison_linear_graphs", wgt="weight",
578         titles=["Gen Beta II", "Dagum", "Burr", "Davis",
579                "Log-Normal/Pareto Cutoff", "Log-Normal/Pareto Mix"])
```

Figure 10: comparison_loglog_graphs.pdf—Same thing with loglog scaling.

```
580     print("Making comparison_loglog_graphs.pdf")
581     make_figures.loglog_graphs(*[data[2023]]*6, "income",
582         [F_CSS_InvG, F_GB2],
583         [{"c": "black", "lw": 0.5, "label": "Inverse-Gamma"},
584          {"c": "blue", "lw": 0.7, "ls": "--", "label": "Gen Beta II"}],
585         [F_CSS_InvG, F_Dagum],
586         [{"c": "black", "lw": 0.5, "label": "Inverse-Gamma"},
587          {"c": "blue", "lw": 0.7, "ls": "--", "label": "Dagum"}],
588         [F_CSS_InvG, F_Burr],
589         [{"c": "black", "lw": 0.5, "label": "Inverse-Gamma"},
590          {"c": "blue", "lw": 0.7, "ls": "--", "label": "Burr"}],
591         [F_CSS_InvG, F_Davis],
592         [{"c": "black", "lw": 0.5, "label": "Inverse-Gamma"},
593          {"c": "blue", "lw": 0.7, "ls": "--", "label": "Davis"}],
594         [F_CSS_InvG, F_LogN_P_cut],
595         [{"c": "black", "lw": 0.5, "label": "Inverse-Gamma"},
596          {"c": "blue", "lw": 0.7, "ls": "--", "label": "Cutoff"}],
597         [F_CSS_InvG, F_LogN_P_mix],
598         [{"c": "black", "lw": 0.5, "label": "Inverse-Gamma"},
```

```

599         {"c": "blue", "lw": 0.7, "ls": "--", "label": "Mixture"}],
600     [1000, 1100000], [1e-10, 2e-4],
601     filename="comparison_loglog_graphs", wgt="weight",
602     titles=["Gen Beta II", "Dagum", "Burr", "Davis",
603         "Log-Normal/Pareto Cutoff", "Log-Normal/Pareto Mix"])

```

Figure 11: `gini_compare_graphs.pdf`—For fun we make a final graphic about the Gini coefficient, which is

$$\frac{\beta}{\beta + (\alpha - 1)c} \frac{\Gamma(\alpha - \frac{1}{2})}{\sqrt{\pi}\Gamma(\alpha)} = \frac{1}{1 + (\alpha - 1)\phi} \frac{\Gamma(\alpha - \frac{1}{2})}{\sqrt{\pi}\Gamma(\alpha)}$$

We will use `four_graphs()` again. The `gini` Series contains values of the Gini coefficient in each year. The `gini_unshift` is the value of the shape parameter in the unshifted Gini coefficient formula. The third and fourth graphs show the Gini coefficient for different values of α .

```

604     print("Making gini_compare_graphs.pdf")
605     gini = 1 / (1 + (CSS_InvG_lin_parameters["alpha"] - 1) * phi_lin) * \
606         est.G(CSS_InvG_lin_parameters["alpha"] - 0.5) / (np.sqrt(np.pi) *
607             est.G(CSS_InvG_lin_parameters["alpha"]))
608     gini_unshift = \
609         est.G(CSS_InvG_lin_parameters["alpha"] - 0.5) / (np.sqrt(np.pi) *
610             est.G(CSS_InvG_lin_parameters["alpha"]))
611     make_figures.single_graph(pd.DataFrame({"gini": gini}), "gini",
612         filename="gini_graph")
613     make_figures.single_graph(pd.DataFrame({"gini": gini_unshift}), "gini",
614         filename="gini_unshift_graph")

```

Now make the Gini coefficient for different values of α . The graph with Gini coefficient of unshifted inverse-gamma distribution is straightforward. For shifted inverse-gamma distribution, we have a minimum value around 3. Setting the derivative equal to 0 gives us

$$y = \frac{1}{1 + (x - 1)\phi} \frac{\Gamma(x - \frac{1}{2})}{\sqrt{\pi}\Gamma(x)}$$

$$\log y = -\log(1 + (x - 1)\phi) + \log \Gamma\left(x - \frac{1}{2}\right) - \log \Gamma(x) - \log \sqrt{\pi}$$

$$\frac{y'}{y} = -\frac{\phi}{1 + (x - 1)\phi} + \psi\left(x - \frac{1}{2}\right) - \psi(x) = 0$$

We can solve this equation numerically for α . The `right_singularity` variable is the x -value where the Gini coefficient has a singularity because the denominator blows up, in other words $1 + (\alpha - 1)\phi \rightarrow 0$. The `alpha_vals`, `gini_vals_shift`, and `gini_vals_unshift` lists are lists of x and y -values for plotting the Gini coefficient.

```

615     right_singularity = 1 - 1 / phi_lin
616     alpha_vals = np.linspace(0.6, right_singularity - 0.2, 200)
617     gini_vals_shift = 1 / (1 + (alpha_vals - 1) * phi_lin) * \
618         est.G(alpha_vals - 0.5) / (np.sqrt(np.pi) *
619             est.G(alpha_vals))

```

```

620 gini_vals_unshift = \
621     est.G(alpha_vals - 0.5) / (np.sqrt(np.pi) *
622     est.G(alpha_vals))

```

The `D_gini` function is the log-derivative y'/y of the Gini coefficient for shifted inverse-gamma distribution. As noted above, we can find the α corresponding to minimum Gini coefficient when $D_gini = 0$. We set `min_alpha` and `max_alpha` to be the minimum and maximum values of α in the data.

```

623 def D_gini(x):
624     return est.psi(x - 0.5) - est.psi(x) - phi_lin / (1 + (x - 1) * phi_lin)
625 min_gini = est.root(D_gini, bracket=[1, right_singularity - 0.2]).root
626 min_alpha = CSS_InvG_lin_parameters["alpha"].min()
627 max_alpha = CSS_InvG_lin_parameters["alpha"].max()

```

Now make the figure. We have several graphs on each subplot:

1. Shape parameter and horizontal lines for minimum and maximum values
2. Gini coefficient calculated from parameter estimates
3. Gini coefficient as a function of α under a shifted inverse-gamma distribution as well as vertical lines at the minimum and maximum values of α observed in the data
4. Gini coefficient as a function of α under an unshifted inverse-gamma distribution as well as vertical lines at the minimum and maximum values of α observed in the data

The file will be `gini_compare_graphs.pdf`.

```

628 make_figures.four_graphs(
629     [years, (years[0], years[-1]), (years[0], years[-1])],
630     [CSS_InvG_lin_parameters["alpha"],
631      (min_alpha, min_alpha), (max_alpha, max_alpha)],
632     [plt.plot, plt.plot, plt.plot],
633     [{"c": "black", "lw": 0.5},
634      {"c": "red", "lw": 0.7, "ls": "--"},
635      {"c": "blue", "lw": 0.7, "ls": "--"}],
636     [years], [gini], [plt.plot],
637     [{"c": "black", "lw": 0.5}],
638     [alpha_vals, (min_alpha, min_alpha), (max_alpha, max_alpha)],
639     [gini_vals_shift, (0,2), (0,2)], [plt.plot, plt.plot, plt.plot],
640     [{"c": "black", "lw": 0.5},
641      {"c": "red", "lw": 0.7, "ls": "--"},
642      {"c": "blue", "lw": 0.7, "ls": "--"}],
643     [alpha_vals, (min_alpha, min_alpha), (max_alpha, max_alpha)],
644     [gini_vals_unshift, (0,2), (0,2)], [plt.plot, plt.plot, plt.plot],
645     [{"c": "black", "lw": 0.5},
646      {"c": "red", "lw": 0.7, "ls": "--"},
647      {"c": "blue", "lw": 0.7, "ls": "--"}],
648     filename="gini_compare_graphs",
649     titles=["Shape Parameter",
650            "Gini Coefficient from Parameters",
651            "Gini Coefficient Function (With Shift)",
652            "Gini Coefficient Function (No Shift)"],

```

The `extra_code` for this function sets the vertical limits for the third and fourth subplots and adds a light blue rectangular patch to the third graph to denote the portion of the domain where the Gini coefficient is increasing.

```
653     extra_code= \
654     """plt.subplot(2, 2, 3)
655     plt.ylim([0,2])
656     plt.gca().add_patch(mpl.patches.Rectangle([min_gini,0],
657         {max_alpha} - {min_gini}, 2, color="aliceblue"))
658     plt.subplot(2, 2, 4)
659     plt.ylim([0,2])""".format(min_gini=min_gini, max_alpha=alpha_vals[-1]))
    Done with figures.

660 else:
661     print("Skipping figures")
662     print()
    Random code to test things.

663 ## Code for testing
664 if do_test:
665     the_time()
666     print("Testing...")
667
668     # This section is intentionally blank
669
670     print()
    End of code!

671 the_time()
672 print("End of main.py")
673 print()
```