

AnalysisNote8

November 15, 2022

1 Real distributions vs Permutation

As we have seen, most tricksters exist where corresponding real animals are observed. However, because of cultural transmission, distributions are more dense in realized tricksters distributions than a null hypothesis (trickster randomly distributed where real animals are observed).

To avoid the problems of difference in research efforts, we only consider presence/absence in each hex grid.

```
[74]: import numpy as np
import seaborn as sns
import pandas as pd
import matplotlib.pyplot as plt
import geopandas as gpd
import sys
#import rioxtarray as rxr we used this function to read tiff, but could cause
↳ conflict with plotting hex grids
import h3
from geopy.distance import geodesic
#from geojson import Feature, Point, FeatureCollection, Polygon
import plotly.express as px
import random
import scipy as sp
from scipy.integrate import cumtrapz
from statsmodels.stats import multitest
import math
import statannot
df=pd.read_csv('TrickSter_data3.csv')
df
```

```
[74]:
```

	nid	lat	lng	TrickSter	Annu_Mean_Temp	Annu_Prec	\
0	103918	68.72000	158.70000	raven	-13.231120	173.0	
1	103817	64.95187	64.64844	raven	-4.123444	473.0	
2	3926	57.82907	-152.98389	raven	3.271220	1499.0	
3	2434	51.00000	-112.50000	raven	3.957418	342.0	
4	2439	40.00000	-103.00000	raven	9.747086	395.0	
..	
512	2464	52.00000	-120.00000	wren	1.460037	964.0	

513	2466	50.50000	-122.80000	wren	4.090998	822.0
514	2472	47.11975	-123.53846	wren	9.940903	1873.0
515	2477	45.50000	-124.00000	wren	10.183416	2005.0
516	2220	51.50000	14.40000	wren	9.163620	593.0

	hex_index	presence	Norm_Annu_Mean_Temp	Norm_Annu_Prec
0	8104ffffffffffff	1	0.050405	0.025701
1	81107ffffffffffff	1	0.249975	0.072716
2	810c7ffffffffffff	1	0.412010	0.233506
3	8112ffffffffffff	1	0.427046	0.052186
4	8126bffffffffffff	1	0.553911	0.060492
..
512	8112ffffffffffff	1	0.372323	0.149663
513	8128ffffffffffff	1	0.429973	0.127409
514	8128ffffffffffff	1	0.558158	0.292117
515	8128ffffffffffff	1	0.563472	0.312804
516	811f3ffffffffffff	1	0.541126	0.091522

[517 rows x 10 columns]

```
[78]: def MedianDistance(data):
    # calculate median distance given data points
    center_lng=[]
    center_lat=[]
    for i in range(len(data)):
        x=h3.h3_to_geo_boundary(data[i], True)
        x=np.mean(x, axis=0)
        center_lng.append(x[0])
        center_lat.append(x[1])
    df=pd.DataFrame({'latitude':center_lat, 'longitude':center_lng})
    Distance=[] # realized distances between tricksters
    for i in range(len(df)):
        place1=df.iloc[i, :]
        for j in range(i+1, len(df)):
            place2=df.iloc[j, :]
            Distance.append(geodesic(place1, place2).km)
    return np.median(Distance)

def Distribution_Test(df, target, rep=10000):
    # Test whether TS is more dense than random distributions on RA
    # df: data frame of tricksters. See above as an example
    # target: str of target species
    # rep: int of replications to generate a median distributions under the
    # null hypothesis
    if target == 'water bird' or target == "monkey" or target == '
    ground_squirrel':
        print ("We ignore this species")
```

```

        return np.nan
    else:
        □
    ↪#-----
        # Step 1: calculate distances between tricksters
        hex_TS=np.unique(df[df['TrickSter']==target]['hex_index']) # hex_grids_□
    ↪of focal Tricksters
        median_TS=MedianDistance(hex_TS)
        □
    ↪#-----
        # Step 2: calculate distances under null hypothesis (random_□
    ↪distirbutions of TS given RA exist)
        df_meta=pd.read_csv('./GBIF/For_gbif_trickstar.xlsx - Sheet1.csv') #□
    ↪meta file
        taxa=df_meta[df_meta['Category']==target]['Taxa'].reset_index(drop=True)
        for i in range(len(taxa)):
            if i==0:
                data=pd.read_csv('./GBIF/'+target+'/' +taxa[i]+'_cleaned.csv')
            else:
                dd=pd.read_csv('./GBIF/'+target+'/' +taxa[i]+'_cleaned.csv')
                data=pd.concat([data,dd])
        #print(data)
        hex_RA=np.unique(data['hex_index'])
        Median_RA=[]
        for k in range(rep):
            hex_extract=np.random.choice(hex_RA, len(hex_TS), replace=False)
            Median_RA.append(MedianDistance(hex_extract))
        kernel=sp.stats.gaussian_kde(Median_RA)
        x=np.linspace(0, max(Median_RA), 5000)
        y=kernel(x)
        cum_y= cumtrapz(y, x)
        idx_d= np.searchsorted(cum_y, 0.05) # index of x that gives cum~0.05
        p_val= cum_y[np.searchsorted(x, median_TS, side='right')]
        plt.plot(x, y, color='#66c2a5')
        plt.vlines(x=median_TS, ymin=0, ymax=max(y)*0.8, color='k',□
    ↪linestyle='--')
        plt.ylim(0, max(y)*1.05)
        plt.fill_between(x[:idx_d], np.zeros([np.size(x[:idx_d])]), y[:idx_d],□
    ↪color='#66c2a5')
        plt.xlabel("Median distance (km)", fontsize=20)
        plt.xticks([0, 4000, 8000, 12000], fontsize=16)
        plt.ylabel("Probability", fontsize=20)
        plt.yticks(fontsize=16)
        plt.text(x=median_TS*0.9, y=0.8*(max(y)), s='Median of tricksters',□
    ↪fontsize=16)
        plt.title(target, fontsize=20)

```

```

plt.savefig("CompareDistance_Earth"+target+".pdf", bbox_inches='tight',
pad_inches=0.05)
plt.show()
return p_val # significance

```

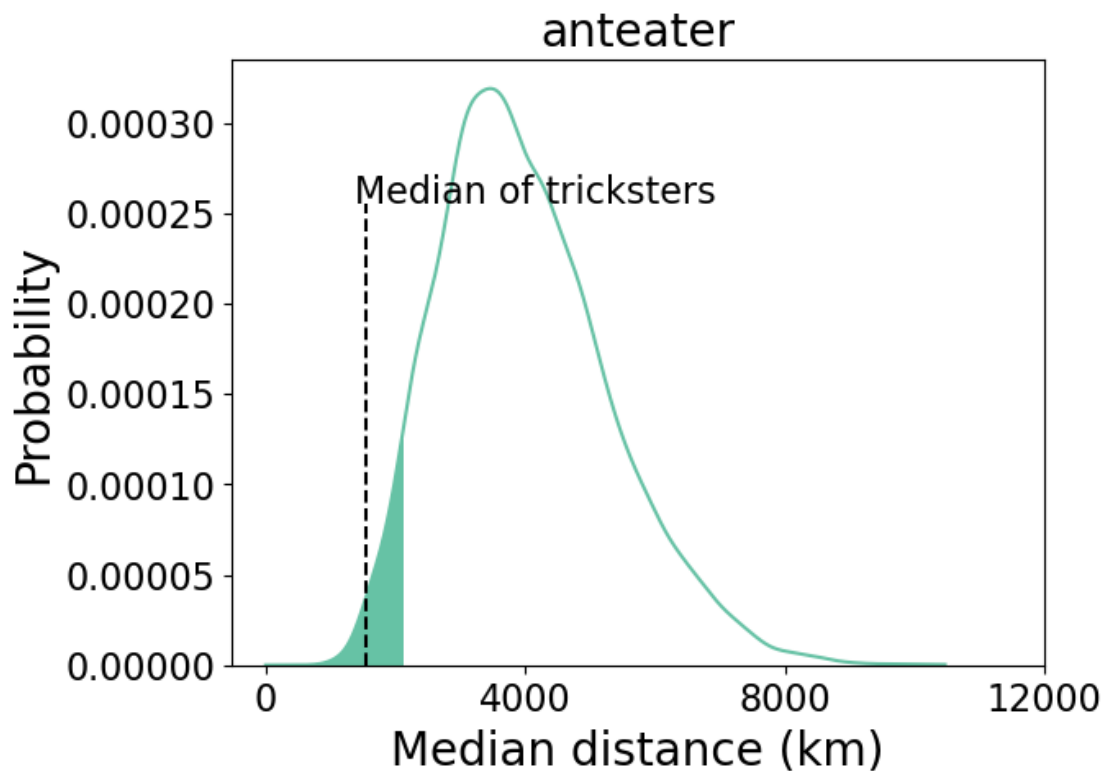
```

[79]: species=np.sort(np.unique(df["Trickster"]))
P_vals=[]
for i in range(len(species)):
    target=species[i]
    print(i)
    P_vals.append(Distribution_Test(df, target, rep=10000))

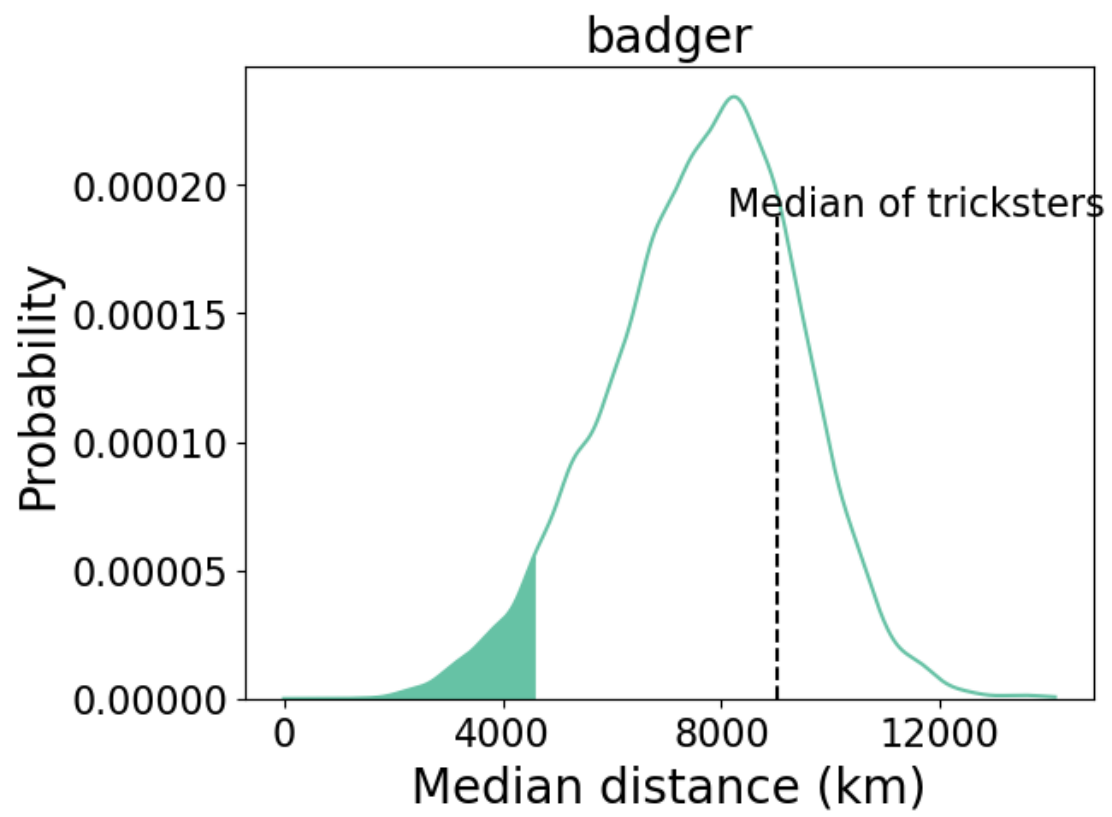
df_result=pd.DataFrame({'Tricksters':species, 'P values':P_vals})
df_result

```

0



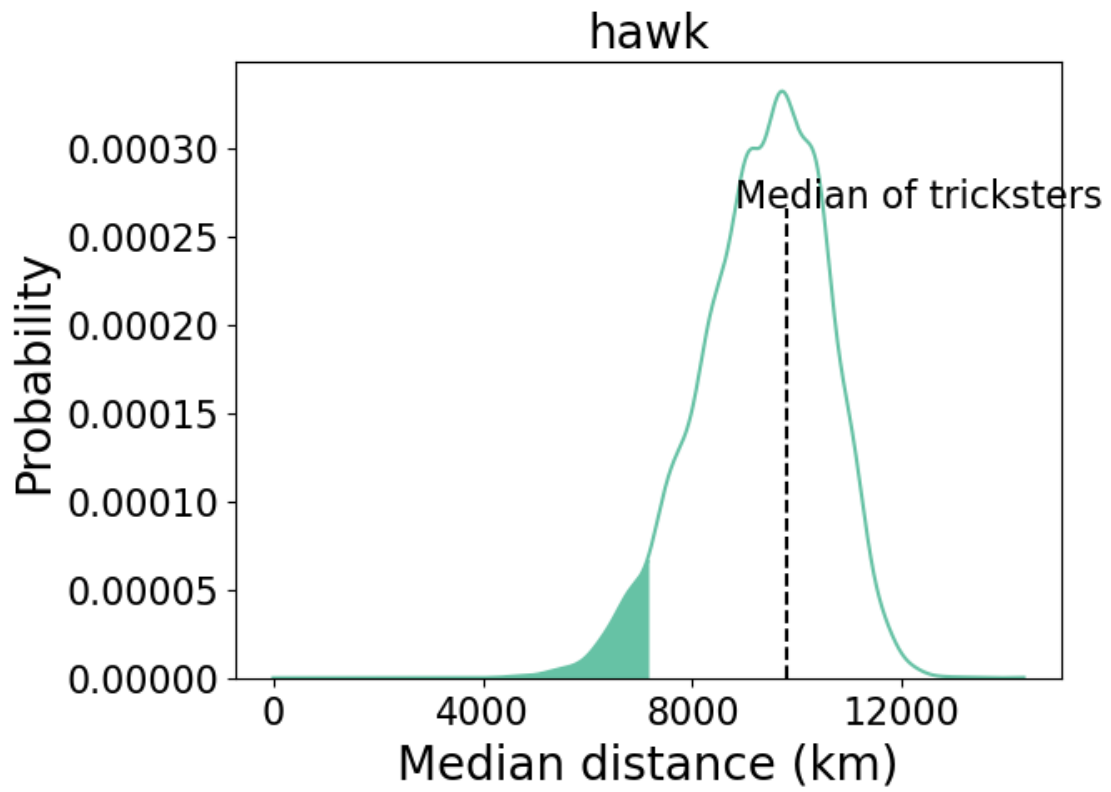
1

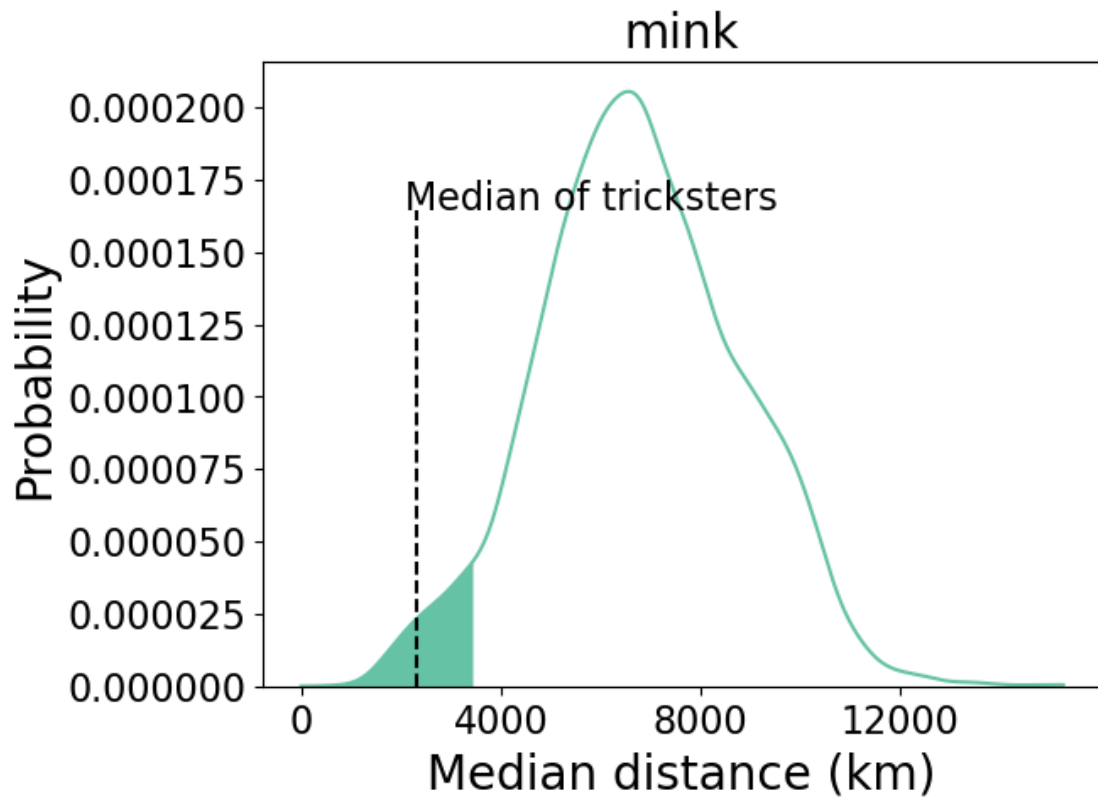


2

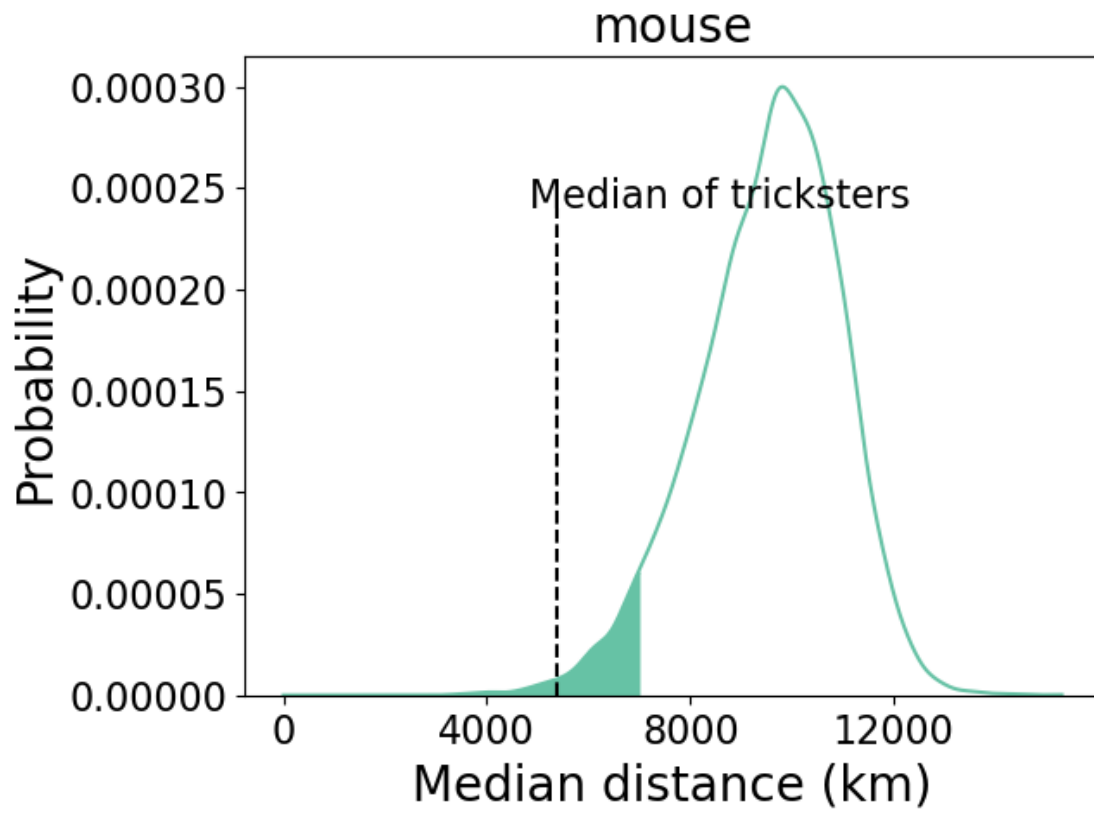
We ignore this species

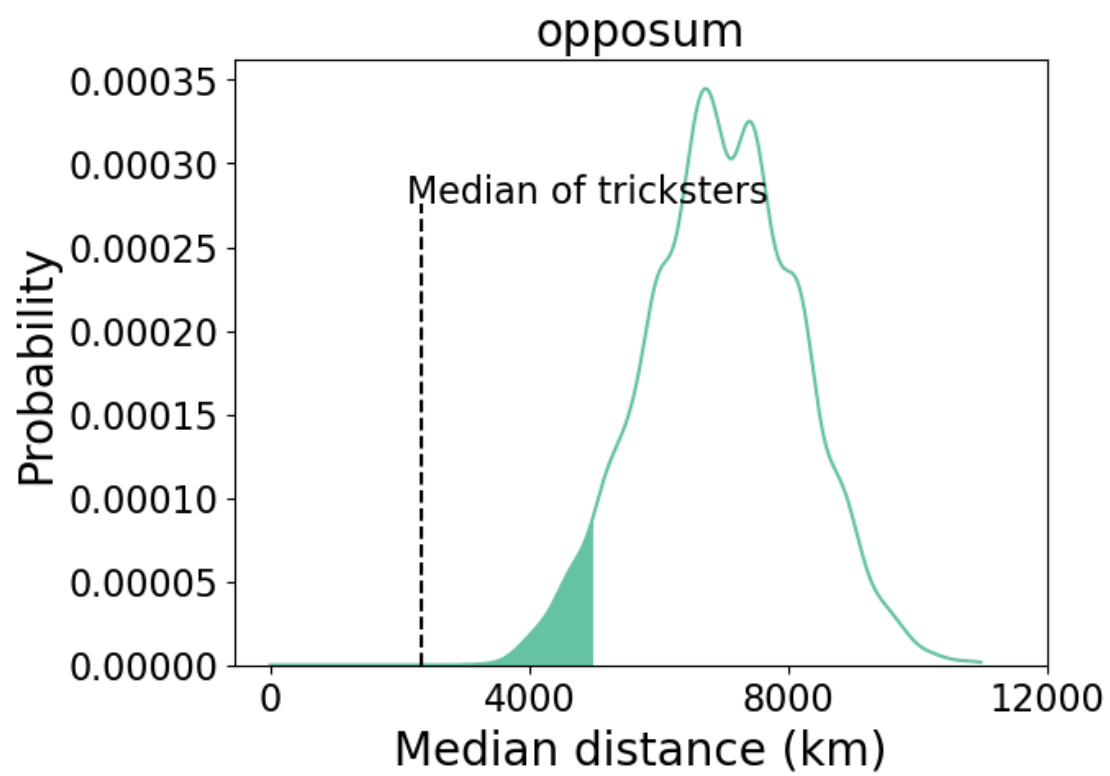
3

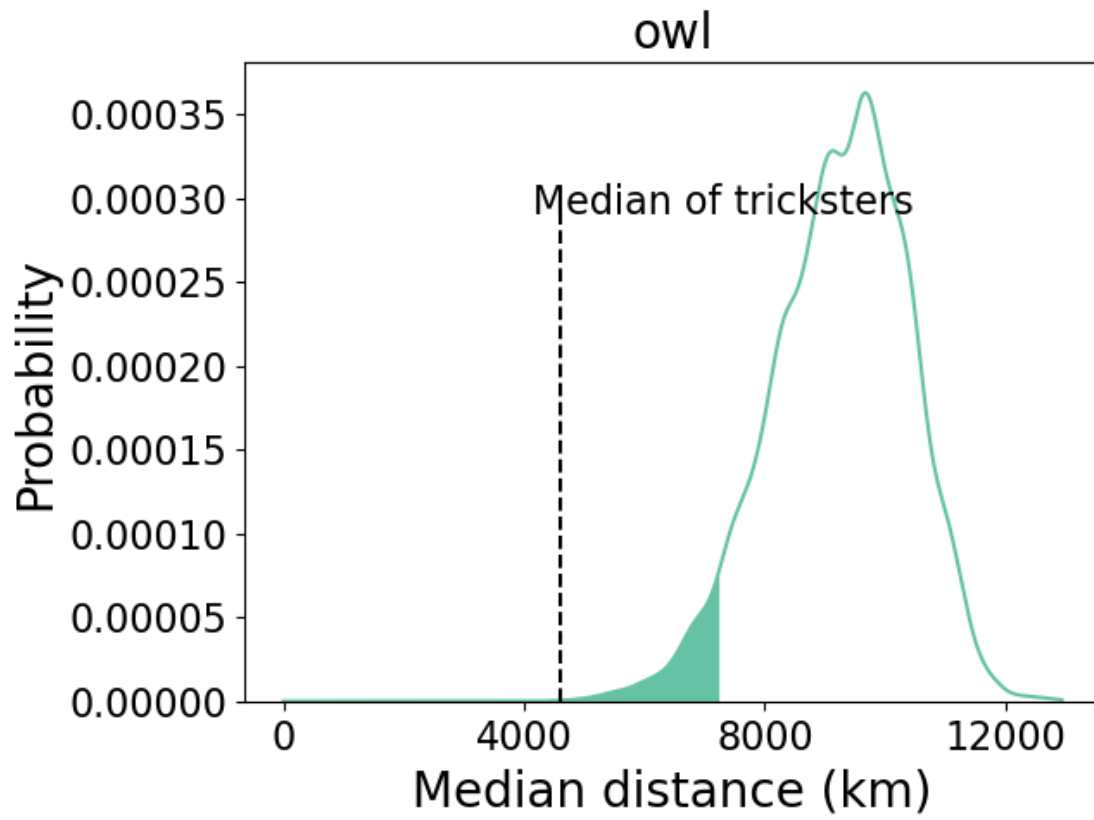


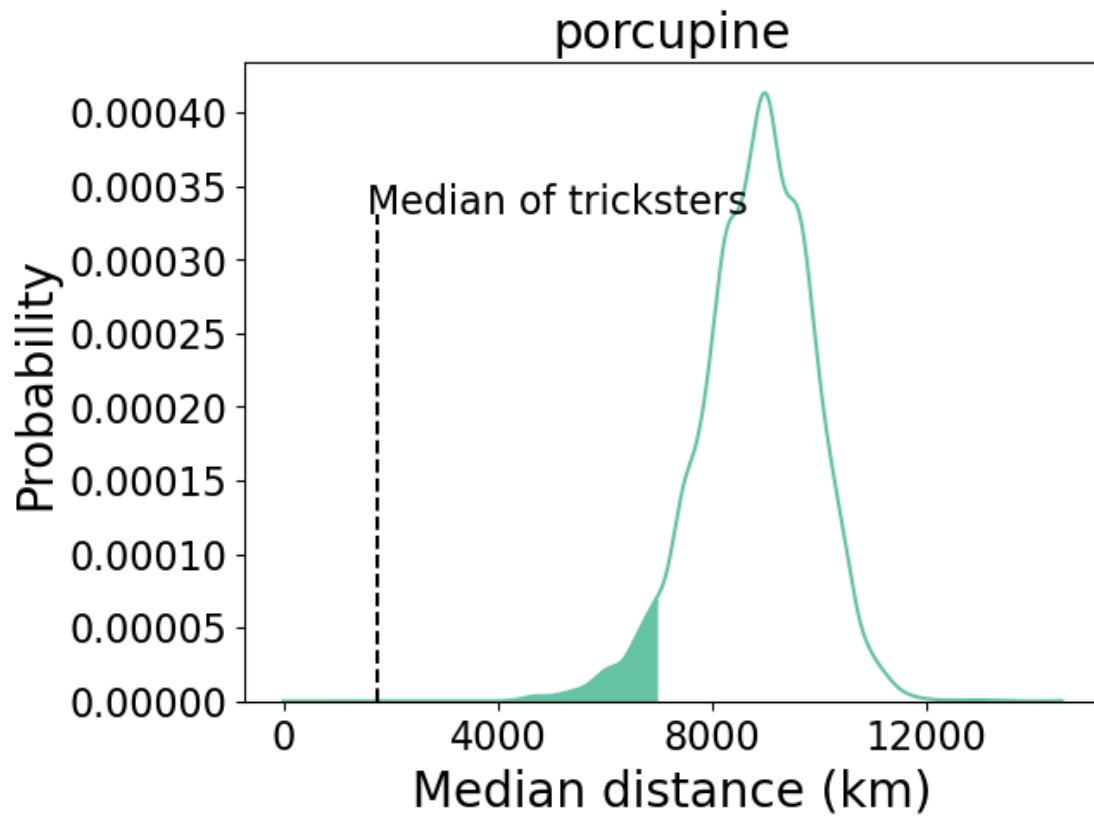


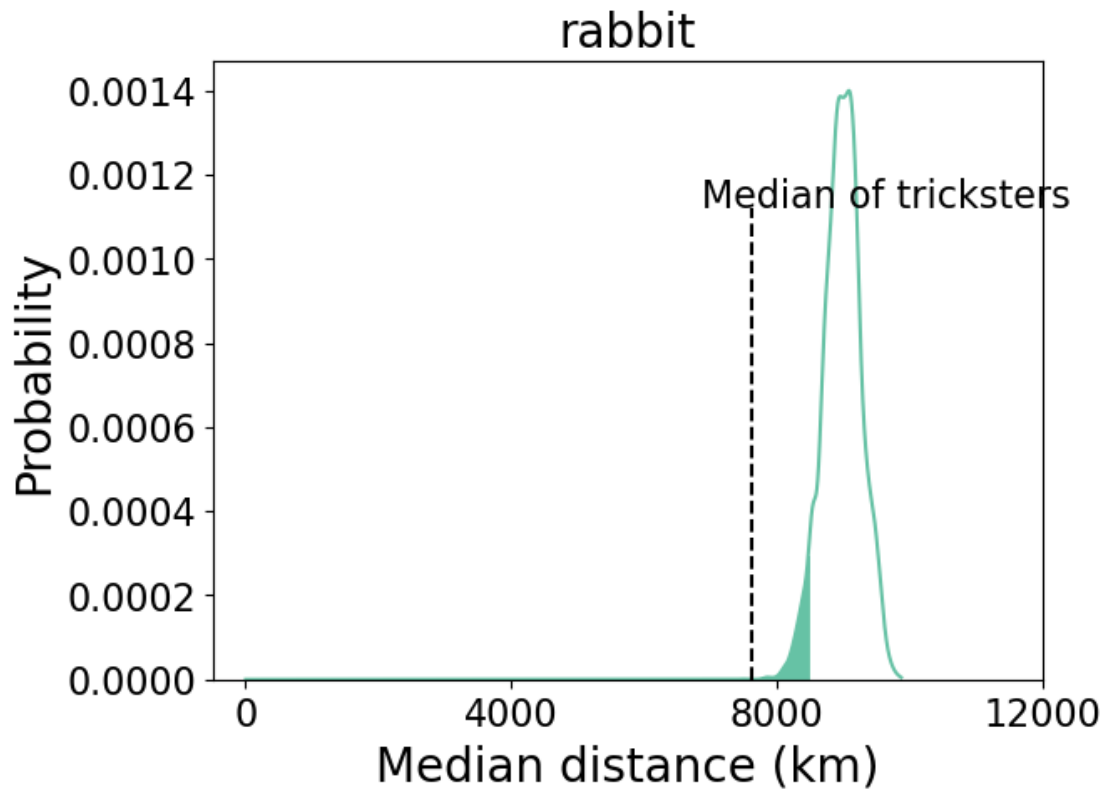
5
We ignore this species
6

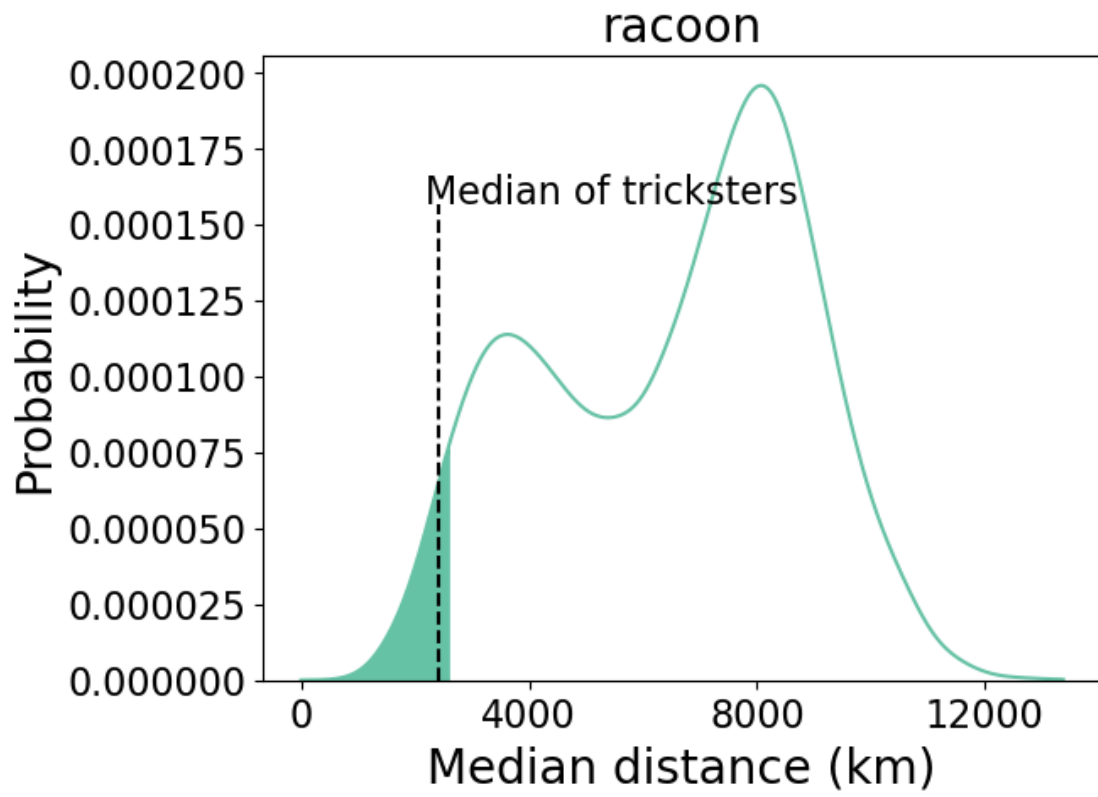


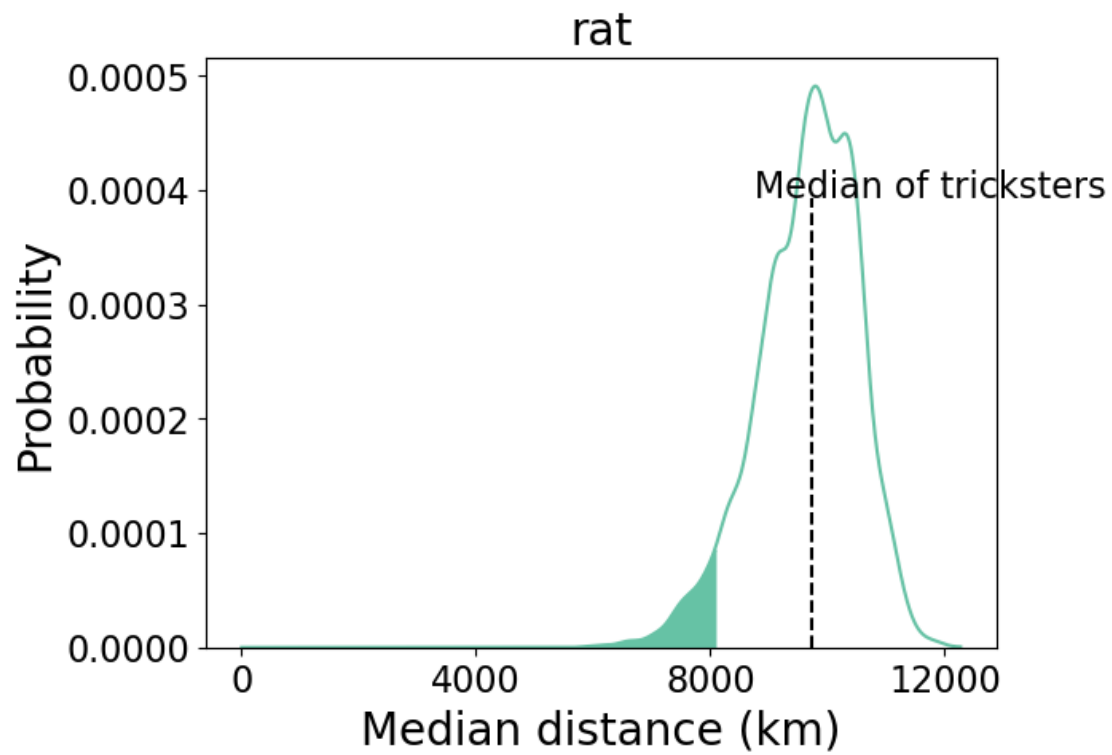


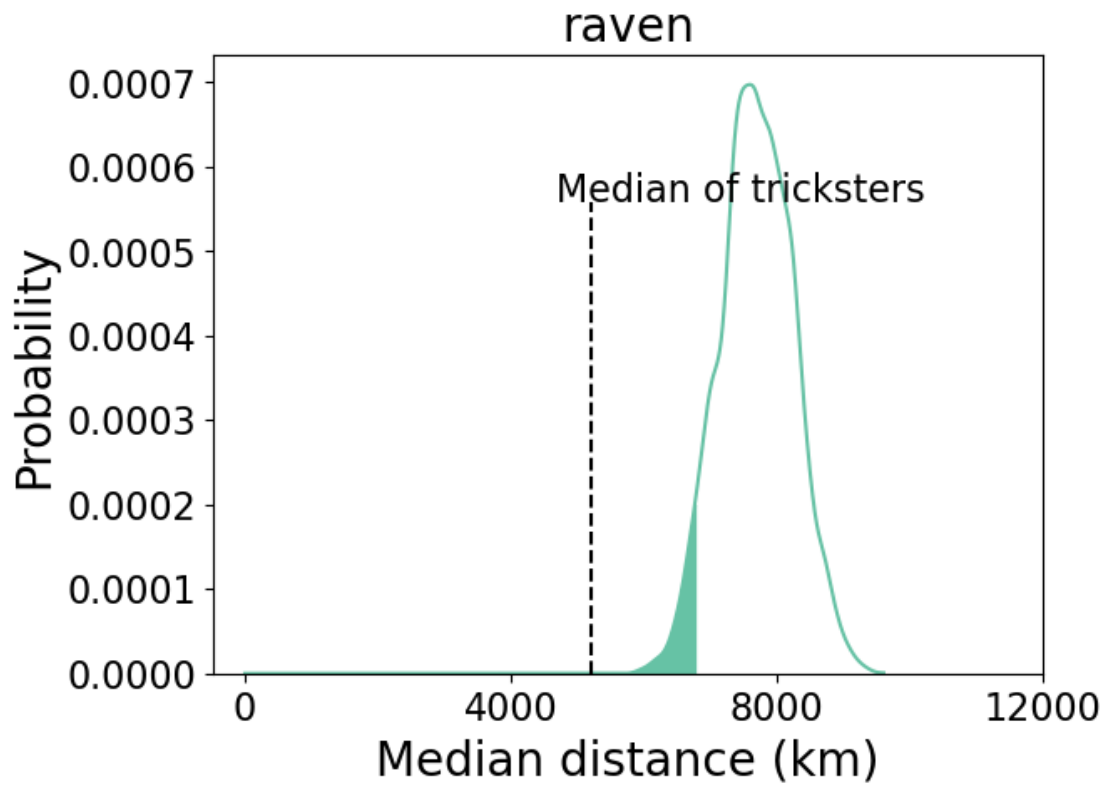


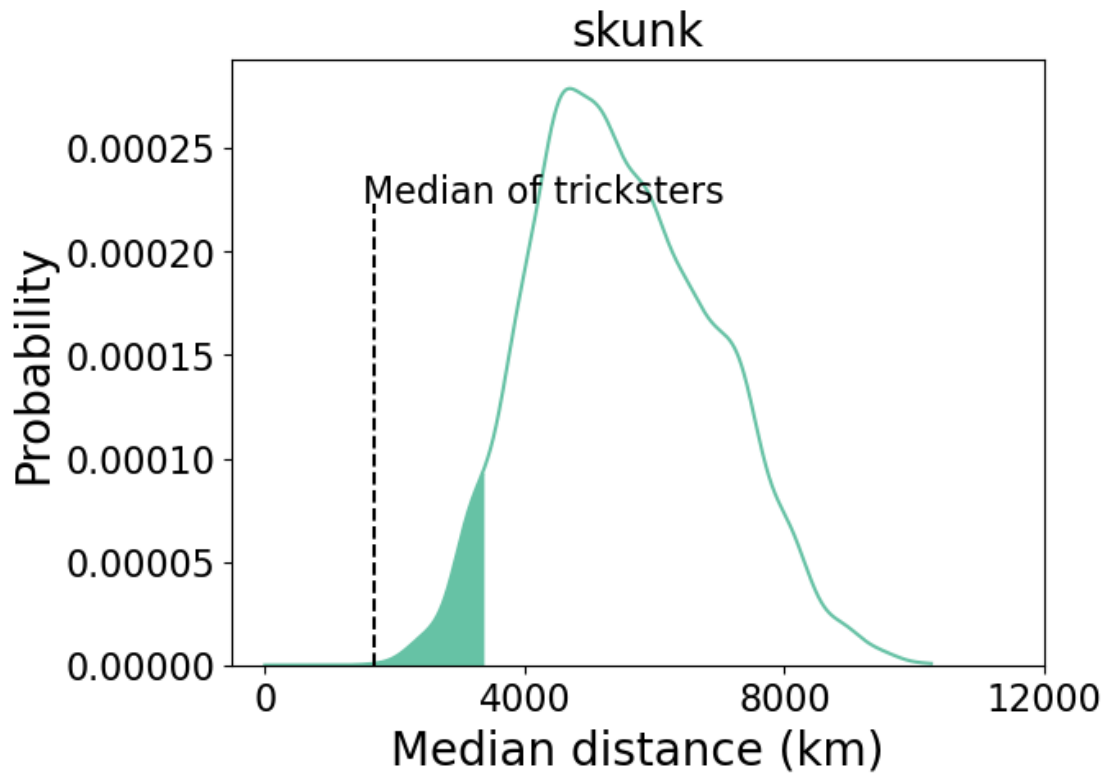


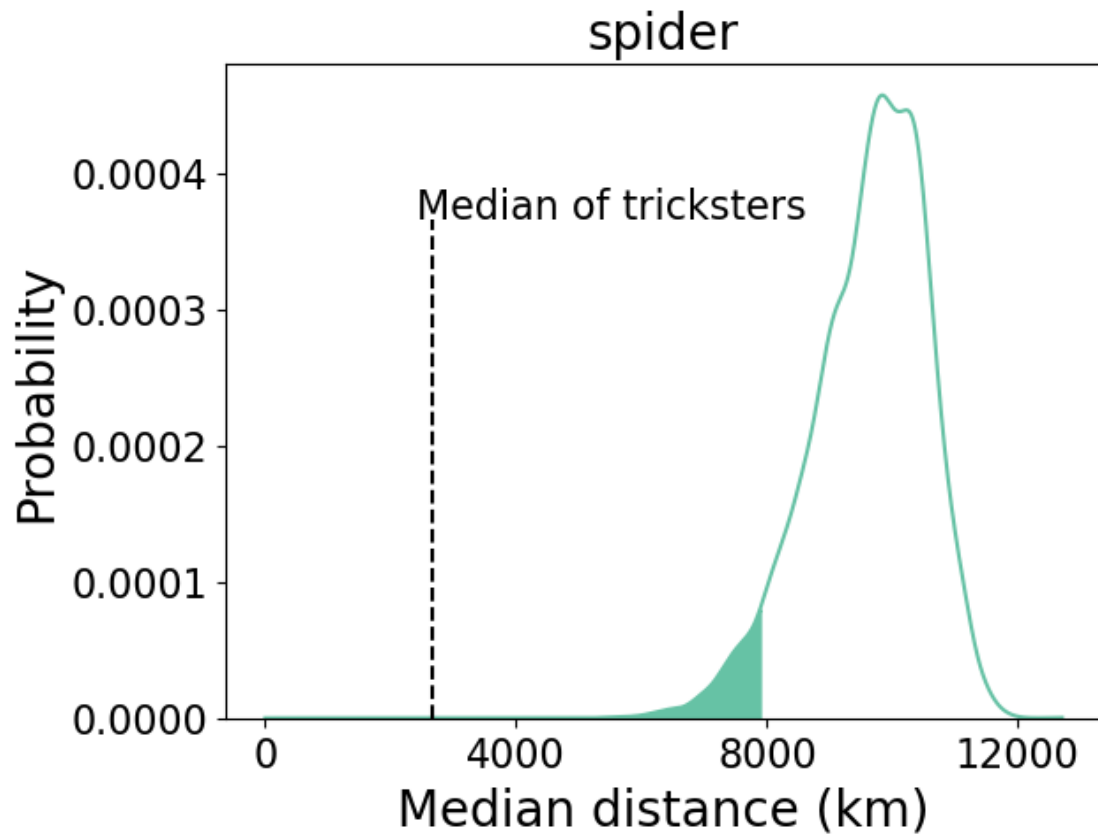








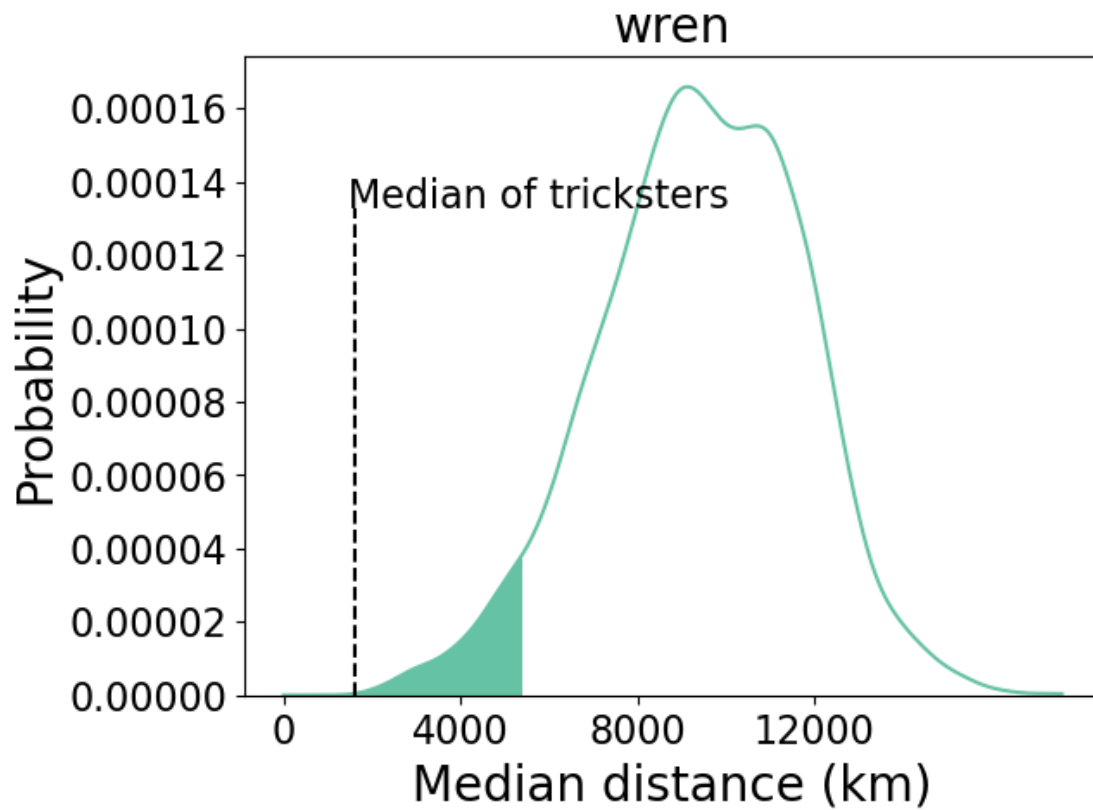




16

We ignore this species

17



```
[79]:
```

	Tricksters	P values
0	anteater	8.103057e-03
1	badger	7.761161e-01
2	ground_squirrel	NaN
3	hawk	6.139158e-01
4	mink	1.487048e-02
5	monkey	NaN
6	mouse	5.394848e-03
7	opposum	7.457305e-10
8	owl	8.808659e-05
9	porcupine	4.552271e-56
10	rabbit	9.994088e-05
11	raccoon	3.886512e-02
12	rat	4.906571e-01
13	raven	5.871440e-21
14	skunk	1.285825e-04
15	spider	1.480809e-85
16	water bird	NaN
17	wren	8.685610e-05

```
[82]: a, b=multitest.fdr correction(df_result['P values'])
df_result['FDR']=a
df_result
```

```
[82]:
```

	Tricksters	P values	FDR
0	anteater	8.103057e-03	True
1	badger	7.761161e-01	False
2	ground_squirrel	NaN	False
3	hawk	6.139158e-01	False
4	mink	1.487048e-02	True
5	monkey	NaN	False
6	mouse	5.394848e-03	True
7	opposum	7.457305e-10	True
8	owl	8.808659e-05	True
9	porcupine	4.552271e-56	True
10	rabbit	9.994088e-05	True
11	raccoon	3.886512e-02	False
12	rat	4.906571e-01	False
13	raven	5.871440e-21	True
14	skunk	1.285825e-04	True
15	spider	1.480809e-85	True
16	water bird	NaN	False
17	wren	8.685610e-05	True

```
[83]: 11/15
```

```
[83]: 0.7333333333333333
```

As we can see above, about 70% data (where we remove three species from the analysis) suggest that TS is more clugged than the null hypothesis. The exceptions are badger, hawk, racoonm, and rat. RA distribution of racoon is clugged while that of badger hawk, and rat are as broad as those of TS, respectively (see Analysis Note 6). In the above figures, these three species show median >8000km

```
[ ]:
```