

In [1]:

```
1 # =====
2 # Importation of Libraries
3 #
4 # Michael Collins, Flatiron School
5 # 2020-10-30_0639_MDT
6 #
7 # =====
8 import math
9 from IPython.display import display
10 from openlocationcode import openlocationcode as olc
11 from geographiclib.geodesic import Geodesic
12 import numpy as np
13 import pandas as pd
14 from datetime import datetime
15 import itertools
16 import warnings
17 warnings.filterwarnings('ignore')
18 import scipy.stats as scs
19 from scipy.stats import median_abs_deviation as MAD
20 import statsmodels.api as sm
21 from sklearn.datasets import make_regression
22 import matplotlib.pyplot as plt
23 import seaborn as sns
24 import flatiron.statsFunctions as sf
25 from importlib import reload
26 reload(sf)
27 import plotly.express as px
28 import requests
29
```

In [2]:

```
1 # =====
2 # K-Nearest-Neighbor Evaluation
3 #
4 # Suppose we are given:
5 #   * a number of events "N"
6 #   * a list of event indices "event_i"
7 #   * a list of event coordinates "event_latlng"
8 #   * some number of events "K" < N
9 #   * the coordinates "Q_Latlng" of some arbitrary
10 #      point Q on earth's surface
11 #
12 # The code in this cell returns [the list of indices for
13 # the K events whose proximity is nearest to point Q]
14 # for each of the N events.
15 #
16 # Michael Collins, Flatiron School
17 # 2020-10-28_0524_MDT
18 #
19 # =====
20
21 OLC_CODE_PLANET = "00000000+"
22 OLC_EARTH_CIRCUMFERENCE = np.round(2 * np.pi * Geodesic.WGS84.a, 3)
23 OLC_LATITUDE_SPAN = 180.0
24 OLC_LATITUDE_SOUTH_POLE = -90.0
25 OLC_LATITUDE_EQATOR = 0.0
26 OLC_LATITUDE_NORTH_POLE = 90.0
27 OLC_LONGITUDE_SPAN = 360.0
28 OLC_LONGITUDE_DATELINE_WEST = -180.0
29 OLC_LONGITUDE_GREENWICH = 0.0
30 OLC_LONGITUDE_DATELINE_EAST = 180.0
31
32 OLC_PYLON_MIDPOINT = "midpoint"
33 OLC_PYLON_O = OLC_PYLON_MIDPOINT
34 OLC_PYLON_N = "pylon_N"
35 OLC_PYLON_NE = "pylon_NE"
36 OLC_PYLON_E = "pylon_E"
37 OLC_PYLON_SE = "pylon_SE"
38 OLC_PYLON_S = "pylon_S"
39 OLC_PYLON_SW = "pylon_SW"
40 OLC_PYLON_W = "pylon_W"
41 OLC_PYLON_NW = "pylon_NW"
42 OLC_BORDER_PYLONS = [OLC_PYLON_N, OLC_PYLON_NE, OLC_PYLON_E, OLC_PYLON_SE,
43                      OLC_PYLON_S, OLC_PYLON_SW, OLC_PYLON_W, OLC_PYLON_NW]
44 OLC_PYLONS = [OLC_PYLON_O] + OLC_BORDER_PYLONS
45
46 OLC_OFFSET_N = "offset_N"
47 OLC_OFFSET_NE = "offset_NE"
48 OLC_OFFSET_E = "offset_E"
49 OLC_OFFSET_SE = "offset_SE"
50 OLC_OFFSET_S = "offset_S"
51 OLC_OFFSET_SW = "offset_SW"
52 OLC_OFFSET_W = "offset_W"
53 OLC_OFFSET_NW = "offset_NW"
54 OLC_CARDINAL_OFFSETS = [OLC_OFFSET_N, OLC_OFFSET_E, OLC_OFFSET_S, OLC_OFFSET_W]
55 OLC_DIAGONAL_OFFSETS = [OLC_OFFSET_NE, OLC_OFFSET_SE, OLC_OFFSET_SW, OLC_OFFSET_NW]
56 OLC_OFFSETS = [OLC_OFFSET_N, OLC_OFFSET_NE, OLC_OFFSET_E, OLC_OFFSET_SE,
57                 OLC_OFFSET_S, OLC_OFFSET_SW, OLC_OFFSET_W, OLC_OFFSET_NW]
58
59 OLC_TAU_PLANET = 0    # (180)-degree by (360)-degree resolution
60 OLC_TAU_DOMAIN = 1    # (20)-degree by (20)-degree resolution
61 OLC_TAU_SHIRE = 2     # (1)-degree by (1)-degree resolution
62 OLC_TAU_TOWNSHIP = 3  # (1/20)-degree by (1/20)-degree resolution
63 OLC_TAU_HOLLER = 4    # (1/400)-degree by (1/400)-degree resolution
64 OLC_TAU_ROOFTOP = 5   # (1/8000)-degree by (1/8000)-degree resolution
65 OLC_TAU_PARLOR = 6    # (1/4000)-degree by (1/32000)-degree resolution
66 OLC_TAU_DESK = 7       # subdivides the above into a 5-tall by 4-wide grid
67 OLC_TAU_ENVELOPE = 8   # subdivides the above into a 5-tall by 4-wide grid
68 OLC_TAU_ERASER = 9    # subdivides the above into a 5-tall by 4-wide grid
69 OLC_TAU_STAPLE = 10   # subdivides the above into a 5-tall by 4-wide grid
70 OLC_TAUS = [OLC_TAU_PLANET, OLC_TAU_DOMAIN, OLC_TAU_SHIRE,
71             OLC_TAU_TOWNSHIP, OLC_TAU_HOLLER, OLC_TAU_ROOFTOP,
```

```

72         OLC_TAU_PARLOR, OLC_TAU_DESK, OLC_TAU_ENVELOPE,
73         OLC_TAU_ERASER, OLC_TAU_STAPLE]
74
75 OLC_TAU_BEDROCK = OLC_TAU_ROOFTOP
76 OLC_GRID_NORTH_MAX = np.prod([9, 20, 20, 20, 20][0:OLC_TAU_BEDROCK])
77 OLC_GRID_EAST_MAX = np.prod([18, 20, 20, 20][0:OLC_TAU_BEDROCK])
78
79
80 OLC_TIER_PLANET = "planet"      # (180)-degree by (360)-degree resolution
81 OLC_TIER_DOMAIN = "domain"     # (20)-degree by (20)-degree resolution
82 OLC_TIER_SHIRE = "shire"       # (1)-degree by (1)-degree resolution
83 OLC_TIER_TOWNSHIP = "township" # (1/20)-degree by (1/20)-degree resolution
84 OLC_TIER_HOLLER = "holler"     # (1/400)-degree by (1/400)-degree resolution
85 OLC_TIER_ROOFTOP = "rooftop"   # (1/8000)-degree by (1/8000)-degree resolution
86 OLC_TIER_PARLOR = "parlor"     # (1/40000)-degree by (1/32000)-degree resolution
87 OLC_TIER_DESK = "desk"         # subdivides the above into a 5-tall by 4-wide grid
88 OLC_TIER_ENVELOPE = "envelope" # subdivides the above into a 5-tall by 4-wide grid
89 OLC_TIER_ERASER = "eraser"     # subdivides the above into a 5-tall by 4-wide grid
90 OLC_TIER_STAPLE = "staple"     # subdivides the above into a 5-tall by 4-wide grid
91 OLC_TIERS = [OLC_TIER_PLANET, OLC_TIER_DOMAIN, OLC_TIER_SHIRE,
92             OLC_TIER_TOWNSHIP, OLC_TIER_HOLLER, OLC_TIER_ROOFTOP,
93             OLC_TIER_PARLOR, OLC_TIER_DESK, OLC_TIER_ENVELOPE,
94             OLC_TIER_ERASER, OLC_TIER_STAPLE]
95
96 OLC_CODELEN_PLANET = 0          # (180)-degree by (360)-degree resolution
97 OLC_CODELEN_DOMAIN = 2          # (20)-degree by (20)-degree resolution
98 OLC_CODELEN_SHIRE = 4           # (1)-degree by (1)-degree resolution
99 OLC_CODELEN_TOWNSHIP = 6         # (1/20)-degree by (1/20)-degree resolution
100 OLC_CODELEN_HOLLER = 8          # (1/400)-degree by (1/400)-degree resolution
101 OLC_CODELEN_ROOFTOP = 10        # (1/8000)-degree by (1/8000)-degree resolution
102 OLC_CODELEN_PARLOR = 11         # (1/40000)-degree by (1/32000)-degree resolution
103 OLC_CODELEN_DESK = 12           # subdivides the above into a 5-tall by 4-wide grid
104 OLC_CODELEN_ENVELOPE = 13        # subdivides the above into a 5-tall by 4-wide grid
105 OLC_CODELEN_ERASER = 14          # subdivides the above into a 5-tall by 4-wide grid
106 OLC_CODELEN_STAPLE = 15          # subdivides the above into a 5-tall by 4-wide grid
107 OLC_CODELENS = [OLC_CODELEN_PLANET, OLC_CODELEN_DOMAIN, OLC_CODELEN_SHIRE,
108                 OLC_CODELEN_TOWNSHIP, OLC_CODELEN_HOLLER, OLC_CODELEN_ROOFTOP,
109                 OLC_CODELEN_PARLOR, OLC_CODELEN_DESK, OLC_CODELEN_ENVELOPE,
110                 OLC_CODELEN_ERASER, OLC_CODELEN_STAPLE]
111
112 OLC_tau_to_tier = dict(zip(OLC_TAUS, OLC_TIERS))
113 OLC_tier_to_tau = dict(zip(OLC_TIERS, OLC_TAUS))
114 OLC_TIER_BEDROCK = OLC_tau_to_tier[OLC_TAU_BEDROCK]
115
116 OLC_tau_to_codeLen = dict(zip(OLC_TAUS, OLC_CODELENS))
117
118 OLC_codeLen_to_tau = dict(zip(OLC_CODELENS, OLC_TAUS))
119
120 OLC_tier_to_codeLen = dict(zip(OLC_TIERS, OLC_CODELENS))
121
122 OLC_codeLen_to_tier = dict(zip(OLC_CODELENS, OLC_TIERS))
123
124 def sci_round(x, n):
125     s = np.format_float_scientific(x, precision=n, trim='-')
126     return float(s)
127
128 def olc_encode_tau(latlng, tau):
129     if not tau in OLC_TAUS:
130         return None
131
132     if tau == OLC_TAU_PLANET:
133         code = OLC_CODE_PLANET
134     else:
135         code = olc.encode(*latlng, OLC_tau_to_codeLen[tau])
136
137
138 def olc_add_zone_to_archive(this_zone):
139     raw_dict = {}
140     raw_dict[this_zone.code] = this_zone
141     for child_zone in this_zone.children.values():
142         child_dict = olc_add_zone_to_archive(child_zone)
143         raw_dict.update(child_dict)

```

```

144     return dict(sorted(raw_dict.items()))
145
146 def olc_latitude_to_iNorth_fraCN(latitude=None):
147     if not isinstance(latitude, float):
148         print(f"the latitude must be specified.")
149     return None
150     if not(OLC_LATITUDE_SOUTH_POLE <= latitude <= OLC_LATITUDE_NORTH_POLE):
151         print(f"the latitude {latitude} is not valid.")
152     return None
153     if latitude == OLC_LATITUDE_SOUTH_POLE:
154         return (0, 0.0)
155     if latitude == OLC_LATITUDE_NORTH_POLE:
156         return (OLC_GRID_NORTH_MAX - 1, 1.0)
157     uNorth = sci_roun((latitude - OLC_LATITUDE_SOUTH_POLE)*OLC_GRID_NORTH_MAX/OLC_LATITUDE_SPAN, 10)
158     fN, iN = np.modf(uNorth)
159     iNorth = int(iN)
160     fraCN = float(fN)
161     return (iNorth, fraCN)
162
163 def olc_longitude_to_iEast_fraCE(longitude=None):
164     if not isinstance(longitude, float):
165         print(f"the longitude must be specified.")
166     return None
167     if not(OLC_LONGITUDE_DATELINE_WEST <= longitude <= OLC_LONGITUDE_DATELINE_EAST):
168         print(f"the longitude {longitude} is not valid.")
169     return None
170     if longitude == OLC_LONGITUDE_DATELINE_WEST:
171         return (0, 0.0)
172     if longitude == OLC_LONGITUDE_DATELINE_EAST:
173         return (OLC_GRID_EAST_MAX - 1, 1.0)
174     uEast = sci_roun((longitude - OLC_LONGITUDE_DATELINE_WEST)*OLC_GRID_EAST_MAX/OLC_LONGITUDE_SPAN, 10)
175     fE, iE = np.modf(uEast)
176     iEast = int(iE)
177     fraCE = float(fE)
178     return (iEast, fraCE)
179
180 class OLC_Offset:
181     def __init__(self, lat1=None, lng1=None, lat2=None, lng2=None, mR_est=None, mN_est=None, mE_est=None):
182         gripes = []
183
184         has_est = True
185         if mR_est is None:
186             has_est = False
187         if mN_est is None:
188             has_est = False
189         if mE_est is None:
190             has_est = False
191         if not has_est:
192             mR_est, mN_est, mE_est = None, None, None
193
194         if len(gripes) == 0:
195             if has_est:
196                 if not isinstance(mR_est, float):
197                     gripes.append(f"mR_est={mR_est} must be a floating-point number.")
198                 if not isinstance(mN_est, float):
199                     gripes.append(f"mN_est={mN_est} must be a floating-point number.")
200                 if not isinstance(mE_est, float):
201                     gripes.append(f"mE_est={mE_est} must be a floating-point number.")
202
203         if len(gripes) == 0:
204             if not isinstance(lat1, float):
205                 gripes.append(f"lat1={lat1} must be a floating point number.")
206             elif not(OLC_LATITUDE_SOUTH_POLE <= lat1 <= OLC_LATITUDE_NORTH_POLE):
207                 gripes.append(f"lat1={lat1} is not valid.")
208
209         if len(gripes) == 0:
210             if not isinstance(lat2, float):
211                 gripes.append(f"lat2={lat2} must be a floating point number.")
212             elif not(OLC_LATITUDE_SOUTH_POLE <= lat2 <= OLC_LATITUDE_NORTH_POLE):
213                 gripes.append(f"lat2={lat2} is not valid.")
214
215         if len(gripes) == 0:

```

```

216     if not isinstance(lng1, float):
217         gripes.append(f"lng1={lng1} must be a floating point number.")
218     elif not(OLC_LONGITUDE_DATELINE_WEST <= lng1 <= OLC_LONGITUDE_DATELINE_EAST):
219         gripes.append(f"lng1={lng1} is not valid.")
220
221     if len(gripes) == 0:
222         if not isinstance(lng2, float):
223             gripes.append(f"lng2={lng2} must be a floating point number.")
224         elif not(OLC_LONGITUDE_DATELINE_WEST <= lng2 <= OLC_LONGITUDE_DATELINE_EAST):
225             gripes.append(f"lng2={lng2} is not valid.")
226
227     if len(gripes) == 0:
228         try:
229             gdict = Geodesic.WGS84.Inverse(lat1, lng1, lat2, lng2) # use the WGS84 ellipsoid
230
231             # range in meters from latlng1 to latlng2
232             #   (rounded to nearest millimeter)
233             range_meters_raw = gdict['s12']
234
235             # components of forward azimuth at latlng1
236             azi12_rad = np.radians(gdict['azi1'])
237             azi12_n = sci_round(np.cos(azi12_rad), 10)
238             azi12_e = sci_round(np.sin(azi12_rad), 10)
239         except:
240             gripes.append(f"failure during Geodesic.WGS84.Inverse({lat1}, {lng1}, {lat2}, {lng2}).")
241
242     if len(gripes) == 0:
243         mR = np.round(range_meters_raw, 3)
244         mN = np.round(range_meters_raw * azi12_n, 3) # this is admittedly an estimate
245         mE = np.round(range_meters_raw * azi12_e, 3) # this is admittedly an estimate
246
247         if has_est:
248             delta_mR = np.round(mR_est - mR, 3)
249             delta_mN = np.round(mN_est - mN, 3)
250             delta_mE = np.round(mE_est - mE, 3)
251
252         else:
253             delta_mR = None
254             delta_mN = None
255             delta_mE = None
256
257         self.valid = True
258         self.latlng1 = (lat1, lng1)
259         self.latlng2 = (lat2, lng2)
260         self.mR = mR
261         self.delta_mR = delta_mR
262         self.mN = mN
263         self.delta_mN = delta_mN
264         self.mE = mE
265         self.delta_mE = delta_mE
266         self.azi12_n = azi12_n
267         self.azi12_e = azi12_e
268
269     else:
270         self.valid = False
271         self.latlng1 = (lat1, lng1)
272         self.latlng2 = (lat2, lng2)
273         self.gripes = gripes
274
275     def __str__(self):
276         return f"OLC_Offset({repr(self.latlng1)},{repr(self.latlng2)})|"
277
278     def __repr__(self):
279         s = []
280         for k in self.__dict__:
281             if k in ['mR', 'delta_mR', 'mN', 'delta_mN', 'mE', 'delta_mE']:
282                 s.append(f".{k}={repr(self.__dict__[k])}")
283
284         return '[' + ', '.join(s) + ']'
285
286 class OLC_Point:
287     def __init__(self, latitude=None, longitude=None):
288         gripes = []
289
290         if len(gripes) == 0:
291             if not isinstance(latitude, float):

```

```

288         gripes.append(f"the latitude must be specified.")
289     elif not(OLC_LATITUDE_SOUTH_POLE <= latitude <= OLC_LATITUDE_NORTH_POLE):
290         gripes.append(f"the latitude {latitude} is not valid.")
291
292     if len(gripes) == 0:
293         if not isinstance(longitude, float):
294             gripes.append(f"the longitude must be specified.")
295         elif not(OLC_LONGITUDE_DATELINE_WEST <= longitude <= OLC_LONGITUDE_DATELINE_EAST):
296             gripes.append(f"the longitude {longitude} is not valid.")
297
298     if len(gripes) == 0:
299         iN_fn = olc_latitude_to_iNorth_fracN(latitude)
300         if iN_fn is None:
301             gripes.append(f"failure during olc_latitude_to_iNorth_fracN({latitude}).")
302         else:
303             iNorth, fracN = iN_fn
304
305     if len(gripes) == 0:
306         iE_fE = olc_longitude_to_iEast_fracE(longitude)
307         if iE_fE is None:
308             gripes.append(f"failure during olc_longitude_to_iEast_fracE({longitude}).")
309         else:
310             iEast, fracE = iE_fE
311
312     if len(gripes) == 0:
313         code_bedrock = olc_encode_tau((latitude, longitude), OLC_TAU_BEDROCK)
314         if code_bedrock is None:
315             gripes.append(f"failure during olc_encode_tau(({latitude}, {longitude}), {OLC_TAU_BEDROCK}).")
316
317     if len(gripes) == 0:
318         self.valid = True
319         self.lat = latitude
320         self.lng = longitude
321         self.latlng = (latitude, longitude)
322         self.uNorth = iNorth + fracN
323         self.iNorth = iNorth
324         self.fracN = fracN
325         self.uEast = iEast + fracE
326         self.iEast = iEast
327         self.fracE = fracE
328         self.code_bedrock = code_bedrock
329         self.has_codes = False
330     else:
331         self.valid = False
332         self.latlng = (latitude, longitude)
333         self.has_codes = False
334         self.gripes = gripes
335
336     def code(self, tau=None, tier=None):
337         gripes = []
338
339         this_code = None
340
341         if tier == 'bedrock':
342             # process tier-specific request for code_bedrock
343             return self.code_bedrock
344
345         if tau == OLC_TAU_BEDROCK:
346             # process tau-specific request for code_bedrock
347             return self.code_bedrock
348
349         if not self.has_codes:
350             # generate this zone's encoding-related information
351             valid_taus = [t for t in OLC_TAUS if t <= OLC_TAU_BEDROCK]
352             valid_tiers = [OLC_tau_to_tier[t] for t in valid_taus]
353             valid_codes = []
354
355             for t in valid_taus:
356                 try:
357                     code_t = olc_encode_tau(self.latlng, t)
358                     valid_codes.append(code_t)
359                 except:
360                     gripes.append(f"failure during olc_encode_tau({self.latlng}, {t}).")

```

```

360
361     if len(gripes) == 0:
362         tau_to_code = dict(zip(valid_taus, valid_codes))
363         tier_to_code = dict(zip(valid_tiers, valid_codes))
364
365         # store this zone's encoding-related information
366         #   for possible future use.
367         self.has_codes = True
368         self.valid_taus = valid_taus
369         self.valid_tiers = valid_tiers
370         self.valid_codes = valid_codes
371         self.tau_to_code = tau_to_code
372         self.tier_to_code = tier_to_code
373
374     else:
375         valid_taus = self.valid_taus
376         valid_tiers = self.valid_tiers
377         valid_codes = self.valid_codes
378         tau_to_code = self.tau_to_code
379         tier_to_code = self.tier_to_code
380
381     if len(gripes) == 0:
382         if tau in valid_taus:
383             code_tau = tau_to_code[tau]
384             if tier in valid_tiers:
385                 if code_tau == tier_to_code[tier]:
386                     # there is no conflict between tau and tier, so use tau
387                     this_code = code_tau
388                 else:
389                     gripes.append(f"parameters tau={tau} and tier={tier} have conflicting values.")
390             else:
391                 this_code = code_tau
392         else:
393             if tier in valid_tiers:
394                 this_code = tier_to_code[tier]
395             else:
396                 gripes.append(f"neither tau={tau} nor tier={tier} is valid.")
397
398     if len(gripes) > 0:
399         self.gripes = gripes
400
401     return this_code
402
403     def __str__(self):
404         return f"OLC_Point{repr(self.latlng)}|"
405
406     def __repr__(self):
407         s = []
408         for k in self.__dict__:
409             s.append(f".{k}={repr(self.__dict__[k])}")
410         return '|'.join(s) + '|'
411
412 class OLC_Zone:
413     def __init__(self, code=None, parent=None):
414         gripes = []
415
416         if (code is None) or (code == OLC_CODE_PLANET):
417             midpoint = (OLC_LATITUDE_EQUATOR, OLC_LONGITUDE_GREENWICH)
418             origin = (OLC_LATITUDE_SOUTH_POLE, OLC_LONGITUDE_DATELINE_WEST)
419             zone_dict = {}
420             zone_dict['parent'] = None
421             zone_dict['tau'] = OLC_TAU_PLANET
422             zone_dict['tier'] = OLC_TIER_PLANET
423             zone_dict['code'] = OLC_CODE_PLANET
424             zone_dict['origin'] = OLC_Point(*origin)
425             zone_dict['midpoint'] = OLC_Point(*midpoint)
426             zone_dict['offset_from_parent'] = None
427             zone_dict[OLC_OFFSET_N] = OLC_Offset(*midpoint, OLC_LATITUDE_NORTH_POLE, OLC_LONGITUDE_GREENWICH)
428             zone_dict[OLC_OFFSET_NE] = OLC_Offset(*midpoint, OLC_LATITUDE_NORTH_POLE, OLC_LONGITUDE_DATELINE_EAST)
429             zone_dict[OLC_OFFSET_E] = OLC_Offset(*midpoint, OLC_LATITUDE_EQUATOR, OLC_LONGITUDE_DATELINE_EAST)
430             zone_dict[OLC_OFFSET_SE] = OLC_Offset(*midpoint, OLC_LATITUDE_SOUTH_POLE, OLC_LONGITUDE_DATELINE_EAST)
431             zone_dict[OLC_OFFSET_S] = OLC_Offset(*midpoint, OLC_LATITUDE_SOUTH_POLE, OLC_LONGITUDE_GREENWICH)
432             zone_dict[OLC_OFFSET_SW] = OLC_Offset(*midpoint, OLC_LATITUDE_SOUTH_POLE, OLC_LONGITUDE_DATELINE_WEST)
433             zone_dict[OLC_OFFSET_W] = OLC_Offset(*midpoint, OLC_LATITUDE_EQUATOR, OLC_LONGITUDE_DATELINE_WEST)

```

```

432     zone_dict[OLC_OFFSET_NW] = OLC_Offset(*midpoint, OLC_LATITUDE_NORTH_POLE, OLC_LONGITUDE_DATELINE_WEST)
433     this_offsets = [zone_dict[offset] for offset in OLC_OFFSETs]
434 else:
435     if isinstance(parent, type(self)):
436         this_parent = parent
437     else:
438         gripes.append(f"\"a parent zone was not specified.\"")
439
440     if len(gripes) == 0:
441         try:
442             codeArea = olc.decode(code)
443             this_code = code
444         except:
445             gripes.append(f"\"failed to decode OLC code {repr(code)}\"")
446
447     if len(gripes) == 0:
448         this_codeLen = codeArea.codeLength
449         if not this_codeLen in OLC_CODELENS[1:]:
450             gripes.append(f"\"the derived code length {this_codeLen} of OLC code {repr(this_code)} is not valid.\"")
451
452     if len(gripes) == 0:
453         try:
454             this_tau = OLC_codeLen_to_tau[this_codeLen]
455         except:
456             gripes.append(f"\"failed to evaluate OLC_codeLen_to_tau[{this_codeLen}]\"")
457
458     if len(gripes) == 0:
459         try:
460             this_tier = OLC_tau_to_tier[this_tau]
461         except:
462             gripes.append(f"\"failed to evaluate OLC_tau_to_tier[{this_tau}]\"")
463
464     if len(gripes) == 0:
465         lat_hi = round(codeArea.latitudeHi, 12)
466         lat_md = round(codeArea.latitudeCenter, 12)
467         lat_lo = round(codeArea.latitudeLo, 12)
468         lng_hi = round(codeArea.longitudeHi, 12)
469         lng_md = round(codeArea.longitudeCenter, 12)
470         lng_lo = round(codeArea.longitudeLo, 12)
471         lat_parent = this_parent.midpoint.lat
472         lng_parent = this_parent.midpoint.lng
473         zone_dict = {}
474         zone_dict['parent'] = this_parent
475         zone_dict['tau'] = this_tau
476         zone_dict['tier'] = this_tier
477         zone_dict['code'] = this_code
478         zone_dict['origin'] = OLC_Point(lat_lo, lng_lo)
479         zone_dict['midpoint'] = OLC_Point(lat_md, lng_md)
480         zone_dict['offset_from_parent'] = OLC_Offset(lat_parent, lng_parent, lat_md, lng_md)
481         zone_dict[OLC_OFFSET_N] = OLC_Offset(lat_md, lng_md, lat_hi, lng_md)
482         zone_dict[OLC_OFFSET_NE] = OLC_Offset(lat_md, lng_md, lat_hi, lng_hi)
483         zone_dict[OLC_OFFSET_E] = OLC_Offset(lat_md, lng_md, lat_md, lng_hi)
484         zone_dict[OLC_OFFSET_SE] = OLC_Offset(lat_md, lng_md, lat_lo, lng_hi)
485         zone_dict[OLC_OFFSET_S] = OLC_Offset(lat_md, lng_md, lat_lo, lng_md)
486         zone_dict[OLC_OFFSET_SW] = OLC_Offset(lat_md, lng_md, lat_lo, lng_lo)
487         zone_dict[OLC_OFFSET_W] = OLC_Offset(lat_md, lng_md, lat_md, lng_lo)
488         zone_dict[OLC_OFFSET_NW] = OLC_Offset(lat_md, lng_md, lat_hi, lng_lo)
489         this_offsets = [zone_dict[offset] for offset in OLC_DIAGONAL_OFFSETs]
490
491     if len(gripes) == 0:
492         this_tolerances = [zone_offset.mR for zone_offset in this_offsets]
493         this_tolerance = np.round(np.max(this_tolerances), 3)
494         zone_dict['tolerance'] = this_tolerance
495
496     if len(gripes) == 0:
497         self.valid = True
498         self.is_pending = True
499         self._dict_.update(zone_dict)
500         self.children = {}
501     else:
502         self.valid = False
503         self.gripes = gripes

```

```

504
505     def __str__(self):
506         return f"OLC_Zone({repr(self.code)})"
507
508     def __repr__(self):
509         s = []
510         for k in self.__dict__:
511             if k == "parent":
512                 s.append(f".{k}={str(self.parent)}")
513             elif k == "children":
514                 child_zones = ", ".join([str(c_z) for c_z in self.children.values()])
515                 s.append(f".{k}=[{child_zones}]")
516             else:
517                 s.append(f".{k}={repr(self.__dict__[k])}")
518         return '|'.join(s)
519
520     def event_add(self, event_index, event_latlng):
521         gripes = []
522         if self.tau == OLC_TAU_BEDROCK:
523             event_lat, event_lng = event_latlng
524             try:
525                 self.j_index.append(event_index)
526             except:
527                 self.j_index = [event_index]
528             try:
529                 self.j_lat.append(event_lat)
530             except:
531                 self.j_lat = [event_lat]
532             try:
533                 self.j_lng.append(event_lng)
534             except:
535                 self.j_lng = [event_lng]
536         else:
537             child_code = olc_encode_tau(event_latlng, self.tau + 1)
538             child_zone = None
539             try:
540                 child_zone = self.children[child_code]
541             except:
542                 pass
543             if not isinstance(child_zone, type(self)):
544                 child_zone = OLC_Zone(code=child_code, parent=self)
545                 try:
546                     self.children[child_code] = child_zone
547                 except:
548                     gripes.append(f"failed to assign self.children[{repr(child_code)}] = {str(child_zone)}")
549
550             if not isinstance(child_zone, type(self)):
551                 gripes.append(f"child_zone is not an instance of type {type(self)}")
552
553         if len(gripes) == 0:
554             child_zone.event_add(event_index, event_latlng)
555
556         if len(gripes) > 0:
557             for gripe in gripes:
558                 print(gripe)
559
560     return
561
562     def flip_to_pending(self):
563         num_flips = 0
564
565         child_zones = list(self.children.values())
566         for zone_c in child_zones:
567             num_flips += zone_c.flip_to_pending()
568
569         if not self.is_pending:
570             self.is_pending = True
571             num_flips += 1
572
573         return num_flips
574
575     def search(self, Q, K, slate_previous=None, mR_est=None, mN_est=None, mE_est=None):

```

```

576
577     def search_radius(slate, N):
578         try:
579             # if slate already contains a number iNum
580             # of nearest-neighbor candidates, the distance
581             # to the iNumth neighbor (SO FAR) will be
582             # the NEW (working) search radius.
583             radius = slate[N - 1][0]
584         except:
585             # slate does NOT yet contain the needed minimum
586             # number K of nearest-neighbor candidates.
587             # Let the search radius be implausibly large
588             radius = OLC_EARTH_CIRCUMFERENCE
589         return radius
590
591     def slate_display(slate, comment):
592         if len(slate) == 0:
593             print(f"{comment} is empty.")
594             print()
595         else:
596             print(f"{comment} is as follows:")
597             nomDossiers = [nomDoss_j for R_j, code_j, nomDoss_j in slate]
598             for i, nomDoss in enumerate(nomDossiers):
599                 print(f"i={i}: {repr(nomDoss)}")
600             print()
601         return
602
603     if self.tau == OLC_TAU_PLANET:
604         # the current zone represents a planet.
605         #      reset the "is_pending" flags to True
606         #      for the current zone and all its subzones
607         num_flip = self.flip_to_pending()
608
609         #      start with a "clean" slate of nearest neighbors
610         slate_previous = []
611     else:
612         if not self.is_pending:
613             # Nothing to see here. Keep moving.
614             return slate_previous
615
616     slate_b = slate_previous
617
618     gripes = []
619
620     tau_b = self.tau
621     tier_b = self.tier
622     code_b = self.code
623     theta_b = self.tolerance
624
625     if tau_b == OLC_TAU_BEDROCK:
626         # zone_b exists within the BEDROCK tier
627         # therefore zone_b cannot have child zones
628
629         if code_b == Q.code_bedrock:
630             # Bedrock-tier zone_b CONTAINS point Q
631             #      we purposely EXCLUDE zone_b from placement
632             #      onto (its own) K-nearest-neighbors list.
633             self.is_pending = False
634             return slate_b
635         else:
636             # Bedrock-tier zone_b DOES NOT CONTAIN point Q
637             offset_Qb = OLC_Offset(*Q.latlng, *self.midpoint.latlng, mR_est, mN_est, mE_est)
638             self.offset_Qb = offset_Qb
639
640             mN_Qb = offset_Qb.mN
641             delta_mN_Qb = offset_Qb.delta_mN
642             mE_Qb = offset_Qb.mE
643             delta_mE_Qb = offset_Qb.delta_mE
644             mR_Qb = offset_Qb.mR
645             delta_mR_Qb = offset_Qb.delta_mR
646
647             R_Qb_min = np.round(np.max([0, np.subtract(mR_Qb, theta_b)]), 3)

```

```

648         R_Qb_max = np.round(mR_Qb + theta_b, 3)
649
650     if R_Qb_min > search_radius(slate_b, K):
651         # zone_b lies completely OUTSIDE the current search radius.
652         # Ignore it.
653         self.is_pending = False
654         return slate_b
655     else:
656         # parts of zone_b lie WITHIN the prescribed search radius.
657         # Regard zone_b as a potential member of the
658         # K-nearest-neighbors list.
659         nomDoss = {}
660         nomDoss['code'] = code_b
661         nomDoss['event_indices'] = self.j_index
662         nomDoss['R'] = mR_Qb
663         nomDoss['dR'] = delta_mR_Qb
664         nomDoss['R_min'] = R_Qb_min
665         nomDoss['R_max'] = R_Qb_max
666
667         nominee = (mR_Qb, code_b, nomDoss)
668         slate_b = list(sorted([*slate_b, nominee]))[0:K]
669         self.is_pending = False
670         return slate_b
671
672     # At this point, we know that zone_b is ABOVE the bedrock tier
673
674     # Check whether zone_b has any child zones
675     child_codes_raw = list(self.children.keys())
676     if len(child_codes_raw) == 0:
677         # zone_b has no child zones
678         self.is_pending = False
679         return slate_b
680
681     # calculate the geodesic distance from Point Q to the
682     # midpoint of zone_b (this gets used in the estimation of
683     # the distance from Point Q to the center of child zone c, later)
684     offset_Qb = OLC_Offset(*Q.latlng, *self.midpoint.latlng, mR_est, mN_est, mE_est)
685     self.offset_Qb = offset_Qb
686     mN_Qb = offset_Qb.mN
687     mE_Qb = offset_Qb.mE
688
689     # Check whether zone_b has a child zone that CONTAINS point Q
690     code_Qc = Q.code(tau_b + 1)
691     if code_Qc in child_codes_raw:
692         # zone_b has a child zone that CONTAINS point Q.
693         # search that child zone.
694         zone_c = self.children[code_Qc]
695         tier_c = zone_c.tier
696         mN_bc = zone_c.offset_from_parent.mN
697         mE_bc = zone_c.offset_from_parent.mE
698         mN_Qc_est = np.round(mN_Qb + mN_bc, 3)
699         mE_Qc_est = np.round(mE_Qb + mE_bc, 3)
700         mR_Qc_est = np.round(np.sqrt(mN_Qc_est**2 + mE_Qc_est**2), 3)
701         slate_new = zone_c.search(Q, K, slate_b, mR_Qc_est, mN_Qc_est, mE_Qc_est)
702         slate_b = slate_new
703         child_codes = [code_c for code_c in child_codes_raw if not code_c == code_Qc]
704     else:
705         child_codes = child_codes_raw
706
707     if len(child_codes) == 0:
708         # we have processed all child codes of interest
709         # within zone_b
710         self.is_pending = False
711         return slate_b
712
713     child_zones = [self.children[code_c] for code_c in child_codes]
714     zones_pending = [zone_c for zone_c in child_zones if zone_c.is_pending]
715     if len(zones_pending) == 0:
716         # zone_b has no PENDING child zones
717         self.is_pending = False
718         return slate_b
719

```

```

720     # At this point, zone_b still has one or more PENDING child zones
721
722     # Estimation phase: calculate ESTIMATED geodesic distance between Q
723     # and the midpoint of each child zone, zone_c
724     docket_raw = []
725     for zone_c in zones_pending:
726         theta_c = zone_c.tolerance
727         mN_bc = zone_c.offset_from_parent.mN
728         mE_bc = zone_c.offset_from_parent.mE
729         mN_Qc_est = np.round(mN_Qb + mN_bc, 3)
730         mE_Qc_est = np.round(mE_Qb + mE_bc, 3)
731         mR_Qc_est = np.round(np.sqrt(mN_Qc_est**2 + mE_Qc_est**2), 3)
732         R_Qc_min = np.round(np.max([0, np.subtract(mR_Qc_est, theta_c)]), 3)
733         R_Qc_max = np.round(mR_Qc_est + theta_c, 3)
734         if R_Qc_min > search_radius(slate_b, K):
735             # all parts of zone_c are completely OUTSIDE
736             # the prescribed search radius
737             # Exclude this child zone from the "docket"
738             zone_c.is_pending = False
739         else:
740             # some parts of zone_c lie WITHIN
741             # the prescribed search radius
742             # put zone_c on the "docket"
743             task_c = (mR_Qc_est, mN_Qc_est, mE_Qc_est, zone_c)
744             docket_raw.append(task_c)
745
746     if len(docket_raw) == 0:
747         # None of the child zones warranted a detailed search
748         self.is_pending = False
749         return slate_b
750
751     docket = sorted(docket_raw)
752     for mR_Qc_est, mN_Qc_est, mE_Qc_est, zone_c in docket:
753         slate_c = zone_c.search(Q, K, slate_b, mR_Qc_est, mN_Qc_est, mE_Qc_est)
754         slate_b = slate_c
755
756     self.is_pending = False
757
758     return slate_b
759
760 def olc_K_nearest_neighbors(events_index, events_latlng, K):
761
762     planet = OLC_Zone(OLC_CODE_PLANET)
763     print("AFTER __init__...")
764     print(repr(planet))
765     print()
766
767     num_events = len(events_index)
768
769     events_i_latlng = list(zip(events_i, events_latlng))
770
771     for i, latlng in events_i_latlng:
772         planet.event_add(event_index=i, event_latlng=latlng)
773     print()
774
775     print("AFTER ADDING ALL EVENTS...")
776     print(repr(planet))
777     print()
778
779     events_KNN = []
780     print("====")
781     print(f"STARTED KNN EVALUATION AT {datetime.now()}")
782     print()
783     print("====")
784     for i, Q_latlng in events_i_latlng:
785         write_line = False
786         if i == 0:
787             write_line = True
788         if i == num_events - 1:
789             write_line = True
790         if i % 50 == 0:
791             write_line = True

```

```

792
793     if write_line:
794         print(f"event_i={i}, time={datetime.now()}")
795     Q_i = OLC_Point(*Q_latlng)
796     slate_i = planet.search(Q_i, K)
797     neighbor_eventIndices_i = []
798     for R_ij, code_ij, nomDoss_ij in slate_i:
799         indices_ij = nomDoss_ij['event_indices']
800         for neighbor_eventIndex in indices_ij:
801             neighbor_eventIndices_i.append(neighbor_eventIndex)
802     KNN_i = neighbor_eventIndices_i[0:K]
803     events_KNN.append(KNN_i)
804
805 return events_KNN
806
807 test_the_functions = False
808 if test_the_functions:
809     A_latlng = (47.5111, -122.2571) # bedrock code = '84VVGP6V+C5'
810     B_latlng = (47.7210, -122.319) # bedrock code = '84VVPMCJ+CC'
811     P_latlng = (47.7379, -122.233) # bedrock code = '84VVPQQ8+5R'
812     Q_latlng = (47.5208, -122.393) # bedrock code = '84VVGJC4+8R'
813
814     point_A = OLC_Point(*A_latlng)
815     print("point_A after __init__")
816     print(repr(point_A))
817     print()
818
819     point_B = OLC_Point(*B_latlng)
820     print("point_B after __init__")
821     print(repr(point_B))
822     print()
823
824     point_P = OLC_Point(*P_latlng)
825     print("point_P after __init__")
826     print(repr(point_P))
827     print()
828
829     point_Q = OLC_Point(*Q_latlng)
830     print("point_Q after __init__")
831     print(repr(point_Q))
832     print()
833
834     Q_code_bedrock = point_Q.code(tier="bedrock")
835     print(f"Q_code_bedrock={repr(Q_code_bedrock)}")
836     print()
837
838     Q_code_holler = point_Q.code(tier="holler")
839     print(f"Q_code_holler={repr(Q_code_holler)}")
840     print()
841
842     print("after point_Q.code(tier='holler') was invoked, point_Q=")
843     print(repr(point_Q))
844     print()
845
846     offset_PQ = OLC_Offset(*P_latlng, *Q_latlng)
847     print("offset_PQ =")
848     print(repr(offset_PQ))
849     print()
850
851     # Create an empty instance OLC_Zone() called "planet"
852     # this will be the "planet Zone"
853     planet = OLC_Zone()
854     print("AFTER __init__...")
855     print(repr(planet))
856     print()
857
858     # Add some events to the "planet zone"
859     num_events = 4
860     event_i = list(range(num_events))
861     event_latlng = [A_latlng, B_latlng, P_latlng, Q_latlng]
862     event_i_latlng = list(zip(event_i, event_latlng))
863     for i, latlng in event_i_latlng:

```

```
864     planet.event_add(event_index=i, event_latlng=latlng)
865     print(f"AFTER ADDING num_events={num_events}...")
866     print(repr(planet))
867     print()
868
869
870     print("code_to_zone dictionary =")
871     code_to_zone = olc_add_zone_to_archive(planet)
872     display(list(code_to_zone.keys())[0:100])
873
874 print("Please proceed to next cell.")
```

Please proceed to next cell.

```
In [3]: 1 # =====
2 #   IMPORT DATA
3 #
4 #   Imports data from either the preferred csv file,
5 #       or from the alternate csv file.
6 #
7 # Michael Collins, Flatiron School
8 # 2020-10-28_0526_MDT
9 # =====
10
11 GLOBAL_ROW = 0
12 GLOBAL_ROW_LAST = 4
13
14 def housing_convert_view(v_str):
15     v_value = int(0)
16     try:
17         v_numeric = pd.to_numeric(v_str, errors='raise')
18         view_nonzero = [1, 2, 3, 4]
19         if v_numeric in view_nonzero:
20             v_value = int(v_numeric)
21     except:
22         pass
23     return v_value
24
25 def housing_convert_waterfront(v_str):
26     v_value = int(0)
27     try:
28         v_numeric = pd.to_numeric(v_str, errors='raise')
29         if v_numeric > 0:
30             v_value = int(v_numeric)
31     except:
32         pass
33     return v_value
34
35 def housing_convert_yr_renovated(v_str):
36     v_value = int(0)
37     try:
38         v_numeric = pd.to_numeric(v_str, errors='raise')
39         year_first = 1900
40         year_last = datetime.now().year
41         years_nonzero = list(range(year_first, year_last + 1))
42         if v_numeric in years_nonzero:
43             v_value = int(v_numeric)
44     except:
45         pass
46     return v_value
47
48 def housing_convert_sfBasement(v_str):
49     v_value = int(0)
50     try:
51         v_numeric = pd.to_numeric(v_str, errors='raise')
52         if v_numeric > 0:
53             v_value = int(v_numeric)
54     except:
55         pass
56     return v_value
57
58 housing_columns_raw = ['id',           # included in housing_dtype
59                       'date',          # included in housing_parse_dates
60                       'price',
61                       'bedrooms',
62                       'bathrooms',
63                       'sqft_living',
64                       'sqft_lot',
65                       'floors',
66                       'waterfront',    # uses housing_convert_waterfront
67                       'view',          # uses housing_convert_view
68                       'condition',
69                       'grade',
70                       'sqft_above',
71                       'sqft_basement', # uses housing_convert_sfBasement
```

```

72     'yr_built',
73     'yr_renovated', # uses housing_convert_yr_renovated
74     'zipcode',
75     'lat',
76     'long',
77     'sqft_living15',
78     'sqft_lot15']
79 housing_dtype = {'id':str}
80 housing_parse_dates = ['date']
81 housing_converters = {'view':housing_convert_view,
82                       'waterfront':housing_convert_waterfront,
83                       'yr_renovated':housing_convert_yr_renovated,
84                       'sqft_basement':housing_convert_sfBasement}
85
86 csv_folder = './data/'
87 csv_prefix_preferred = 'kc_house_data_KNN'
88 csv_prefix_alternate = 'kc_house_data'
89 csv_suffix = '.csv'
90 csvPath_preferred = csv_folder + csv_prefix_preferred + csv_suffix
91 csvPath_alternate = csv_folder + csv_prefix_alternate + csv_suffix
92
93 df_raw = None
94
95 if not isinstance(df_raw, pd.DataFrame):
96     print(f"Attempting to import data from csv file {repr(csvPath_preferred)}...")
97     try:
98         GLOBAL_ROW = 0
99         df_raw = pd.read_csv(csvPath_preferred,
100                             converters=housing_converters,
101                             parse_dates=housing_parse_dates,
102                             dtype=housing_dtype
103                             ).astype({'yr_renovated':int})
104         print(f"      Imported data from csv file {repr(csvPath_preferred)}.\n")
105     except:
106         print(f"Failed to import data from csv file {repr(csvPath_preferred)}.\n")
107
108 if not isinstance(df_raw, pd.DataFrame):
109     print(f"Attempting to import FILTERED data from csv file {repr(csvPath_alternate)}...")
110     try:
111         GLOBAL_ROW = 0
112         df_raw = pd.read_csv(csvPath_alternate,
113                             converters=housing_converters,
114                             parse_dates=housing_parse_dates,
115                             dtype=housing_dtype
116                             ).astype({'yr_renovated':int})
117         print(f"      Imported FILTERED data from csv file {repr(csvPath_alternate)}.\n")
118     except:
119         print(f"Failed to import FILTERED data from csv file {repr(csvPath_alternate)}.\n")
120
121 if not isinstance(df_raw, pd.DataFrame):
122     print(f"Attempting to import UNFILTERED data from csv file {repr(csvPath_alternate)}...")
123     try:
124         df_raw = pd.read_csv(csvPath_alternate)
125         print(f"      Imported UNFILTERED data from csv file {repr(csvPath_alternate)}.\n")
126     except:
127         print(f"Failed to import UNFILTERED data from csv file {repr(csvPath_alternate)}.\n")
128
129 if not isinstance(df_raw, pd.DataFrame):
130     print("Unable to continue without data.\n")
131     assert False
132
133 print("===== df_raw COLUMNS =====")
134 columns_df_raw = list(df_raw.columns)
135 print(repr(columns_df_raw))
136
137 show_dfraw_value_counts = True
138 if show_dfraw_value_counts:
139     print("===== df_raw Value Counts =====")
140     columns_showVC = ['waterfront', 'view', 'grade', 'yr_built', 'yr_renovated', 'zipcode']
141     print(f"for columns = {repr(columns_showVC)}\n")
142     for c in columns_showVC:
143         print(f"VALUE COUNTS for column = {repr(c)}:")

```

```

144     print(df_raw[c].value_counts())
145     print()
146
147
148
149
150
151
152 df = df_raw.copy()
153 print("===== df (original working copy) =====")
154 display(df)
155 print()
156
157
158 col_lat = 'lat'
159 col_long = 'long'
160 col_price = 'price'
161 col_yr_built = 'yr_built'
162 col_yr_renovated = 'yr_renovated'
163 col_bedrooms = 'bedrooms'
164 col_bathrooms = 'bathrooms'
165 col_floors = 'floors'
166 col_waterfront = 'waterfront'
167 col_view = 'view'
168 col_condition = 'condition'
169 col_grade = 'grade'
170 col_sfLiving = 'sfLiving'
171 col_sfAbove = 'sfAbove'
172 col_sfBasement = 'sfBasement'
173 col_sfYard = 'sfYard'
174
175 events_i = list(df.index.values)
176 num_events = len(events_i)
177 events_lat = [round(lat, 3) for lat in df[col_lat].to_list()]
178 events_lng = [round(lng, 3) for lng in df[col_long].to_list()]
179 events_lating = list(zip(events_lat, events_lng))
180 events_price = [p for p in df[col_price].to_list()]
181 events_bedrooms = [b for b in df[col_bedrooms].to_list()]
182 events_view = [v for v in df[col_view].to_list()]
183 events_bathrooms = [b for b in df[col_bathrooms].to_list()]
184 events_grade = [g for g in df[col_grade].to_list()]
185 events_yr_built = [y for y in df[col_yr_built].to_list()]
186 events_yr_renovated = [y for y in df[col_yr_renovated].to_list()]
187 events_sfLiving = [sf for sf in df[col_sfLiving].to_list()]
188 events_sfAbove = [sf for sf in df[col_sfAbove].to_list()]
189 events_sfLiving_sfAbove = list(zip(events_sfLiving, events_sfAbove))
190 events_i_sfLiving_sfAbove = list(zip(events_i, events_sfLiving, events_sfAbove))
191
192
193
194
195 # The quanta_*** are lists of the unique values that appear
196 # in the cleaned-up versions of selected columns
197 quanta_bedrooms = ['1', '2', '3', '4', '5', '6', '7', '8', '9', '10', '11', '33']
198 quanta_bathrooms = ['0.5', '0.75', '1.0', '1.25', '1.5', '1.75', '2.0', '2.25',
199             '2.5', '2.75', '3.0', '3.25', '3.5', '3.75', '4.0', '4.25',
200             '4.5', '4.75', '5.0', '5.25', '5.5', '5.75', '6.0', '6.25',
201             '6.5', '6.75', '7.5', '7.75', '8.0']
202 quanta_floors = ['1.0', '1.5', '2.0', '2.5', '3.0', '3.5']
203 quanta_waterfront = ['0', '1']
204 quanta_view = ['0', '1', '2', '3', '4']
205 quanta_condition = ['1', '2', '3', '4', '5']
206 quanta_grade = ['3', '4', '5', '6', '7', '8', '9', '10', '11', '12', '13']
207
208 col_metascore = 'metascore'
209 events_metascore = []
210 df[col_metascore] = (df[col_bedrooms]*0 +
211                     df[col_bathrooms]*0 +
212                     df[col_floors]*0 +
213                     df[col_waterfront]*0 +
214                     df[col_view]*1 +
215                     df[col_condition]*1 +

```

```

216         df[col_grade]*1)
217 events_metascore = list(df[col_metascore])
218
219 # Ignore certain outliers
220 events_i_bedrooms = list(zip(events_i, events_bedrooms))
221 eventIndices_BEDROOMS_FAIL = [i for i, beds in events_i_bedrooms if beds > 11]
222 events_i_bathrooms = list(zip(events_i, events_bathrooms))
223 eventIndices_BATHROOMS_FAIL = [i for i, baths in events_i_bathrooms if baths < 0.75]
224
225
226 eventIndices_ignore1 = [264, 1166, 1761, 3950, 4365, 4577, 4907, 6096,
227     6825, 8149, 8537, 11544, 14813, 14821, 16404, 16556,
228     17138, 18212, 18261, 18352, 19089, 19173, 20836, 21185]
229 eventIndices_ignore2 = [540, 700, 1313, 1784, 1943, 2624, 3488, 4014, 4020, 5068,
230     5084, 5685, 5961, 6396, 6685, 7070, 7245, 8050, 8172, 8270,
231     8597, 9027, 9096, 9309, 9503, 10416, 10527, 10977, 11674,
232     11736, 11859, 12122, 12764, 12993, 13110, 13451, 13616, 13757,
233     13774, 14625, 15025, 15246, 15363, 17460, 18194, 18273, 18833,
234     19002, 19388, 19668, 20138, 20309, 20436, 20519, 20812, 20958,
235     21034, 21356, 21469, 21532, 21560]
236 eventIndices_ignore3 = []
237 eventIndices_ignore = list(sorted(eventIndices_ignore1 +
238     eventIndices_ignore2 +
239     eventIndices_ignore3 +
240     eventIndices_BEDROOMS_FAIL +
241     eventIndices_BATHROOMS_FAIL))
242
243 events_valid = [not (i in eventIndices_ignore) for i in events_i]
244 events_intValid = [int(ev) for ev in events_valid]
245
246
247 col_eventValid = 'eventValid'
248 df.insert(loc=0, column=col_eventValid, value=events_intValid)
249
250 col_eventIndex = 'eventIndex'
251 df.insert(loc=0, column=col_eventIndex, value=events_i)
252
253
254 # =====
255 # Check whether needed CUSTOM data series (columns)
256 # are already present in the imported data.
257 # If a needed custom series is NOT already present,
258 # construct the data for that series, and then
259 # ADD that series to the data frame.
260 # =====
261
262 print("looking up building_age...")
263 print()
264 col_building_age = 'building_age'
265 if col_building_age in columns_df_raw:
266     events_building_age = df[col_building_age]
267 else:
268     year_now = datetime.now().year
269     events_building_age = [int(year_now - yr) for yr in events_yr_built]
270     df[col_building_age] = events_building_age
271
272 print("looking up interior_age...")
273 print()
274 # Number of years since renovation (if yr_renovated is nonzero), otherwise the
275 # number of years since original construction
276 col_interior_age = 'interior_age'
277 if col_interior_age in columns_df_raw:
278     events_interior_age = df[col_interior_age]
279 else:
280     events_yrBldg_yrReno = list(zip(events_yr_built, events_yr_renovated))
281     events_i_yrBldg_yrReno = list(zip(events_i, events_yr_built, events_yr_renovated))
282     for i, yrBldg, yrReno in events_i_yrBldg_yrReno:
283         if not yrReno == 0:
284             if yrReno < yrBldg:
285                 print(f"row {i} of {num_events}: yr_built={yrBldg}, yr_renovated={yrReno}); yr_renovated is not plausible.")
286     year_now = datetime.now().year
287     events_interior_age = [(int(year_now - yrReno) if yrReno > 0 else int(year_now - yrBldg))]
```

```

288             for yrBldg, yrReno in events_yrBldg_yrReno]
289         df[col_interior_age] = events_interior_age
290
291     print("forcing sfBasement = (sfLiving - sfAbove) ...")
292     print()
293     col_sfBasement = 'sqft_basement'
294     for i, sfLiving, sfAbove in events_i_sfLiving_sfAbove:
295         if sfLiving < sfAbove:
296             print(f"row {i} of {num_events}: sfLiving={sfLiving}, sfAbove={sfAbove}); sfBasement cannot be negative; using zero")
297     events_sfBasement = [np.max([0, int(sfLiving - sfAbove)]) for sfLiving, sfAbove in events_sfLiving_sfAbove]
298     df[col_sfBasement] = events_sfBasement
299
300     print("looking up log10_price...")
301     print()
302     col_log10_price = 'log10_price'
303     if col_log10_price in columns_df_raw:
304         events_log10_price = df[col_log10_price]
305     else:
306         events_log10_price = [round(np.log10(p), 3) for p in events_price]
307     df[col_log10_price] = events_log10_price
308
309
310
311
312     def sim_KNN(iEvent, eventIndices, iNum):
313         eventIndices_X = [ei for ei in eventIndices if not ei == iEvent]
314         indices_drawn = np.random.choice(eventIndices_X, iNum, replace=False)
315         return list(indices_drawn)
316
317     K_table = 50
318     print(f"looking up indices of nearest K={K_table} events...")
319     print()
320     simulate_KNN = False
321     if simulate_KNN:
322         col_KNN = 'sim_KNN'
323     else:
324         col_KNN = 'KNN'
325
326     time_start = datetime.now()
327     if col_KNN in columns_df_raw:
328         print(f"===== FOUND events_KNN information in df_raw ======\n")
329         events_KNN = [[int(s) for s in sl[1:-1].split(',')] for sl in df[col_KNN]]
330     else:
331         if simulate_KNN:
332             print(f"=====SIMULATING olc_K_nearest_neighbors(events_i, events_latlng, K={K_table}) ======\n")
333             events_KNN = []
334             for i in events_i:
335                 write_line = False
336                 if i == 0:
337                     write_line = True
338                 if i == num_events - 1:
339                     write_line = True
340                 if i % 1000 == 0:
341                     write_line = True
342
343                 if write_line:
344                     print(f"i={i}, datetime.now()={datetime.now()}")
345                     KNN_i = sim_KNN(i, events_i, K_table)
346                     events_KNN.append(KNN_i)
347             print(f"=====DONE SIMULATING olc_K_nearest_neighbors(events_i, events_latlng, K={K_table}) ======\n")
348         else:
349             print(f"=====CALLING olc_K_nearest_neighbors(events_i, events_latlng, K={K_table}) ======\n")
350             events_KNN = olc_K_nearest_neighbors(events_i, events_latlng, K_table)
351             print(f"=====DONE WITH olc_K_nearest_neighbors(events_i, events_latlng, K={K_table}) ======\n")
352
353     df[col_KNN] = events_KNN
354
355
356
357 # =====#
358 #   Derive target and features from K-nearest-neighbors lists
359 # =====#

```

```

360
361 columns_v =['price',
362     'bedrooms',
363     'bathrooms',
364     'floors',
365     'waterfront',
366     'view',
367     'condition',
368     'grade',
369     'metascore',
370     'sqft_living',
371     'sqft_above',
372     'sqft_basement',
373     'sqft_lot',
374     'building_age',
375     'interior_age']
376
377 K_use = 15
378 col_price_median = f"{col_price}_median{K_use}"
379 for col_v in columns_v:
380
381     print(f"Generating new columns for {repr(col_v)} using K={K_use} nearest neighbors...")
382
383     events_v = df[col_v]
384     events_v_mean = []
385     events_v_median = []
386     events_v_diff = []
387     events_v_delta = []
388     events_v_sigma = []
389     events_v_MAD = []
390     events_v_zScore = []
391     events_v_zed = []
392     events_v_phi = []
393
394     for i in events_i:
395         if events_valid[i]:
396             p_i = events_price[i]
397             v_i = events_v[i]
398             eventIndices_i = events_KNN[i]
399             eventIndices_i_valid = [ei for ei in eventIndices_i if events_valid[ei]][0:K_use]
400             p_knn_valid = [events_price[eiv] for eiv in eventIndices_i_valid]
401             v_knn_valid = [events_v[eiv] for eiv in eventIndices_i_valid]
402             p_mean = np.mean(p_knn_valid)
403             p_median = np.median(p_knn_valid)
404             p_diff = p_i - p_mean
405             v_mean = np.mean(v_knn_valid)
406             v_median = np.median(v_knn_valid)
407             v_diff = v_i - v_mean
408             v_delta = v_i - v_median
409             v_sigma = np.std(v_knn_valid)
410             v_MAD = MAD(v_knn_valid)
411
412             # Calculate "zScore"
413             if v_sigma == 0:
414                 v_zScore = 0
415             else:
416                 v_zScore = v_diff / v_sigma
417                 if v_zScore < -5:
418                     v_zScore = -5
419                 if v_zScore > 5:
420                     v_zScore = 5
421
422             # Calculate "zed"
423             if col_v == col_price:
424                 # Calculate "zed" for the price
425                 if p_mean == 0:
426                     # This is unlikely, because houses have positive prices
427                     v_zed = 0
428                 else:
429                     v_zed = p_diff / p_mean
430                     if v_zed < -5:
431                         v_zed = -5

```

```

432         if v_zed > 5:
433             v_zed = 5
434     else:
435         # Calculate "zed" for other quantities
436         v_zed = (1 / 100) * v_diff
437
438         # Calculate "phi"
439         if p_median == 0:
440             v_phi = 0
441         else:
442             v_phi = (p_median / 100) * (v_delta)
443             if v_phi < -5:
444                 v_phi = -5
445             if v_phi > 5:
446                 v_phi = 5
447
448     else:
449         v_mean = None
450         v_median = None
451         v_diff = None
452         v_delta = None
453         v_sigma = None
454         v_MAD = None
455         v_zScore = None
456         v_zed = None
457         v_phi = None
458
459         events_v_mean.append(v_mean)
460         events_v_median.append(v_median)
461         events_v_diff.append(v_diff)
462         events_v_delta.append(v_delta)
463         events_v_sigma.append(v_sigma)
464         events_v_MAD.append(v_MAD)
465         events_v_zScore.append(v_zScore)
466         events_v_zed.append(v_zed)
467         events_v_phi.append(v_phi)
468
469 df["{col_v}_mean{K_use}"] = events_v_mean
470 df["{col_v}_median{K_use}"] = events_v_median
471 df["{col_v}_diff{K_use}"] = events_v_diff
472 df["{col_v}_delta{K_use}"] = events_v_delta
473 df["{col_v}_sigma{K_use}"] = events_v_sigma
474 df["{col_v}_MAD{K_use}"] = events_v_MAD
475 df["{col_v}_zScore{K_use}"] = events_v_zScore
476 df["{col_v}_zed{K_use}"] = events_v_zed
477 df["{col_v}_phi{K_use}"] = events_v_phi
478
479 print(f"Done generating columns of MEDIANs, DELTAs, MADs, and PHIs using K={K_use} nearest neighbors.\n")
480
481
482 # =====
483 # Summarize any changes that were made to the columns
484 # that appear in the data frame.
485 # =====
486
487 columns_df = list(df.columns)
488 print("===== BEGIN Summary of Modifications to Data Frame =====")
489 print()
490 print("_____ all original columns in df_raw:")
491 print(f"columns_df_raw = {repr(columns_df_raw)}")
492 print()
493 print("_____ original columns removed from df:")
494 original_columns_removed_from_df = [cc for cc in columns_df_raw if not cc in columns_df]
495 print(f"original_columns_removed_from_df = {repr(original_columns_removed_from_df)}")
496 print()
497 print("_____ original columns preserved in df:")
498 original_columns_preserved_in_df = [cc for cc in columns_df if cc in columns_df_raw]
499 print(f"original_columns_preserved_in_df = {repr(original_columns_preserved_in_df)}")
500 print()
501 print("_____ new columns added to df:")
502 new_columns_added_to_df = [cc for cc in columns_df if not cc in columns_df_raw]
503 print(f"new_columns_added_to_df = {repr(new_columns_added_to_df)}")

```

```

504 print()
505 print("      all columns in MODIFIED df:")
506 print(f"columns_df = {repr(columns_df)}")
507 print()
508 print("===== END Summary of Modifications to Data Frame =====")
509 print()
510
511 # =====
512 #   Auto-save the current data frame to a csv file,
513 #       regardless of whether it differs from df_raw.
514 # =====
515
516 stamp_now = datetime.now().strftime("%Y-%m-%d_%H%M%S")
517 print(f"===== df (as of {stamp_now}) =====")
518 display(df)
519 print()
520 modified_csv = './data/kc_house_data_modified_' + stamp_now + '.csv'
521 try:
522     df.to_csv(modified_csv, index=False)
523     print(f"df (with modified columns) was written to file [{modified_csv}].")
524 except:
525     print(f"FAILED to write df (with modified columns) to file [{modified_csv}].")
526 print()
527
528
529
530 # =====
531 #   Save a subset of the dataframe to a new csv file.
532 #
533 # =====
534
535 ext_m = "_median"
536 # ext_d = "_delta"
537 ext_f = "_zed"
538 ext_K = K_use
539 columns_subset =['eventIndex',
540                  # 'price',
541                  f"price{ext_m}{ext_K}",
542                  # f"price{ext_d}{ext_K}",
543                  f"price{ext_f}{ext_K}",
544                  f"bedrooms{ext_f}{ext_K}",
545                  f"bathrooms{ext_f}{ext_K}",
546                  f"floors{ext_f}{ext_K}",
547                  f"waterfront{ext_f}{ext_K}",
548                  f"view{ext_f}{ext_K}",
549                  f"condition{ext_f}{ext_K}",
550                  f"grade{ext_f}{ext_K}",
551                  f"metascore{ext_f}{ext_K}",
552                  f"sqft_living{ext_f}{ext_K}",
553                  f"sqft_above{ext_f}{ext_K}",
554                  f"sqft_basement{ext_f}{ext_K}",
555                  f"sqft_lot{ext_f}{ext_K}",
556                  f"building_age{ext_f}{ext_K}",
557                  f"interior_age{ext_f}{ext_K}"]
558
559 df_subset = df[columns_subset].loc[df[col_eventValid] == True]
560
561 stamp_now = datetime.now().strftime("%Y-%m-%d_%H%M%S")
562 print(f"===== df_subset (as of {stamp_now}) =====")
563 display(df_subset)
564 print()
565 subset_csv = f"./data/kc_house_data_subset_K={K_use}.csv"
566 try:
567     df_subset.to_csv(subset_csv, index=False)
568     print(f"df (with subset columns) was written to file [{subset_csv}].")
569 except:
570     print(f"FAILED to write df (with modified columns) to file [{subset_csv}].")
571 print()
572
573 df_clean = df_subset # df_subset.drop(columns=['price',
574                                     #                 f"price{ext_m}{K_use}",
575                                     #                 f"price{ext_d}{K_use}"]

```

```

576 #          J)
577
578 stamp_now = datetime.now().strftime("%Y-%m-%d %H%M%S")
579 print(f"===== df_clean (as of {stamp_now}) =====")
580 display(df_clean)
581 print()
582

      2      2    519900.0   -0.617753   -0.013333   -0.012667   -0.007333   0.000000   0.000000   -0.000667   -0.020667   -0.021333   -14.767333   -14.3406
      3      3    560000.0   -0.046466   0.010000   0.006500   -0.005667   -0.000667   -0.009333   0.017333   -0.006000   0.002000   -1.240000   -5.0533
      4      4    525000.0   0.006215   -0.002000   -0.003000   -0.006000   0.000000   0.000000   -0.000667   -0.000667   -0.001333   -2.500000   -0.9000
     ...
     ...
     ...
21592 21592 400000.0   -0.122729   0.000000   0.003667   0.008667   0.000000   0.000000   -0.002000   0.006000   0.004000   0.733333   1.4800
21593 21593 340000.0   0.196342   0.006667   0.005500   0.006000   0.000000   0.000000   -0.000667   0.008000   0.007333   5.526667   7.3866
21594 21594 387000.0   -0.014780   -0.006667   -0.012167   0.001333   0.000000   0.000000   -0.001333   -0.000667   -0.002000   -3.500000   -1.8933
21595 21595 387000.0   -0.205077   -0.002667   -0.001333   0.000000   0.000000   -0.001333   0.000000   0.000667   -0.000667   -4.320000   -1.7733
21596 21596 387000.0   -0.203692   -0.006667   -0.012167   0.001333   0.000000   0.000000   -0.001333   -0.000667   -0.002000   -3.500000   -1.8933

21507 rows × 17 columns

```

```

In [4]: 1 # =====
2 #      Ancillary functions related to OLS
3 #
4 # Michael Collins, Flatiron School
5 # 2020-10-30_0639_MDT
6 #
7 # =====
8
9 def plot_corr(df):
10     corr = df.corr().abs()
11     plt.figure(figsize=(8, 5))
12     sns.heatmap(corr, fmt='0.2g', annot=True, cmap=sns.color_palette('Blues'))
13     plt.show()
14
15
16 # VIF test is a great test for multicollinearity
17 def calculate_vif(df, target_col, show_res=False):
18     x = df.drop(columns=[target_col])
19     y = df[target_col]
20     ols = sm.OLS(y, x).fit()
21     if show_res:
22         display(ols.summary())
23     vif = 1 / (1 - ols.rsquared)
24     return vif
25
26 print("Please proceed to the next cell.")

```

Please proceed to the next cell.

```
In [5]: 1 # =====
2 #       Selection of Features to Include in OLS
3 #
4 # Michael Collins, Flatiron School
5 # 2020-10-30_0639_MDT
6 #
7 # =====
8
9 stamp_now = datetime.now().strftime("%Y-%m-%d %H%M%S")
10 print(f"===== df_clean (as of {stamp_now}) =====")
11 display(df_clean)
12 print()
13
14 ext_f = "_zed"
15 ext_K = K_use
16 cols_every = list(df_clean.columns)
17 #       Features that are commented out, below,
18 #       are excluded from consideration
19 cols_feature = [
20     f"bedrooms{ext_f}{ext_K}",
21     f"bathrooms{ext_f}{ext_K}",
22     f"floors{ext_f}{ext_K}",
23     f"waterfront{ext_f}{ext_K}",
24     # f"view{ext_f}{ext_K}",
25     # f"condition{ext_f}{ext_K}",
26     # f"grade{ext_f}{ext_K}",
27     f"metascore{ext_f}{ext_K}",
28     f"sqft_living{ext_f}{ext_K}",
29     # f"sqft_above{ext_f}{ext_K}",
30     # f"sqft_basement{ext_f}{ext_K}",
31     # f"sqft_lot{ext_f}{ext_K}",
32     f"building_age{ext_f}{ext_K}",
33     # f"interior_age{ext_f}{ext_K}"
34   ]
35 cols_non_feature = [c for c in cols_every if not c in cols_feature]
36 col_target = f"price{ext_f}{ext_K}"
37
38 print(f"===== cols_every =====")
39 print(f"cols_every = {repr(cols_every)}\n")
40 print(f"===== cols_feature =====")
41 print(f"cols_feature = {repr(cols_feature)}\n")
42 print(f"===== cols_non_feature =====")
43 print(f"cols_non_feature = {repr(cols_non_feature)}\n")
44 print(f"===== col_target =====")
45 print(f"col_target = {repr(col_target)}\n")

```

===== df\_clean (as of 2020-10-31\_083203) =====

	eventIndex	price_median15	price_zed15	bedrooms_zed15	bathrooms_zed15	floors_zed15	waterfront_zed15	view_zed15	condition_zed15	grade_zed15	metascore_zed15	sqft_living_zed15	sqft_above_zed
0	0	400000.0	-0.420485	-0.009333	-0.007667	-0.005000	0.000000	-0.008000	-0.006667	-0.001333	-0.016000	-8.614667	-3.8280
1	1	402000.0	0.388877	0.001333	0.008167	0.006667	0.000000	0.000000	-0.006667	-0.001333	-0.008000	11.451333	9.3480
2	2	519900.0	-0.617753	-0.013333	-0.012667	-0.007333	0.000000	0.000000	-0.000667	-0.020667	-0.021333	-14.767333	-14.3406
3	3	560000.0	-0.046466	0.010000	0.006500	-0.005667	-0.000667	-0.009333	0.017333	-0.006000	0.002000	-1.240000	-5.0533
4	4	525000.0	0.006215	-0.002000	-0.003000	-0.006000	0.000000	0.000000	-0.000667	-0.000667	-0.001333	-2.500000	-0.9000
...	...	...	...	...	...	...	...	...	...	...	...	...	...
21592	21592	400000.0	-0.122729	0.000000	0.003667	0.008667	0.000000	0.000000	-0.002000	0.006000	0.004000	0.733333	1.4800
21593	21593	340000.0	0.196342	0.006667	0.005500	0.006000	0.000000	0.000000	-0.000667	0.008000	0.007333	5.526667	7.3866
21594	21594	387000.0	-0.014780	-0.006667	-0.012167	0.001333	0.000000	0.000000	-0.001333	-0.000667	-0.002000	-3.500000	-1.8933
21595	21595	387000.0	-0.205077	-0.002667	-0.001333	0.000000	0.000000	-0.001333	0.000000	0.000667	-0.000667	-4.320000	-1.7733
21596	21596	387000.0	-0.203692	-0.006667	-0.012167	0.001333	0.000000	0.000000	-0.001333	-0.000667	-0.002000	-3.500000	-1.8933

21507 rows × 17 columns

```
===== cols_every =====
cols_every = ['eventIndex', 'price_median15', 'price_zed15', 'bedrooms_zed15', 'bathrooms_zed15', 'floors_zed15', 'waterfront_zed15', 'view_zed15', 'condition_zed15', 'grad
d15', 'sqft_living_zed15', 'sqft_above_zed15', 'sqft_basement_zed15', 'sqft_lot_zed15', 'building_age_zed15', 'interior_age_zed15']

===== cols_feature =====
cols_feature = ['bedrooms_zed15', 'bathrooms_zed15', 'floors_zed15', 'waterfront_zed15', 'metascore_zed15', 'sqft_living_zed15', 'building_age_zed15']

===== cols_non_feature =====
cols_non_feature = ['eventIndex', 'price_median15', 'price_zed15', 'view_zed15', 'condition_zed15', 'grade_zed15', 'sqft_above_zed15', 'sqft_basement_zed15', 'sqft_lot_zed1
5']

===== col_target =====
col_target = 'price_zed15'
```

```
In [6]: 1 df_features = df_clean.drop(columns=cols_non_feature)
2 stamp_now = datetime.now().strftime("%Y-%m-%d %H%M%S")
3 print(f"===== df_features (as of {stamp_now}) =====")
4 display(df_features)
5 print()
6
7
8
9 # lets find some features based on their correlation scores
10
11 columns_correlations = []
12 columns_non_numeric = []
13
14 for col_c in list(df_features.columns):
15     try:
16         corr = np.abs(df_features[col_c].corr(df_clean[col_target]))
17         t = (col_c, corr)
18         columns_correlations.append(t)
19     except:
20         columns_non_numeric.append(col_c)
21
22 # let's get all columns with correlation above 0.20
23
24 correlated_features_above_point2 = [col_c for col_c, corr in columns_correlations if corr >= 0.20]
25 print(f"correlated_features_above_point2 =")
26 print(correlated_features_above_point2)
27 print()
28
29 whole_space = correlated_features_above_point2
30 whole_space.append(col_target)
31
32 correlated_df = df_clean[whole_space]
33
34 correlated_df.head()
```

===== df\_features (as of 2020-10-31\_083203) =====

	bedrooms_zed15	bathrooms_zed15	floors_zed15	waterfront_zed15	metascore_zed15	sqft_living_zed15	building_age_zed15
0	-0.009333	-0.007667	-0.005000	0.000000	-0.016000	-8.614667	-0.160667
1	0.001333	0.008167	0.006667	0.000000	-0.008000	11.451333	0.058667
2	-0.013333	-0.012667	-0.007333	0.000000	-0.021333	-14.767333	0.579333
3	0.010000	0.006500	-0.005667	-0.000667	0.002000	-1.240000	0.050667
4	-0.002000	-0.003000	-0.006000	0.000000	-0.001333	-2.500000	-0.010000
...	...	...	...	...	...	...	...
21592	0.000000	0.003667	0.008667	0.000000	0.004000	0.733333	-0.286667
21593	0.006667	0.005500	0.006000	0.000000	0.007333	5.526667	-0.298000
21594	-0.006667	-0.012167	0.001333	0.000000	-0.002000	-3.500000	-0.351333
21595	-0.002667	-0.001333	0.000000	0.000000	-0.000667	-4.320000	0.004667
21596	-0.006667	-0.012167	0.001333	0.000000	-0.002000	-3.500000	-0.341333

21507 rows × 7 columns

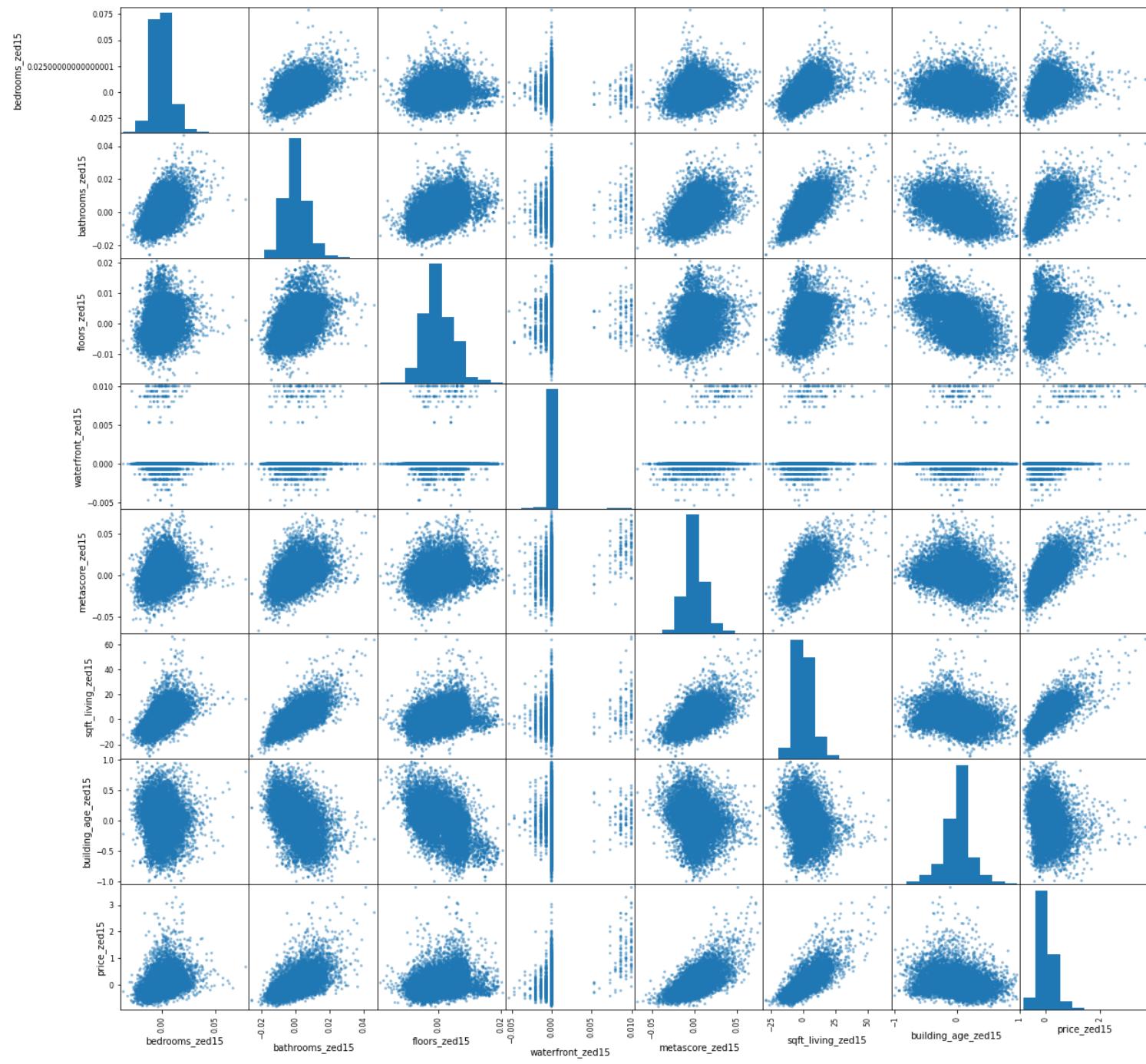
```
correlated_features_above_point2 =
['bedrooms_zed15', 'bathrooms_zed15', 'floors_zed15', 'waterfront_zed15', 'metascore_zed15', 'sqft_living_zed15', 'building_age_zed15']
```

Out[6]:

	bedrooms_zed15	bathrooms_zed15	floors_zed15	waterfront_zed15	metascore_zed15	sqft_living_zed15	building_age_zed15	price_zed15
0	-0.009333	-0.007667	-0.005000	0.000000	-0.016000	-8.614667	-0.160667	-0.420485
1	0.001333	0.008167	0.006667	0.000000	-0.008000	11.451333	0.058667	0.388877
2	-0.013333	-0.012667	-0.007333	0.000000	-0.021333	-14.767333	0.579333	-0.617753

	bedrooms_zed15	bathrooms_zed15	floors_zed15	waterfront_zed15	metascore_zed15	sqft_living_zed15	building_age_zed15	price_zed15
<b>3</b>	0.010000	0.006500	-0.005667	-0.000667	0.002000	-1.240000	0.050667	-0.046466
<b>4</b>	-0.002000	-0.003000	-0.006000	0.000000	-0.001333	-2.500000	-0.010000	0.006215

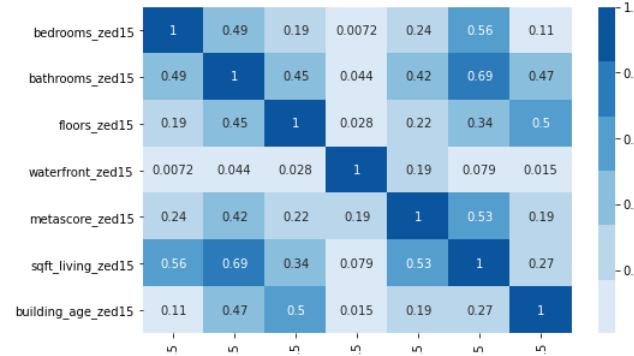
```
In [7]: 1 pd.plotting.scatter_matrix(correlated_df, figsize=(20, 20))
2 plt.show()
```



In [8]:

```
1 print(f"===== df_features CORRELATION PLOT =====")
2 plot_corr(df_features)
3 print()
4
5
6
7 for col_c in df_features.columns:
8     print(f"\n===== vif for column {repr(col_c)} =====")
9     vif = calculate_vif(df=df_features, target_col=col_c, show_res=True)
10    print(f"column {repr(col_c)} has vif of {vif}")
11    print(f"=====\\n")
```

===== df\_features CORRELATION PLOT =====



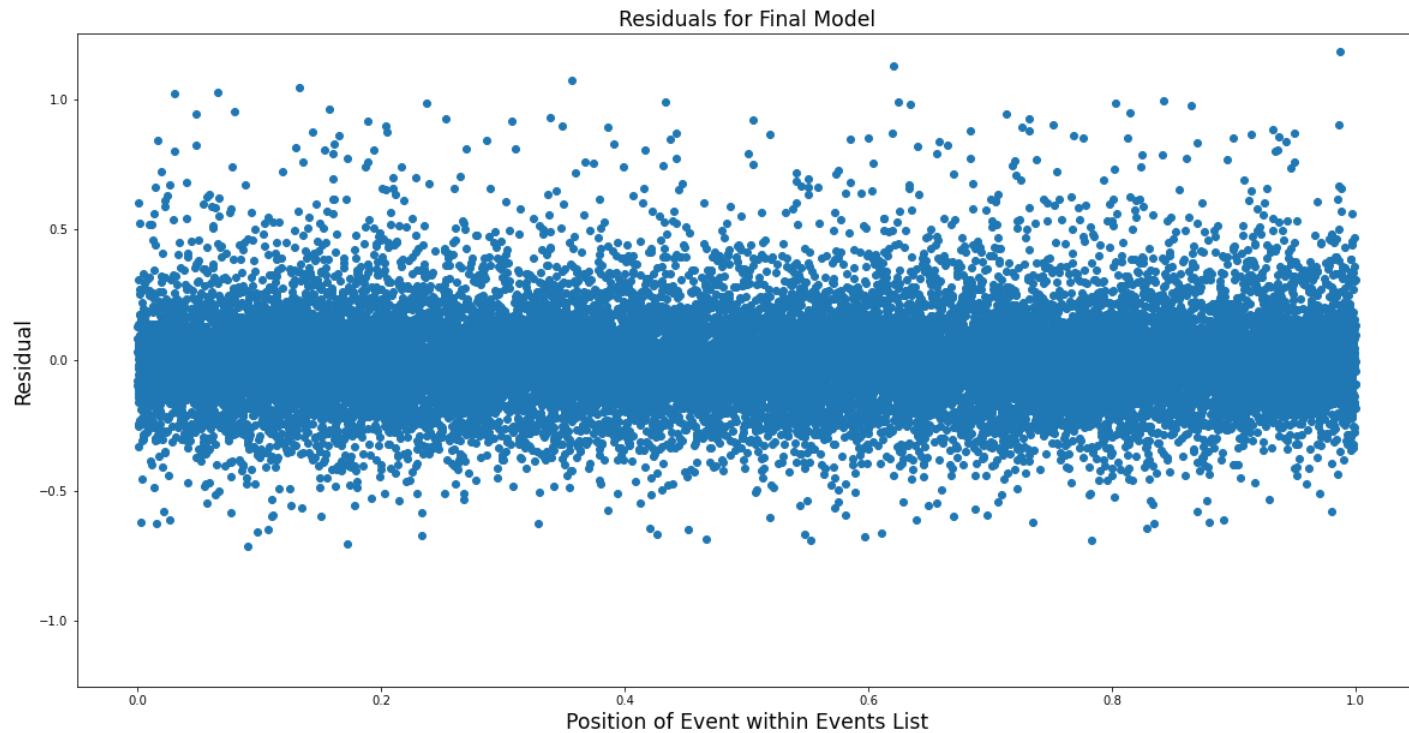
```
In [9]: 1 # =====
2 #          OLS Regression Model
3 #
4 # Michael Collins, Flatiron School
5 # 2020-10-30_0939_MDT
6 #
7 # =====
8
9 ols = sf.check_model(
10     df=df_clean,
11     features_to_use=cols_feature,
12     target_col=col_target,
13     add_constant=False,
14     show_summary=True,
15     vif_threshold=3.0)
```

```
OLS Regression Results
=====
Dep. Variable:      price_zed15   R-squared (uncentered):      0.683
Model:                          OLS   Adj. R-squared (uncentered):      0.683
Method:                     Least Squares   F-statistic:                 6611.
Date:                Sat, 31 Oct 2020   Prob (F-statistic):        0.00
Time:                  08:32:11   Log-Likelihood:            6716.8
No. Observations:      21507   AIC:                  -1.342e+04
Df Residuals:          21500   BIC:                  -1.336e+04
Df Model:                      7
Covariance Type:    nonrobust
=====
            coef    std. err      t    P>|t|      [0.025      0.975]
-----
bedrooms_zed15    -1.1833     0.175    -6.765    0.000    -1.526    -0.840
bathrooms_zed15    2.2330     0.283     7.899    0.000     1.679     2.787
floors_zed15       0.3517     0.315     1.118    0.264    -0.265     0.968
waterfront_zed15   86.2473    1.724    50.027    0.000    82.868    89.626
metascore_zed15     6.3519    0.118    53.798    0.000     6.120     6.583
sqft_living_zed15   0.0255    0.000    95.633    0.000     0.025     0.026
building_age_zed15  -0.0241    0.007    -3.639    0.000    -0.037    -0.011
=====
Omnibus:            3143.009   Durbin-Watson:             2.012
Prob(Omnibus):      0.000   Jarque-Bera (JB):         11168.543
Skew:                  0.722   Prob(JB):                  0.00
Kurtosis:                 6.222   Cond. No.           1.04e+04
=====
Warnings:
[1] Standard Errors assume that the covariance matrix of the errors is correctly specified.
[2] The condition number is large, 1.04e+04. This might indicate that there are
strong multicollinearity or other numerical problems.
Residuals failed test/tests
```

```
In [10]: 1 ols
```

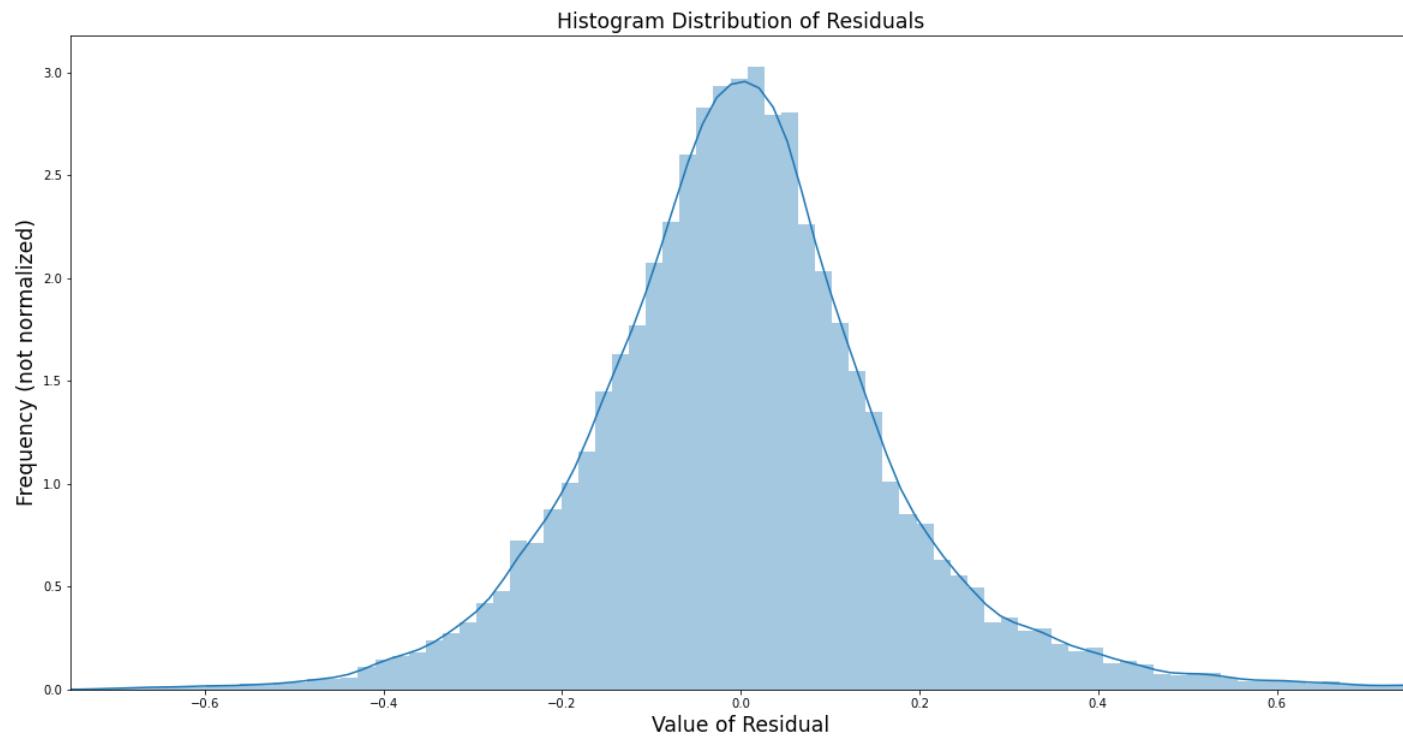
```
Out[10]: <statsmodels.regression.linear_model.RegressionResultsWrapper at 0x16dfa063c50>
```

```
In [11]: 1 # plot residuals
2 ols_resid = ols.resid
3 x_axis = np.linspace(0, 1, len(ols_resid))
4 plt.figure(figsize=(20,10))
5 plt.scatter(x_axis, ols_resid)
6 y_min = -1.25
7 y_max = 1.25
8 plt.ylim(y_min, y_max)
9 plt.ylabel("Residual", fontsize='xx-large')
10 plt.xlabel("Position of Event within Events List", fontsize='xx-large')
11 plt.title("Residuals for Final Model", fontsize='xx-large')
12 plt.show()
13
14
15 print(f"len(ols.resid) = {len(ols.resid)}")
16 ols_eventIndex = df_clean[col_eventIndex]
17 ols_resid = ols.resid
18 ols_eventIndex_resid = list(zip(ols_eventIndex, ols_resid))
19 outlier_eventIndices = []
20 for eventIndex, resid in ols_eventIndex_resid:
21     resid_min = -1
22     resid_max = 1
23     if not (resid_min <= resid <= resid_max):
24         outlier_eventIndices.append(eventIndex)
25 print("===== outlier eventIndices =====")
26 print(f"outlier_eventIndices = {repr(outlier_eventIndices)}\n")
```



```
len(ols.resid) = 21507
===== outlier eventIndices =====
outlier_eventIndices = [656, 1418, 2862, 7693, 13398, 21328]
```

```
In [12]: 1 plt.figure(figsize=(20,10))
2 sns.distplot(ols.resid, bins=100)
3 x_min = np.min(ols.resid)
4 x_max = np.max(ols.resid)
5 x_min = -0.75
6 x_max = 0.75
7 plt.xlim(x_min, x_max)
8 plt.title("Histogram Distribution of Residuals", fontsize='xx-large')
9 plt.ylabel("Frequency (not normalized)", fontsize='xx-large')
10 plt.xlabel("Value of Residual", fontsize='xx-large')
11
12 plt.show()
```



```
In [13]: 1 # How can we test that our residuals are normal?
2 # H0: data is normal
3 # HA: data is not normal
4 scs.shapiro(ols.resid)
5
6 # p = 0.87 -> fail to reject the null, therefore residuals are normally distributed
```

Out[13]: ShapiroResult(statistic=0.9581159353256226, pvalue=0.0)

```
In [14]: 1 df_map = df.loc[df[col_eventValid] == True]
2 display(df_map)
```

	eventIndex	eventValid	id	date	price	bedrooms	bathrooms	sqft_living	sqft_lot	floors	...	building_age_phi15	interior_age_mean15	interior_age_median15	interior_age_diff15	interior_age_
0	0	1	7129300520	2014-10-13	221900.0	3	1.00	1180	5650	1.0	...	-5.0	81.066667	80.0	-16.066667	
1	1	1	6414100192	2014-12-09	538000.0	3	2.25	2570	7242	2.0	...	0.0	63.133333	69.0	-34.133333	
2	2	1	5631500400	2015-02-25	180000.0	2	1.00	770	10000	1.0	...	5.0	29.066667	25.0	57.933333	
3	3	1	2487200875	2014-12-09	604000.0	4	3.00	1960	5000	1.0	...	-5.0	49.933333	65.0	5.066667	
4	4	1	1954400510	2015-02-18	510000.0	3	2.00	1680	8080	1.0	...	0.0	34.000000	33.0	-1.000000	
...	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...	
21592	21592	1	263000018	2014-05-21	360000.0	3	2.50	1530	1131	3.0	...	-5.0	34.733333	21.0	-23.733333	
21593	21593	1	6600060120	2015-02-23	400000.0	4	2.50	2310	5813	2.0	...	-5.0	35.800000	26.0	-29.800000	
21594	21594	1	1523300141	2014-06-23	402101.0	2	0.75	1020	1350	2.0	...	-5.0	46.133333	17.0	-35.133333	
21595	21595	1	291310100	2015-01-16	400000.0	3	2.50	1600	2388	2.0	...	5.0	15.533333	15.0	0.466667	
21596	21596	1	1523300157	2014-10-15	325000.0	2	0.75	1020	1076	2.0	...	-5.0	46.133333	17.0	-34.133333	

21507 rows × 163 columns

```
In [15]: 1 print(f"col_lat = {repr(col_lat)}")
2 print(f"col_long = {repr(col_long)}")
3 print()
4
5
6 col_price_mean = f"price_mean{ext_K}"
7 col_log10_price_mean = f"log10_price_mean{ext_K}"
8 map_price_mean = df_map[col_price_mean]
9 map_log10_price_mean = [np.log10(pm) for pm in map_price_mean]
10 df_map[col_log10_price_mean] = map_log10_price_mean
11 df_map['resid'] = ols.resid
12 df_map['resid_qcut'] = pd.qcut(df_map['resid'], 10, labels=False)
13 display(df_map)

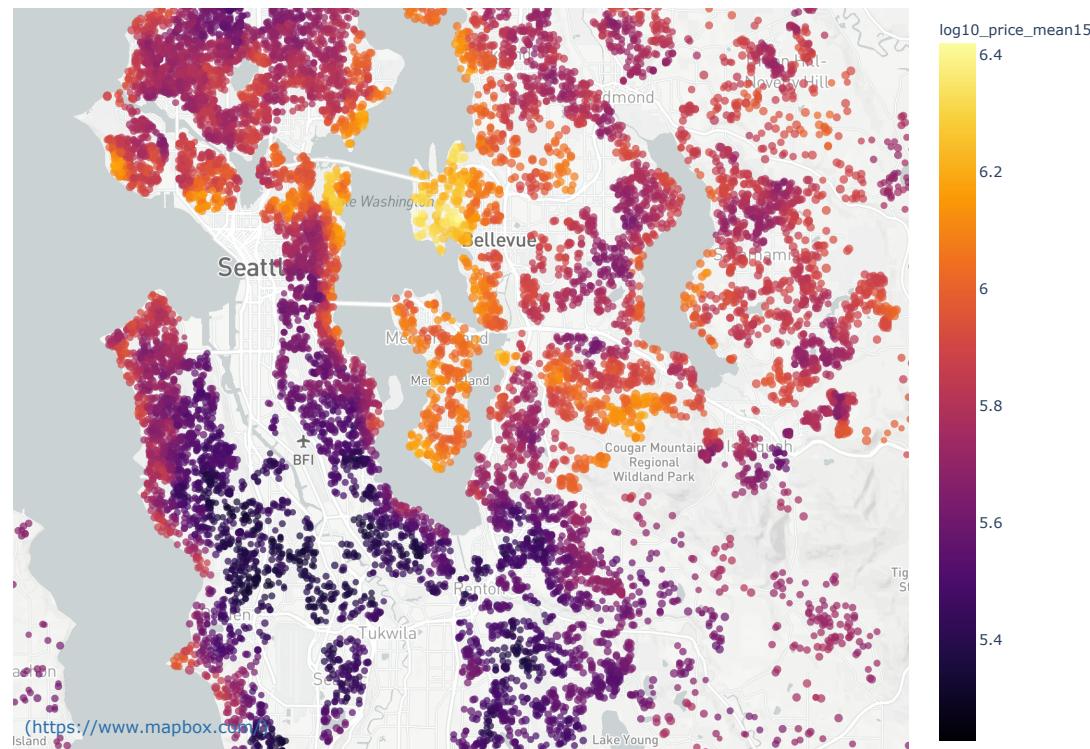
col_lat = 'lat'
col_long = 'long'
```

	eventIndex	eventValid	id	date	price	bedrooms	bathrooms	sqft_living	sqft_lot	floors	...	interior_age_diff15	interior_age_delta15	interior_age_sigma15	interior_age_MAD15	interior_age_z
0	0	1	7129300520	2014-10-13	221900.0	3	1.00	1180	5650	1.0	...	-16.066667	-15.0	27.805195	22.0	-
1	1	1	6414100192	2014-12-09	538000.0	3	2.25	2570	7242	2.0	...	-34.133333	-40.0	15.244088	2.0	-
2	2	1	5631500400	2015-02-25	180000.0	2	1.00	770	10000	1.0	...	57.933333	62.0	14.073458	8.0	-
3	3	1	2487200875	2014-12-09	604000.0	4	3.00	1960	5000	1.0	...	5.066667	-10.0	33.235456	30.0	-
4	4	1	1954400510	2015-02-18	510000.0	3	2.00	1680	8080	1.0	...	-1.000000	0.0	3.425395	1.0	-
...	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...	-
21592	21592	1	263000018	2014-05-21	360000.0	3	2.50	1530	1131	3.0	...	-23.733333	-10.0	31.890368	9.0	-
21593	21593	1	6600060120	2015-02-23	400000.0	4	2.50	2310	5813	2.0	...	-29.800000	-20.0	25.771302	20.0	-
21594	21594	1	1523300141	2014-06-23	402101.0	2	0.75	1020	1350	2.0	...	-35.133333	-6.0	46.173393	9.0	-
21595	21595	1	291310100	2015-01-16	400000.0	3	2.50	1600	2388	2.0	...	0.466667	1.0	1.087300	0.0	-
21596	21596	1	1523300157	2014-10-15	325000.0	2	0.75	1020	1076	2.0	...	-34.133333	-5.0	46.173393	9.0	-

21507 rows × 166 columns

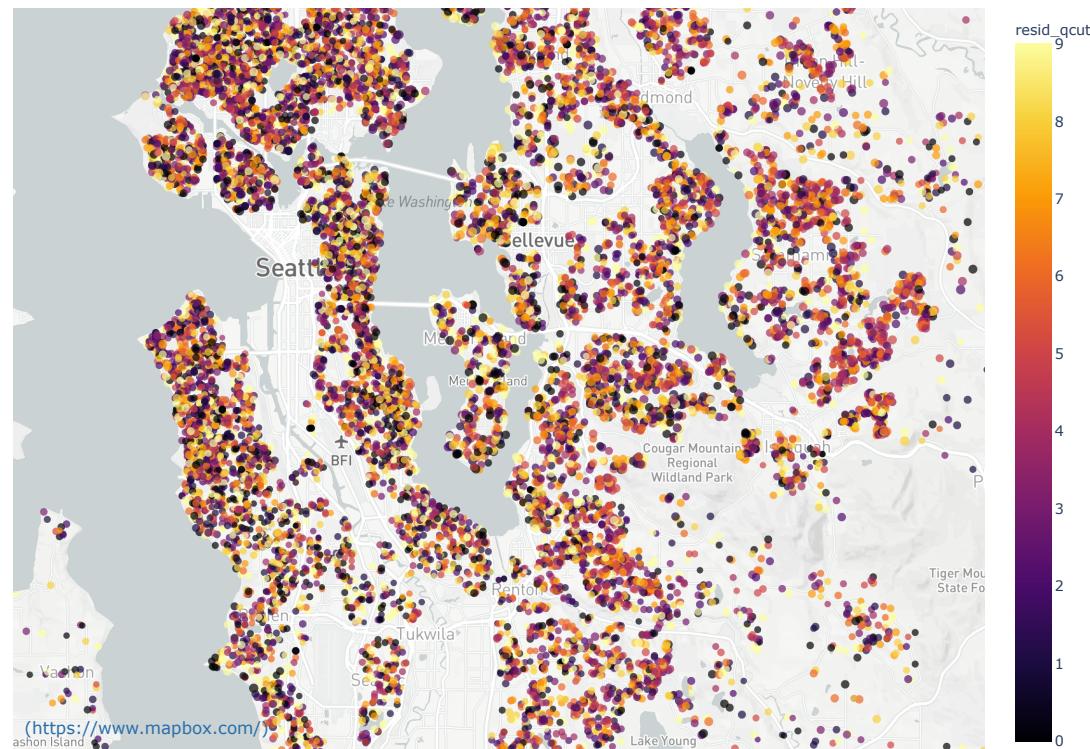
```
In [16]: 1 # =====
2 #   Make a pretty map of selected data from the data frame.
3 # =====
4
5 mapbox_access_token = open(r"c:\.mapbox\access_token").read()
6
7 px.set_mapbox_access_token(mapbox_access_token)
8
9 fig = px.scatter_mapbox(df_map,
10                         lat=col_lat,
11                         lon=col_long,
12                         color=col_log10_price_mean,
13                         size=col_grade,
14                         color_continuous_scale="Inferno",
15                         size_max=6,
16                         zoom=10, title=f"log10(Mean Price of {ext_K} nearest Neighbors)")
17 fig.update_layout(autosize=False,
18                   width=1000,
19                   height=800)
20
21 fig.show()
```

log10(Mean Price of 15 nearest Neighbors)



```
In [17]: 1 # =====
2 #   Make a pretty map of selected data from the data frame.
3 # =====
4
5 mapbox_access_token = open(r"c:\.mapbox\access_token").read()
6
7 px.set_mapbox_access_token(mapbox_access_token)
8
9 fig = px.scatter_mapbox(df_map,
10                         lat=col_lat,
11                         lon=col_long,
12                         color='resid_qcut',
13                         size=col_grade,
14                         color_continuous_scale="Inferno",
15                         size_max=6,
16                         zoom=10, title="residual deciles by location")
17 fig.update_layout(autosize=False,
18                   width=1000,
19                   height=800)
20 fig.show()
```

residual deciles by location



```
In [18]: 1 print(df_raw.columns)
```

```
Index(['id', 'date', 'price', 'bedrooms', 'bathrooms', 'sqft_living',
       'sqft_lot', 'floors', 'waterfront', 'view', 'condition', 'grade',
       'sqft_above', 'sqft_basement', 'yr_builtin', 'yr_renovated', 'zipcode',
       'lat', 'long', 'sqft_living15', 'sqft_lot15', 'KNN', 'building_age',
       'interior_age', 'log10_price'],
      dtype='object')
```

```
In [19]: 1 # =====
2 #     Take a detailed look at the raw values
3 #     in columns whose values may be categorical
4 #
5 #     Michael Collins, Flatiron School
6 #     2020-10-30_0455_MDT
7 # =====
8
9 cols_examine = list(df_raw.columns)
10 cols_summarized = []
11 quanta_all = []
12 for col_c in cols_examine:
13     items_c = sorted(list(dict(df_raw[col_c].value_counts()).items()))
14     num_c = len(items_c)
15     if num_c < 50:
16         cols_summarized.append(col_c)
17         print(f"==== {repr(col_c)} ====")
18         quanta_c = [repr(k) for k, v in items_c]
19         quanta_all.append(f"quanta_{col_c} = {quanta_c}")
20         for k, v in items_c:
21             s_k = str(k) + " " * (8-len(str(k)))
22             print(s_k, v)
23         print()
24
25 print("==== COLUMNS_SUMMARIZED ====")
26 print(f"cols_summarized = {repr(cols_summarized)}\n")
27
28 print("==== QUANTA_SUMMARIZED ====")
29 print(f"quanta_summarized =")
30 for s in quanta_all:
31     print(s)
32 print()
```

===== 'bedrooms' =====

```
1      196
2      2760
3      9824
4      6882
5      1601
6      272
7      38
8      13
9      6
10     3
11     1
33     1
```

===== 'bathrooms' =====

```
0.5      4
0.75     71
1.0      3851
1.25     9
1.5      1445
1.75     3048
2.0      1930
2.25     2047
2.5      5377
2.75     1185
3.0      753
3.25     589
3.5      731
3.75     155
4.0      136
4.25     79
4.5      100
4.75     23
5.0      21
5.25     13
```

```

5.5      10
5.75     4
6.0       6
6.25     2
6.5       2
6.75     2
7.5       1
7.75     1
8.0       2

===== 'floors' =====
1.0      10673
1.5      1910
2.0      8235
2.5      161
3.0      611
3.5      7

===== 'waterfront' =====
0         21451
1         146

===== 'view' =====
0         19485
1         330
2         957
3         508
4         317

===== 'condition' =====
1         29
2         170
3        14020
4         5677
5        1701

===== 'grade' =====
3         1
4         27
5         242
6        2038
7        8974
8        6065
9        2615
10        1134
11        399
12        89
13        13

===== COLUMNS_SUMMARIZED =====
cols_summarized = ['bedrooms', 'bathrooms', 'floors', 'waterfront', 'view', 'condition', 'grade']

===== QUANTA_SUMMARIZED =====
quanta_summarized =
quanta_bedrooms = ['1', '2', '3', '4', '5', '6', '7', '8', '9', '10', '11', '33']
quanta_bathrooms = ['0.5', '0.75', '1.0', '1.25', '1.5', '1.75', '2.0', '2.25', '2.5', '2.75', '3.0', '3.25', '3.5', '3.75', '4.0', '4.25', '4.5', '4.75', '5.0', '5.25', '5.5', '6.75', '7.5', '7.75', '8.0']
quanta_floors = ['1.0', '1.5', '2.0', '2.5', '3.0', '3.5']
quanta_waterfront = ['0', '1']
quanta_view = ['0', '1', '2', '3', '4']
quanta_condition = ['1', '2', '3', '4', '5']
quanta_grade = ['3', '4', '5', '6', '7', '8', '9', '10', '11', '12', '13']

```

In [ ]:

1

