

In [1]:

```
1 import pandas as pd
2 import numpy as np
3 import matplotlib.pyplot as plt
4
5 from sklearn import ensemble
6 from sklearn.model_selection import train_test_split, GridSearchCV, PredefinedS
7 from sklearn.metrics import mean_squared_error, make_scorer
8
9 from xgboost import XGBRegressor
10 from xgboost import plot_tree
```

Data Prep

Repositioning Column "INFECTION RATE" to the end of the dataset

In [2]:

```
1 data = pd.read_csv('final_data_cleaned.csv')
2 data['INFECTION RATE'].isnull().any()
3 # no null value in 'INFECTION RATE' --> treat 0.0 as missing values
```

Out[2]:

False

In [3]:

```
1 infection_rate = data['INFECTION RATE']
2 infection_rate
```

Out[3]:

```
0      0.034420
1      0.031714
2      0.034908
3      0.057709
4      0.020817
...
3142    0.015867
3143    0.229845
3144    0.019409
3145    0.061508
3146    0.000000
Name: INFECTION RATE, Length: 3147, dtype: float64
```

In [4]:

```
1 data = data.drop('INFECTION RATE', axis = 1)
2 data
```

Out[4]:

	GEO_ID	COUNTY	STATE	POPULATION	AGE15_19	AGE20_24	AGE25_29	A
0	0500000US01001	Autauga County	Alabama	55200.0	7.0	6.0	6.6	
1	0500000US01003	Baldwin County	Alabama	208107.0	6.1	5.2	5.7	
2	0500000US01005	Barbour County	Alabama	25782.0	5.6	6.4	7.2	
3	0500000US01007	Bibb County	Alabama	22527.0	6.7	6.6	6.9	
4	0500000US01009	Blount County	Alabama	57645.0	6.4	5.5	5.9	
...
3142	0500000US56037	Sweetwater County	Wyoming	44117.0	7.2	6.2	7.1	
3143	0500000US56039	Teton County	Wyoming	23059.0	4.0	4.8	9.3	
3144	0500000US56041	Uinta County	Wyoming	20609.0	7.2	5.5	5.4	
3145	0500000US56043	Washakie County	Wyoming	8129.0	7.1	5.0	4.6	
3146	0500000US56045	Weston County	Wyoming	7100.0	4.5	6.7	5.1	

3147 rows × 101 columns

In [5]:

```
1 data.insert(101, "INFECTION RATE", infection_rate)
2 data
```

Out[5]:

	GEO_ID	COUNTY	STATE	POPULATION	AGE15_19	AGE20_24	AGE25_29	A
0	0500000US01001	Autauga County	Alabama	55200.0	7.0	6.0	6.6	
1	0500000US01003	Baldwin County	Alabama	208107.0	6.1	5.2	5.7	
2	0500000US01005	Barbour County	Alabama	25782.0	5.6	6.4	7.2	
3	0500000US01007	Bibb County	Alabama	22527.0	6.7	6.6	6.9	
4	0500000US01009	Blount County	Alabama	57645.0	6.4	5.5	5.9	
...
3142	0500000US56037	Sweetwater County	Wyoming	44117.0	7.2	6.2	7.1	
3143	0500000US56039	Teton County	Wyoming	23059.0	4.0	4.8	9.3	
3144	0500000US56041	Uinta County	Wyoming	20609.0	7.2	5.5	5.4	
3145	0500000US56043	Washakie County	Wyoming	8129.0	7.1	5.0	4.6	
3146	0500000US56045	Weston County	Wyoming	7100.0	4.5	6.7	5.1	

3147 rows × 102 columns

Identify Prediction Set (Counties with 0.0% Infection Rate)

In [6]:

```
1 pred_set = data.loc[data['INFECTION RATE'] == 0.0]
2 data = data.drop(pred_set.index)
3 data
```

Out[6]:

	GEO_ID	COUNTY	STATE	POPULATION	AGE15_19	AGE20_24	AGE25_29	A
0	0500000US01001	Autauga County	Alabama	55200.0	7.0	6.0	6.6	
1	0500000US01003	Baldwin County	Alabama	208107.0	6.1	5.2	5.7	
2	0500000US01005	Barbour County	Alabama	25782.0	5.6	6.4	7.2	
3	0500000US01007	Bibb County	Alabama	22527.0	6.7	6.6	6.9	
4	0500000US01009	Blount County	Alabama	57645.0	6.4	5.5	5.9	
...
3141	0500000US56035	Sublette County	Wyoming	9951.0	6.6	3.4	4.7	
3142	0500000US56037	Sweetwater County	Wyoming	44117.0	7.2	6.2	7.1	
3143	0500000US56039	Teton County	Wyoming	23059.0	4.0	4.8	9.3	
3144	0500000US56041	Uinta County	Wyoming	20609.0	7.2	5.5	5.4	
3145	0500000US56043	Washakie County	Wyoming	8129.0	7.1	5.0	4.6	

2618 rows × 102 columns

In [7]:

```
1 pred_set
```

Out[7]:

	GEO_ID	COUNTY	STATE	POPULATION	AGE15_19	AGE20_24	AGE25_29	AC
30	0500000US01061	Geneva County	Alabama	26491.0	6.3	5.3	5.4	
67	0500000US02013	Aleutians East Borough	Alaska	3425.0	5.0	6.9	8.1	
68	0500000US02016	Aleutians West Census Area	Alaska	5750.0	4.8	9.2	8.6	
71	0500000US02060	Bristol Bay Borough	Alaska	890.0	5.1	4.3	5.8	
72	0500000US02068	Denali Borough	Alaska	2232.0	5.6	13.4	5.5	
...
3114	0500000US55123	Vernon County	Wisconsin	30516.0	6.6	4.7	4.6	
3117	0500000US55129	Washburn County	Wisconsin	15689.0	5.0	4.0	3.9	
3125	0500000US56003	Big Horn County	Wyoming	11901.0	7.0	4.5	5.0	
3139	0500000US56031	Platte County	Wyoming	8673.0	5.7	4.9	4.3	
3146	0500000US56045	Weston County	Wyoming	7100.0	4.5	6.7	5.1	

529 rows × 102 columns

Identify the X and y from Data and Prediction Set

In [8]:

```
1 X = data.iloc[:, 4:101]
2 X = X.drop("LAND AREA(SQMI)", axis = 1) # No total population or land area
3 y = data["INFECTION RATE"]
4 X
```

Out[8]:

AGE15_19 AGE20_24 AGE25_29 AGE30_34 AGE55_59 AGE60_64 TOTAL_MALE TOTAL_

	0	1	2	3	4	...	3141	3142	3143	3144	3145
	7.0	6.0	6.6	6.2	7.5	4.9	26874.0	101188.0	13697.0	12152.0	28434.0
	6.1	5.2	5.7	5.5	6.9	7.1	6.6	7.0	6.4	6.1	6.3
	5.6	6.4	7.2	7.0	6.4	6.4	5.367.0	22882.0	11911.0	10505.0	4137.0
	6.7	6.6	6.9	6.7	6.6	5.4	5.367.0	22882.0	11911.0	10505.0	4137.0
	6.4	5.5	5.9	5.9	6.8	6.3	5.367.0	22882.0	11911.0	10505.0	4137.0

	6.6	3.4	4.7	6.0	7.3	6.8	5.367.0	22882.0	11911.0	10505.0	4137.0
	7.2	6.2	7.1	7.7	7.0	6.1	5.367.0	22882.0	11911.0	10505.0	4137.0
	4.0	4.8	9.3	9.1	6.4	6.6	5.367.0	22882.0	11911.0	10505.0	4137.0
	7.2	5.5	5.4	6.5	8.5	5.7	5.367.0	22882.0	11911.0	10505.0	4137.0
	7.1	5.0	4.6	5.0	7.1	7.5	5.367.0	22882.0	11911.0	10505.0	4137.0

2618 rows × 96 columns

In [9]:

```
1 pred_X = pred_set.iloc[:, 4:101]
2 pred_X = pred_X.drop("LAND AREA(SQMI)", axis = 1) # No total population or land
3 pred_X
```

Out[9]:

	AGE15_19	AGE20_24	AGE25_29	AGE30_34	AGE55_59	AGE60_64	TOTAL_MALE	TOTAL_
--	----------	----------	----------	----------	----------	----------	------------	--------

30	6.3	5.3	5.4	5.8	7.7	6.4	12971.0
67	5.0	6.9	8.1	6.3	7.7	7.2	1989.0
68	4.8	9.2	8.6	8.8	7.4	6.6	3432.0
71	5.1	4.3	5.8	4.4	11.7	8.4	526.0
72	5.6	13.4	5.5	6.5	7.6	14.2	1117.0
...
3114	6.6	4.7	4.6	5.0	8.0	7.5	15302.0
3117	5.0	4.0	3.9	4.7	8.8	9.4	7771.0
3125	7.0	4.5	5.0	5.6	6.2	7.4	6033.0
3139	5.7	4.9	4.3	5.4	6.4	9.3	4499.0
3146	4.5	6.7	5.1	6.0	5.3	11.4	3768.0

529 rows × 96 columns

Split Data into 2/3 Training and 1/3 Testing

In [10]:

```
1 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.33, random
```

SKLearn GradientBoostingRegressor

Baseline SKLearn GradientBoostingRegressor Model

In [11]:

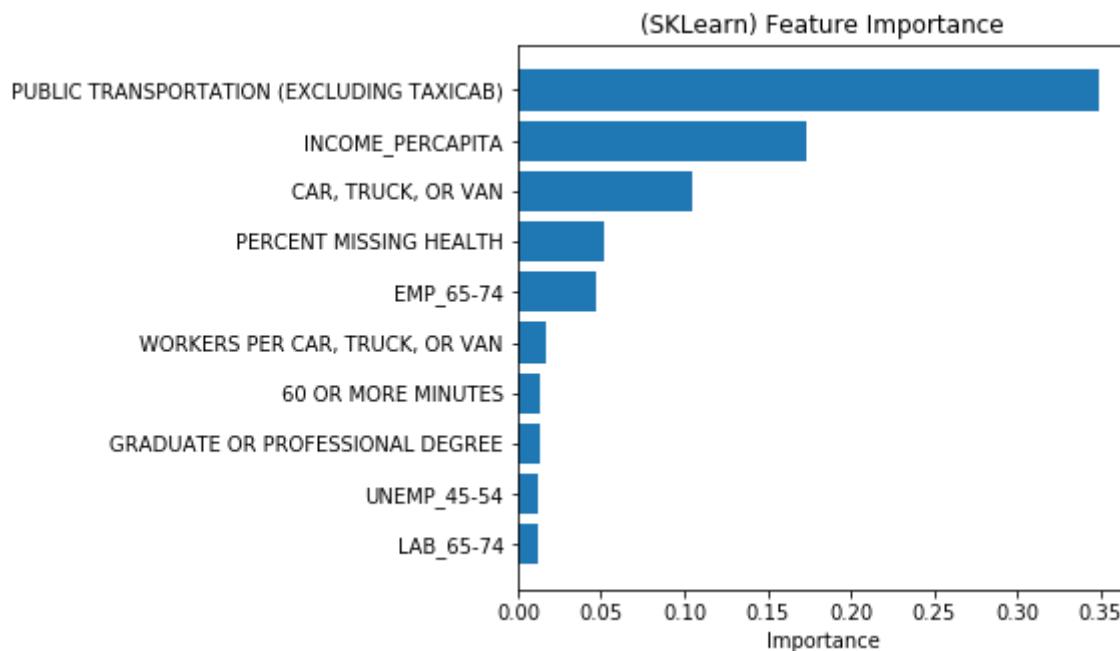
```
1 gbr_model = ensemble.GradientBoostingRegressor()
2 gbr_model.fit(X_train, y_train)
3 mean_squared_error(y_test, gbr_model.predict(X_test))
```

Out[11]:

0.02117833305438754

In [12]:

```
1 feature_importance = gbr_model.feature_importances_
2 #feature_importance = 100.0 * (feature_importance / feature_importance.max())
3 sorted_idx = np.argsort(feature_importance)
4 top_10 = sorted_idx[-10:]
5 pos = np.arange(top_10.shape[0]) + .5
6 plt.figure(figsize=(12, 5))
7 plt.subplot(1, 2, 2)
8 plt.barh(pos, feature_importance[top_10], align='center')
9 plt.yticks(pos, X.columns[top_10])
10 plt.xlabel('Importance')
11 plt.title('(SKLearn) Feature Importance')
12 plt.show()
```



SKLearn: Hyperparameter Tuning

In [14]:

```
1 test_fold = [-1]*len(X_train) + [0]*len(X_test)
2 predefined_split = PredefinedSplit(test_fold=test_fold)
3 predefined_split
```

Out[14]:

PredefinedSplit(test_fold=array([-1, -1, ..., 0, 0]))

In [15]:

```
1 param_grid = {  
2     'max_depth': [3, 5, 7],  
3     'min_samples_split': [2, 4, 6],  
4     'n_estimators': [300, 500, 700],  
5     'learning_rate': [0.001, 0.01, 0.1]  
6 }  
7 gbr = ensemble.GradientBoostingRegressor()  
8 grid = GridSearchCV(  
9     gbr,  
10    param_grid,  
11    scoring = make_scorer(mean_squared_error, greater_is_better = False),  
12    return_train_score=True  
13 )  
14 grid.fit(X_train, y_train)  
15  
16 # pd.set_option('display.max_rows', 5)  
17 df = pd.DataFrame(grid.cv_results_)  
18 # Flip sign of score back, because GridSearchCV likes to maximize,  
19 # so it flips the sign of the score if "greater_is_better = FALSE"  
20 df['mean_test_score'] = -df['mean_test_score']  
21 df['mean_train_score'] = -df['mean_train_score']  
22 cols_to_keep = ["param_learning_rate", "param_max_depth", "param_min_samples_split",  
23                  "mean_test_score", "mean_train_score"]  
24 df_toshow = df[cols_to_keep].fillna('-')  
25 df_toshow = df_toshow.sort_values(by=[ "mean_test_score" ])  
26 df_toshow.head()
```

Out[15]:

	param_learning_rate	param_max_depth	param_min_samples_split	param_n_estimators	mean
78	0.1	7		6	300
60	0.1	3		6	300
65	0.1	5		2	700
64	0.1	5		2	500
79	0.1	7		6	500

In [24]:

```
1 best_params = df_toshow.head(1)
2
3 params = {
4     'n_estimators': best_params['param_n_estimators'].item(),
5     'max_depth': best_params['param_max_depth'].item(),
6     'min_samples_split': best_params['param_min_samples_split'].item(),
7     'learning_rate': best_params['param_learning_rate'].item(),
8     'loss': 'ls', 'criterion': 'mse'
9 }
10 gbr_model = ensemble.GradientBoostingRegressor(**params)
11 gbr_model.fit(X_train, y_train)
12 mean_squared_error(y_test, gbr_model.predict(X_test))
13
14 # Worse than baseline model --> Return to the baseline model
```

Out[24]:

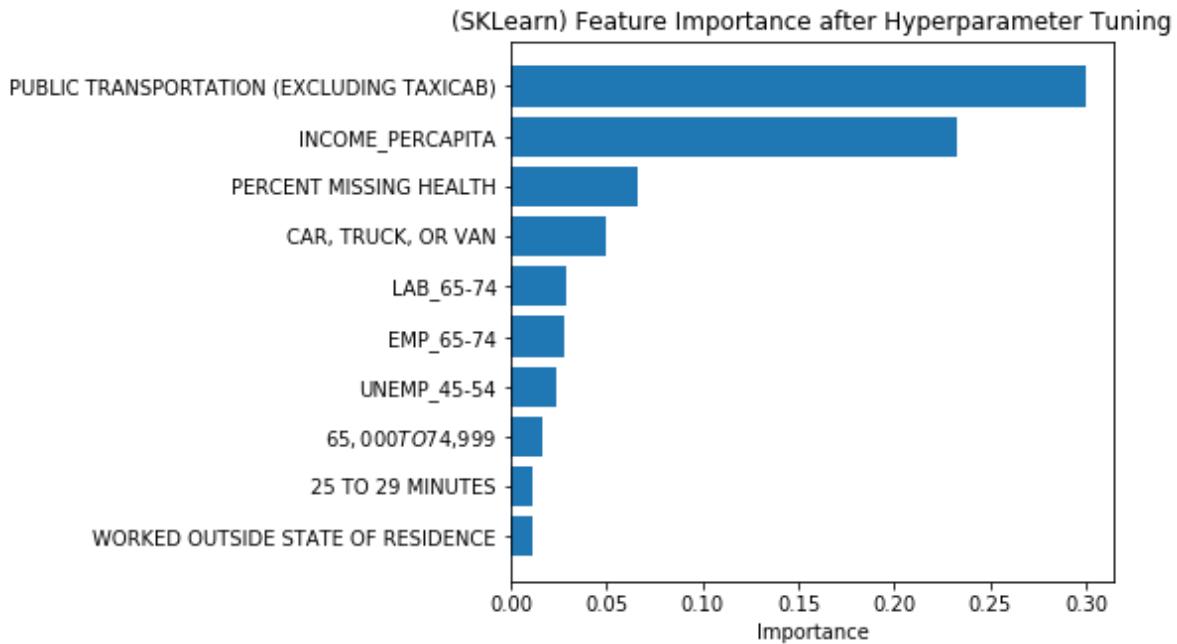
0.02180453198751297

In []:

```
1 # baseline default parameter values
2 params = {
3     'n_estimators': 100,
4     'max_depth': 3,
5     'min_samples_split': 2,
6     'learning_rate': 0.1,
7     'loss': 'ls', 'criterion': 'mse'
8 }
9 gbr_model = ensemble.GradientBoostingRegressor(**params)
10 gbr_model.fit(X_train, y_train)
11 mean_squared_error(y_test, gbr_model.predict(X_test))
```

In [25]:

```
1 feature_importance = gbr_model.feature_importances_
2 #feature_importance = 100.0 * (feature_importance / feature_importance.max())
3 sorted_idx = np.argsort(feature_importance)
4 top_10 = sorted_idx[-10:]
5 pos = np.arange(top_10.shape[0]) + .5
6 plt.figure(figsize=(12, 5))
7 plt.subplot(1, 2, 2)
8 plt.barh(pos, feature_importance[top_10], align='center')
9 plt.yticks(pos, X.columns[top_10])
10 plt.xlabel('Importance')
11 plt.title('(SKLearn) Feature Importance after Hyperparameter Tuning')
12 plt.show()
```



SKLearn Infection Rate Prediction

In [30]:

```
1 infection_rate_predictions = gbr_model.predict(pred_X)
2 pred_set["INFECTION RATE"] = infection_rate_predictions
3 pred_set

//anaconda3/lib/python3.7/site-packages/ipykernel_launcher.py:2: SettingWithCopyWarning:
A value is trying to be set on a copy of a slice from a DataFrame.
Try using .loc[row_indexer,col_indexer] = value instead

See the caveats in the documentation: https://pandas.pydata.org/pandas-docs/stable/user\_guide/indexing.html#returning-a-view-versus-a-copy
(https://pandas.pydata.org/pandas-docs/stable/user\_guide/indexing.html#returning-a-view-versus-a-copy)
```

Out[30]:

	GEO_ID	COUNTY	STATE	POPULATION	AGE15_19	AGE20_24	AGE25_29	AC
30	0500000US01061	Geneva County	Alabama	26491.0	6.3	5.3	5.4	
67	0500000US02013	Aleutians East Borough	Alaska	3425.0	5.0	6.9	8.1	
68	0500000US02016	Aleutians West Census Area	Alaska	5750.0	4.8	9.2	8.6	
71	0500000US02060	Bristol Bay Borough	Alaska	890.0	5.1	4.3	5.8	
72	0500000US02068	Denali Borough	Alaska	2232.0	5.6	13.4	5.5	
...
3114	0500000US55123	Vernon County	Wisconsin	30516.0	6.6	4.7	4.6	
3117	0500000US55129	Washburn County	Wisconsin	15689.0	5.0	4.0	3.9	
3125	0500000US56003	Big Horn County	Wyoming	11901.0	7.0	4.5	5.0	
3139	0500000US56031	Platte County	Wyoming	8673.0	5.7	4.9	4.3	
3146	0500000US56045	Weston County	Wyoming	7100.0	4.5	6.7	5.1	

529 rows × 102 columns

XGBRegressor

Baseline XGBRegressor Model

In [31]:

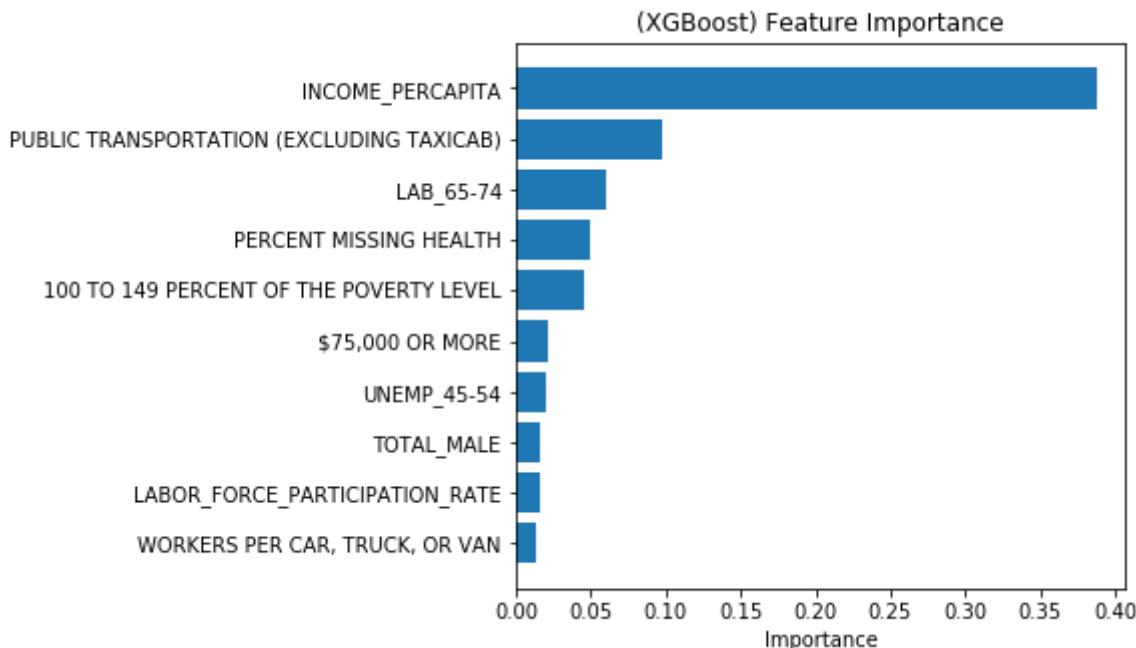
```
1 xgbr_model = XGBRegressor()
2 xgbr_model.fit(X_train, y_train)
3 mean_squared_error(y_test, xgbr_model.predict(X_test))
```

Out[31]:

0.021640949145811064

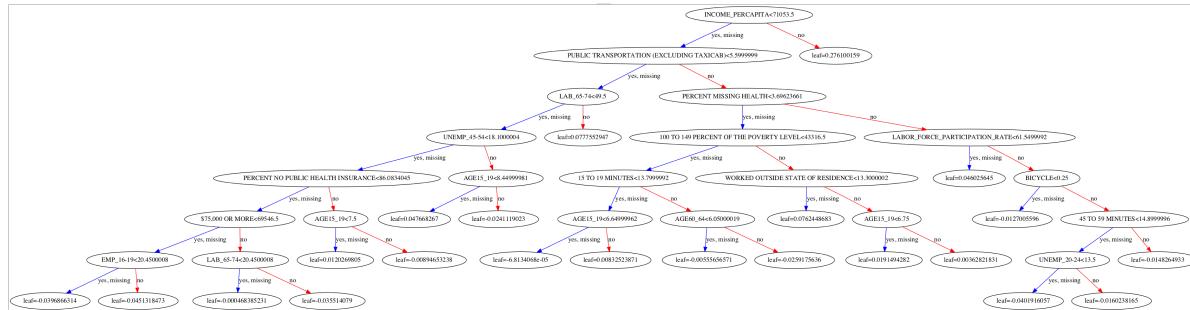
In [32]:

```
1 feature_importance = xgbr_model.feature_importances_
2 sorted_idx = np.argsort(feature_importance)
3 top_10 = sorted_idx[-10:]
4 pos = np.arange(top_10.shape[0]) + .5
5 plt.figure(figsize=(12, 5))
6 plt.subplot(1, 2, 2)
7 plt.barh(pos, feature_importance[top_10], align='center')
8 plt.yticks(pos, X.columns[top_10])
9 plt.xlabel('Importance')
10 plt.title('(XGBoost) Feature Importance')
11 plt.show()
```



In [42]:

```
1 plot_tree(xgbr_model)
2 fig = plt.gcf()
3 fig.set_size_inches(300, 400)
4 plt.title('(XGBoost) Gradient Boosting Tree')
5 plt.show()
```



XGBRegressor: Hyperparameter Tuning

In [34]:

```
1 xparam_grid = {  
2     'max_depth': [3, 5, 7],  
3     'min_child_weight': [2, 4, 6],  
4     'gamma': [0, 0.1, 0.2],  
5     'n_estimators': [300, 500, 700],  
6     'learning_rate': [0.001, 0.01, 0.1]  
7 }  
8 xgbr = XGBRegressor()  
9 xgrid = GridSearchCV(  
10     xgbr,  
11     xparam_grid,  
12     scoring = make_scorer(mean_squared_error, greater_is_better = False),  
13     return_train_score=True  
14 )  
15 xgrid.fit(X_train, y_train)  
16  
17 # pd.set_option('display.max_rows', 20)  
18 xdf = pd.DataFrame(xgrid.cv_results_)  
19 # Flip sign of score back, because GridSearchCV likes to maximize,  
20 # so it flips the sign of the score if "greater_is_better = FALSE"  
21 xdf['mean_test_score'] = -xdf['mean_test_score']  
22 xdf['mean_train_score'] = -xdf['mean_train_score']  
23 cols_to_keep = ["param_gamma", "param_learning_rate", "param_max_depth",  
24                 "param_min_child_weight", "param_n_estimators",  
25                 "mean_test_score", "mean_train_score"]  
26 xdf_toshow = xdf[cols_to_keep].fillna('-')  
27 xdf_toshow = xdf_toshow.sort_values(by=["mean_test_score"])  
28 xdf_toshow.head()
```

Out[34]:

	param_gamma	param_learning_rate	param_max_depth	param_min_child_weight	param_n_e
48	0.0	0.01	7		4
49	0.0	0.01	7		4
129	0.1	0.01	7		4
120	0.1	0.01	5		4
40	0.0	0.01	5		4

In [35]:

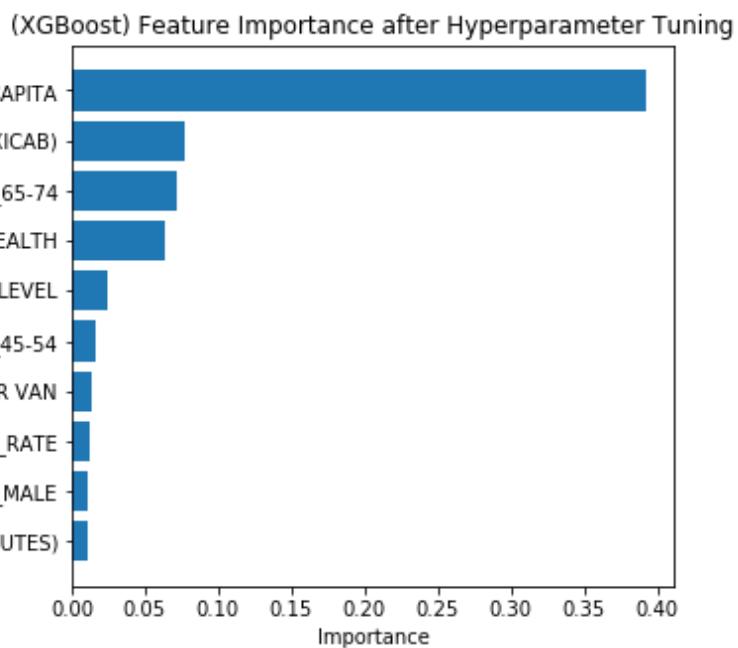
```
1 xbest_params = xdf_toshow.head(1)
2
3 xparams = {
4     'n_estimators': xbest_params['param_n_estimators'].item(),
5     'max_depth': xbest_params['param_max_depth'].item(),
6     'min_child_weight': xbest_params['param_min_child_weight'].item(),
7     'learning_rate': xbest_params['param_learning_rate'].item(),
8     'gamma': xbest_params['param_gamma'].item(),
9     'loss': 'ls', 'criterion': 'mse'
10 }
11 xgbr_model = XGBRegressor(**params)
12 xgbr_model.fit(X_train, y_train)
13 mean_squared_error(y_test, xgbr_model.predict(X_test))
14
15 # Better than baseline model
```

Out[35]:

0.02097232781096379

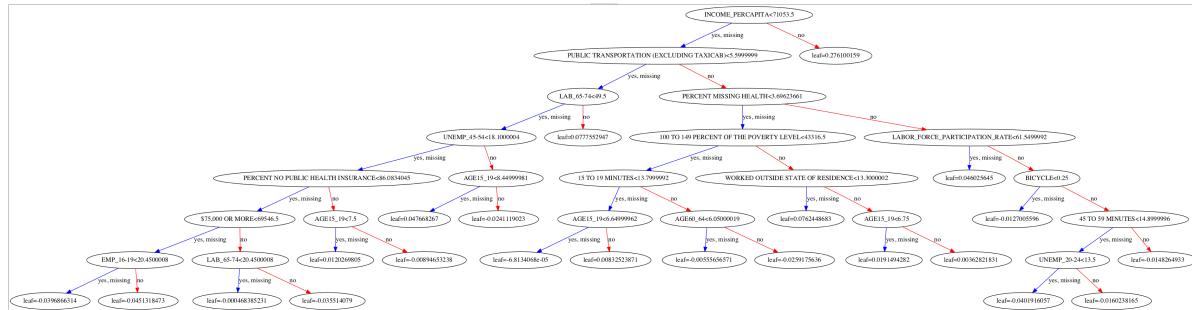
In [36]:

```
1 feature_importance = xgbr_model.feature_importances_
2 sorted_idx = np.argsort(feature_importance)
3 top_10 = sorted_idx[-10:]
4 pos = np.arange(top_10.shape[0]) + .5
5 plt.figure(figsize=(12, 5))
6 plt.subplot(1, 2, 2)
7 plt.barh(pos, feature_importance[top_10], align='center')
8 plt.yticks(pos, X.columns[top_10])
9 plt.xlabel('Importance')
10 plt.title('(XGBoost) Feature Importance after Hyperparameter Tuning')
11 plt.show()
```



In [41]:

```
1 plot_tree(xgbr_model)
2 fig = plt.gcf()
3 fig.set_size_inches(300, 400)
4 plt.title('(XGBoost) Gradient Boosting Tree after Hyperparameter Tuning')
5 plt.show()
```



XGBoost Infection Rate Prediction

In [37]:

```
1 xinfection_rate_predictions = xgbr_model.predict(pred_X)
2 pred_set["XGBOOST INFECTIION RATE"] = xinfection_rate_predictions
3 pred_set
```

//anaconda3/lib/python3.7/site-packages/ipykernel_launcher.py:2: SettingWithCopyWarning:

A value is trying to be set on a copy of a slice from a DataFrame.
Try using .loc[row_indexer,col_indexer] = value instead

See the caveats in the documentation: https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-view-versus-a-copy
(https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-view-versus-a-copy)

Out[37]:

	GEO_ID	COUNTY	STATE	POPULATION	AGE15_19	AGE20_24	AGE25_29	AC
30	0500000US01061	Geneva County	Alabama	26491.0	6.3	5.3	5.4	
67	0500000US02013	Aleutians East Borough	Alaska	3425.0	5.0	6.9	8.1	
68	0500000US02016	Aleutians West Census Area	Alaska	5750.0	4.8	9.2	8.6	
71	0500000US02060	Bristol Bay Borough	Alaska	890.0	5.1	4.3	5.8	
72	0500000US02068	Denali Borough	Alaska	2232.0	5.6	13.4	5.5	
...
3114	0500000US55123	Vernon County	Wisconsin	30516.0	6.6	4.7	4.6	
3117	0500000US55129	Washburn County	Wisconsin	15689.0	5.0	4.0	3.9	
3125	0500000US56003	Big Horn County	Wyoming	11901.0	7.0	4.5	5.0	
3139	0500000US56031	Platte County	Wyoming	8673.0	5.7	4.9	4.3	
3146	0500000US56045	Weston County	Wyoming	7100.0	4.5	6.7	5.1	

529 rows × 103 columns

In [44]:

```
1 sum(pred_set[ "INFECTION RATE" ] - pred_set[ "XGBOOST INFECTION RATE" ]) / len(pred_set)
```

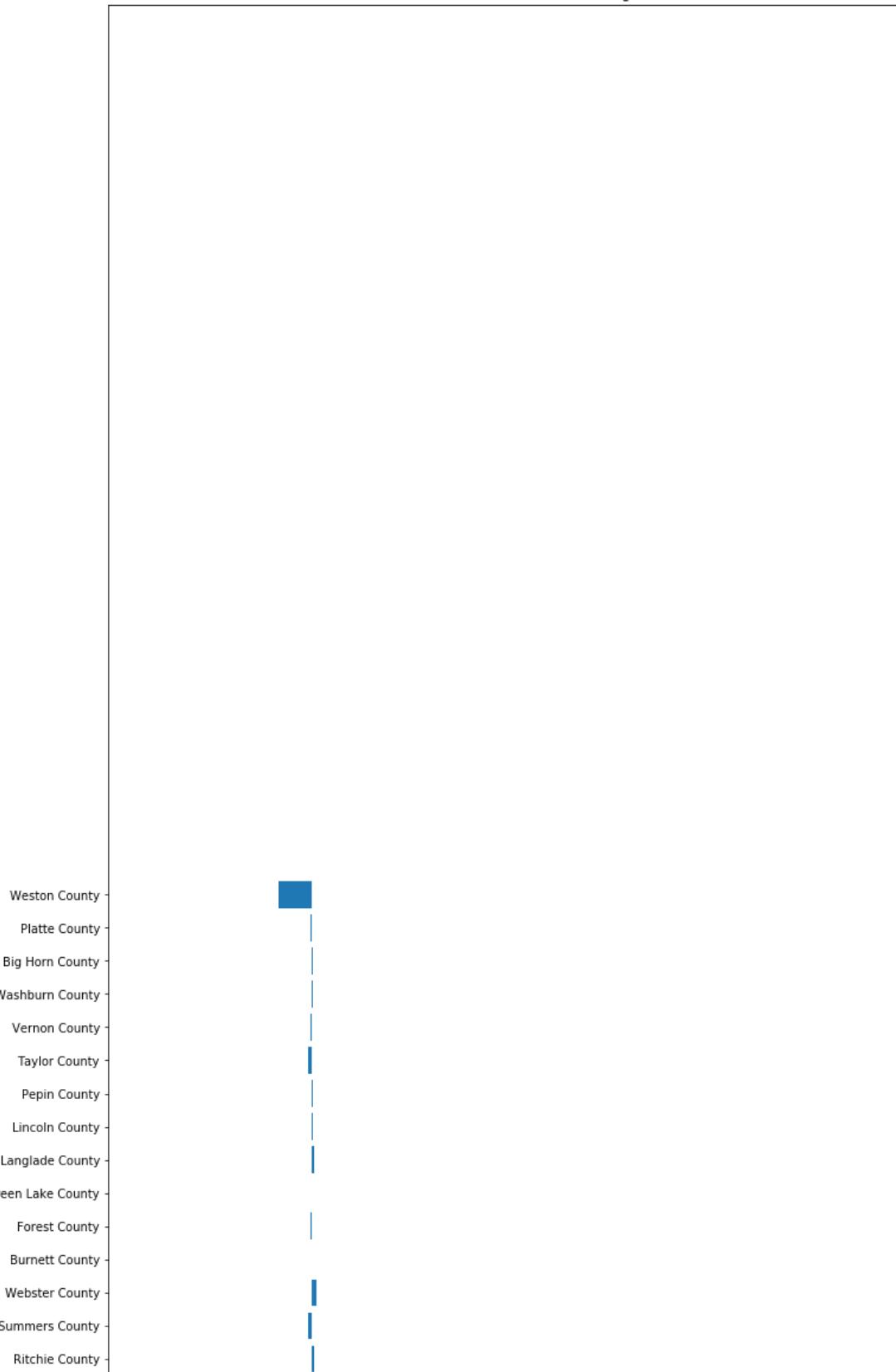
Out[44]:

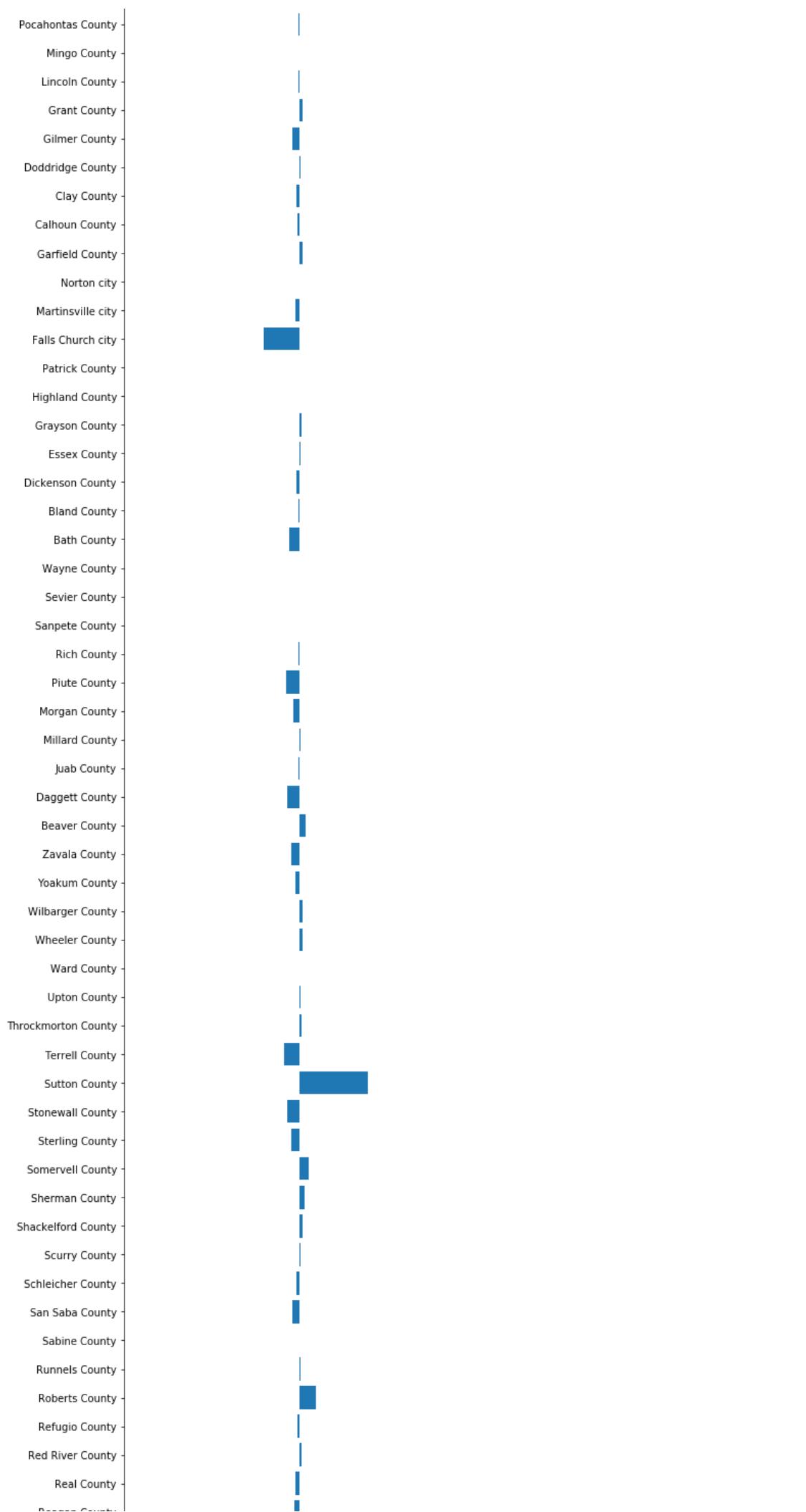
0.03650171111683398

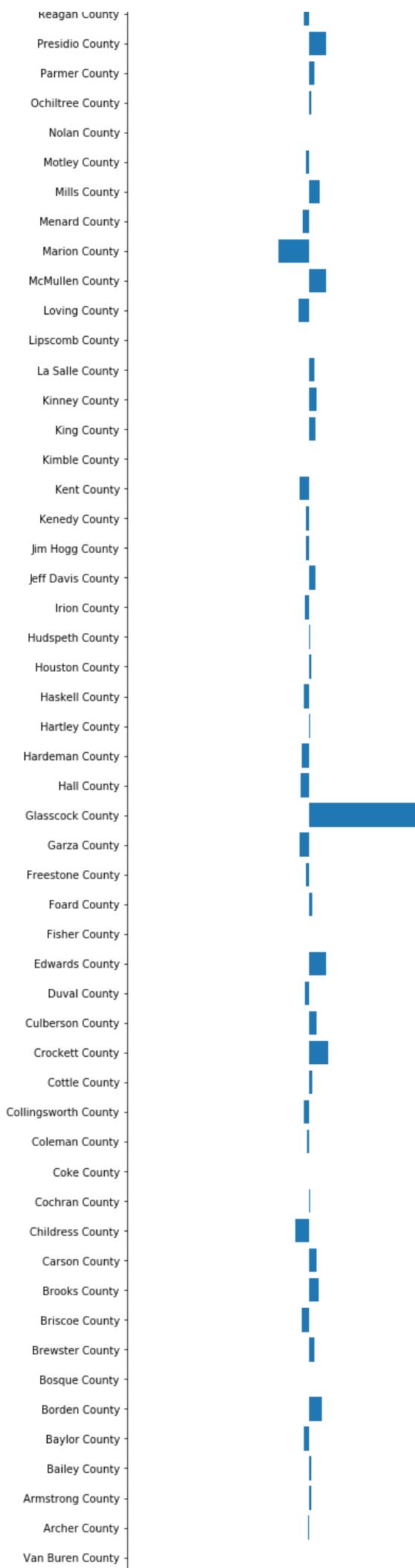
In [64]:

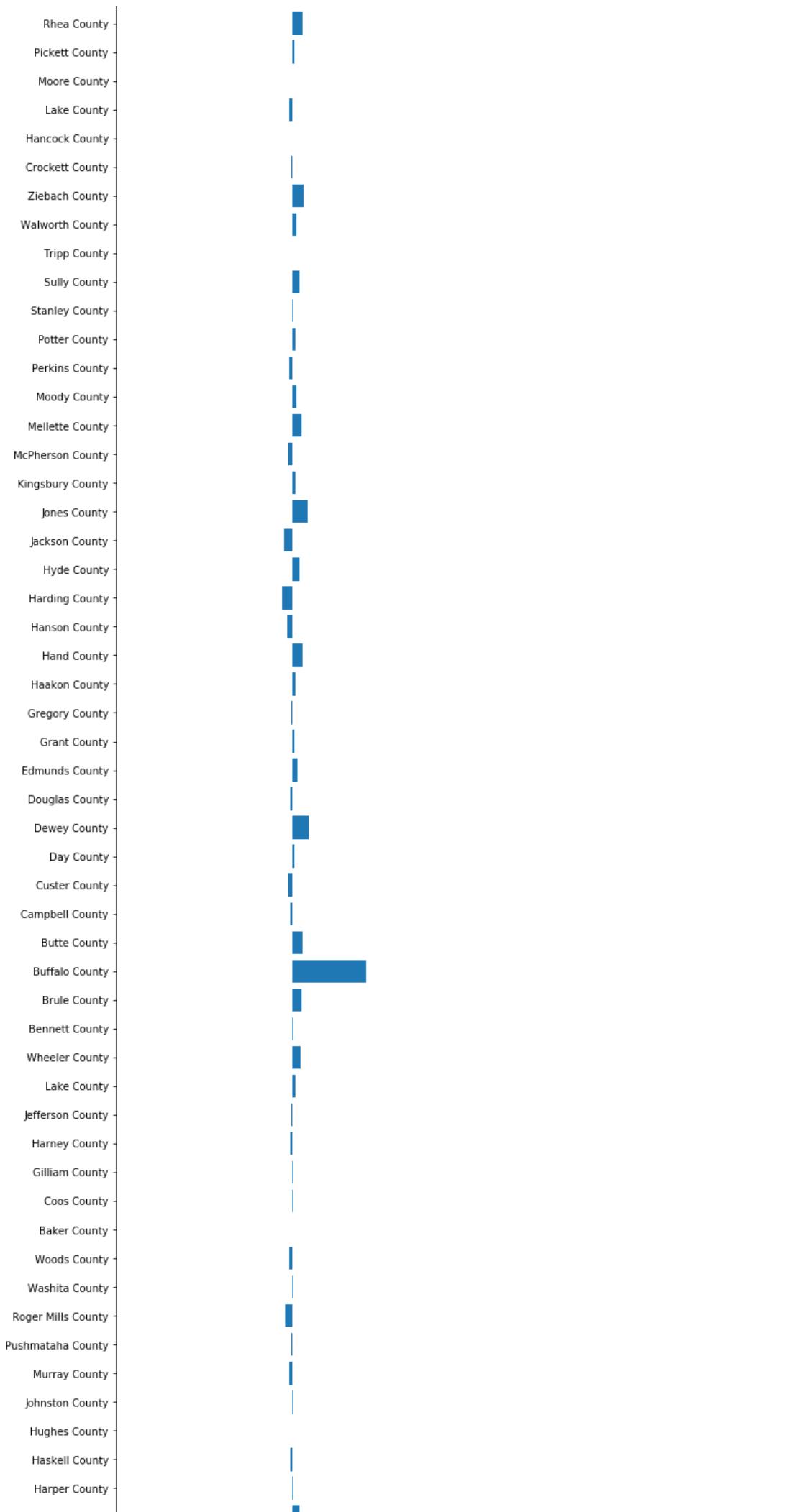
```
1 difference = pred_set["INFECTION RATE"] - pred_set["XGBOOST INFECTION RATE"]
2
3 plt.figure(figsize=(12, 300))
4 plt.barh(np.arange(len(pred_set)), difference, align='center')
5 plt.yticks(np.arange(len(pred_set)), pred_set["COUNTY"])
6 plt.xlabel('Difference')
7 plt.title('Difference between Two Gradient Boosting Models')
8
9 plt.show()
```

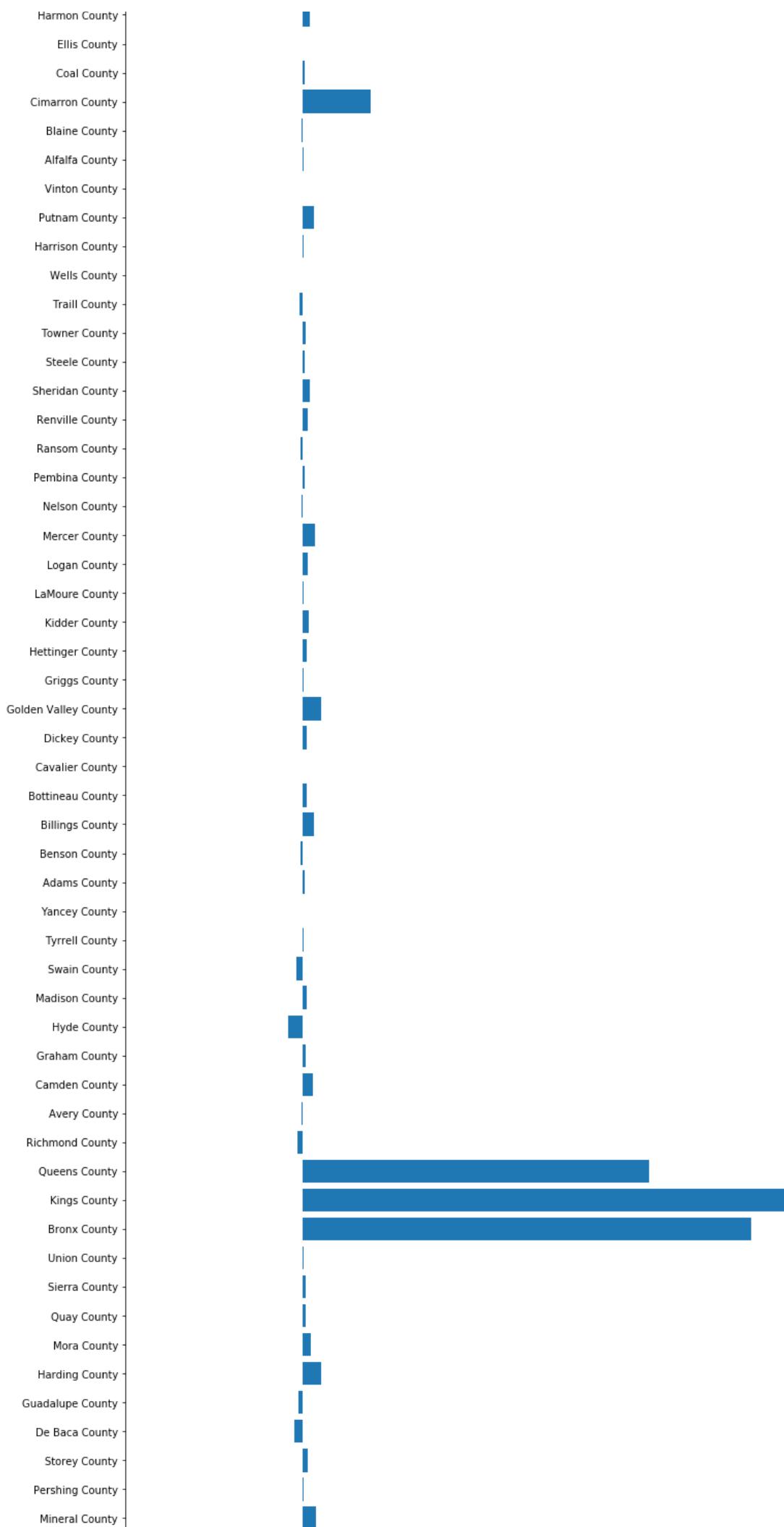
Difference between Two Gradient Boosting Models

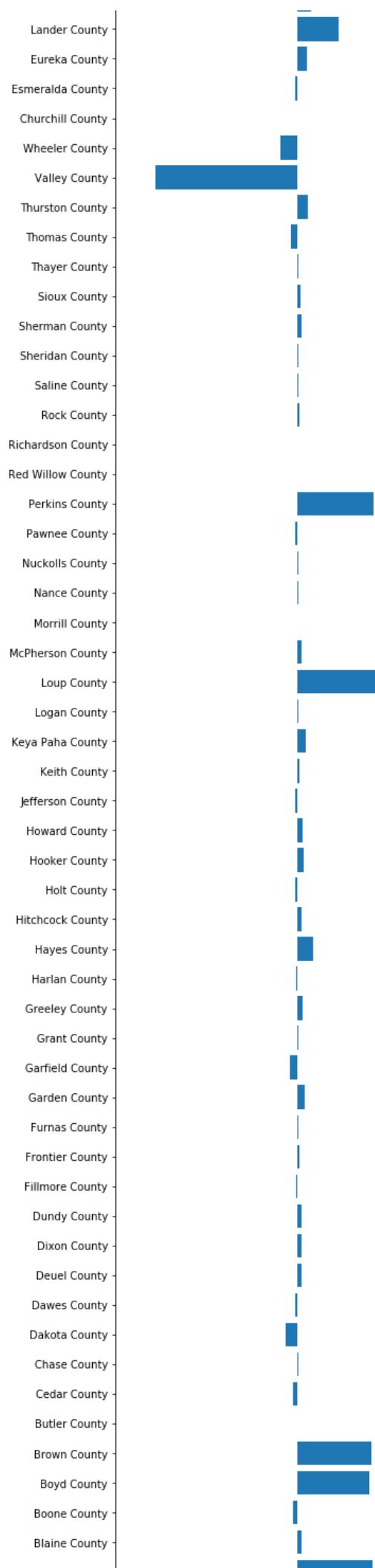


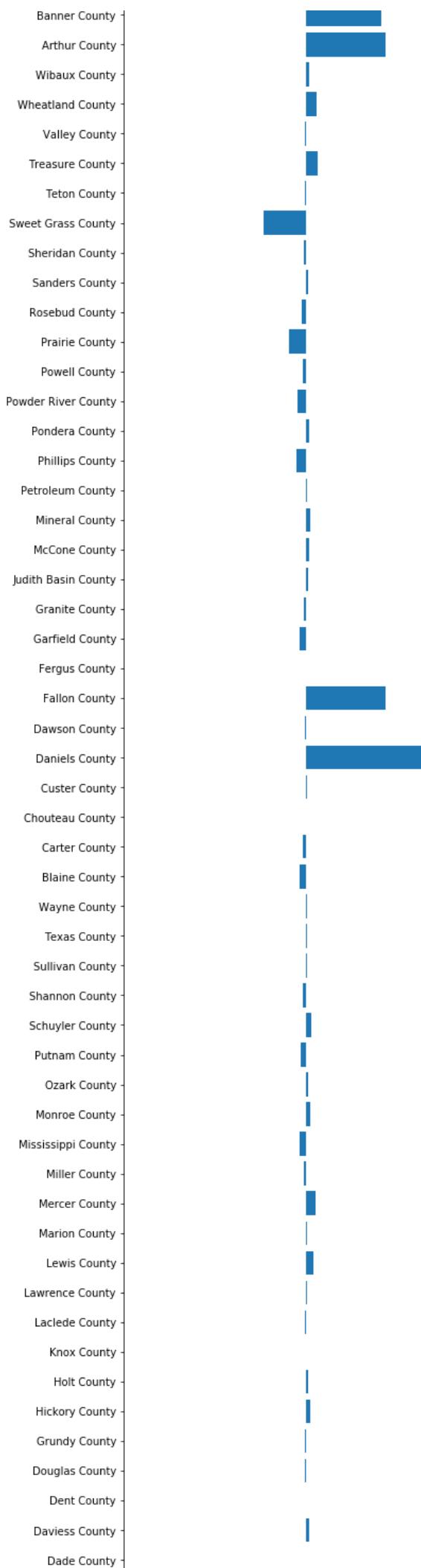


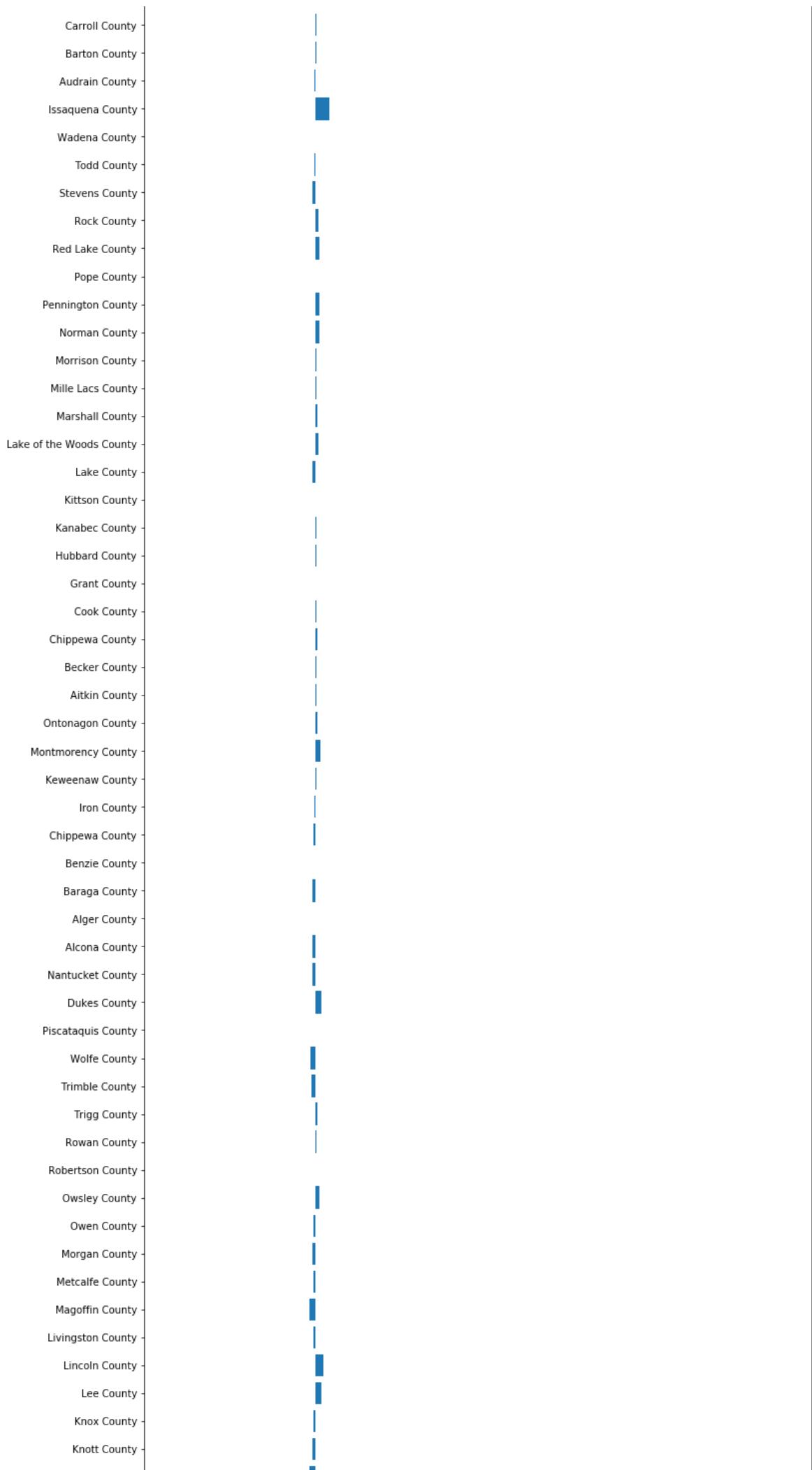


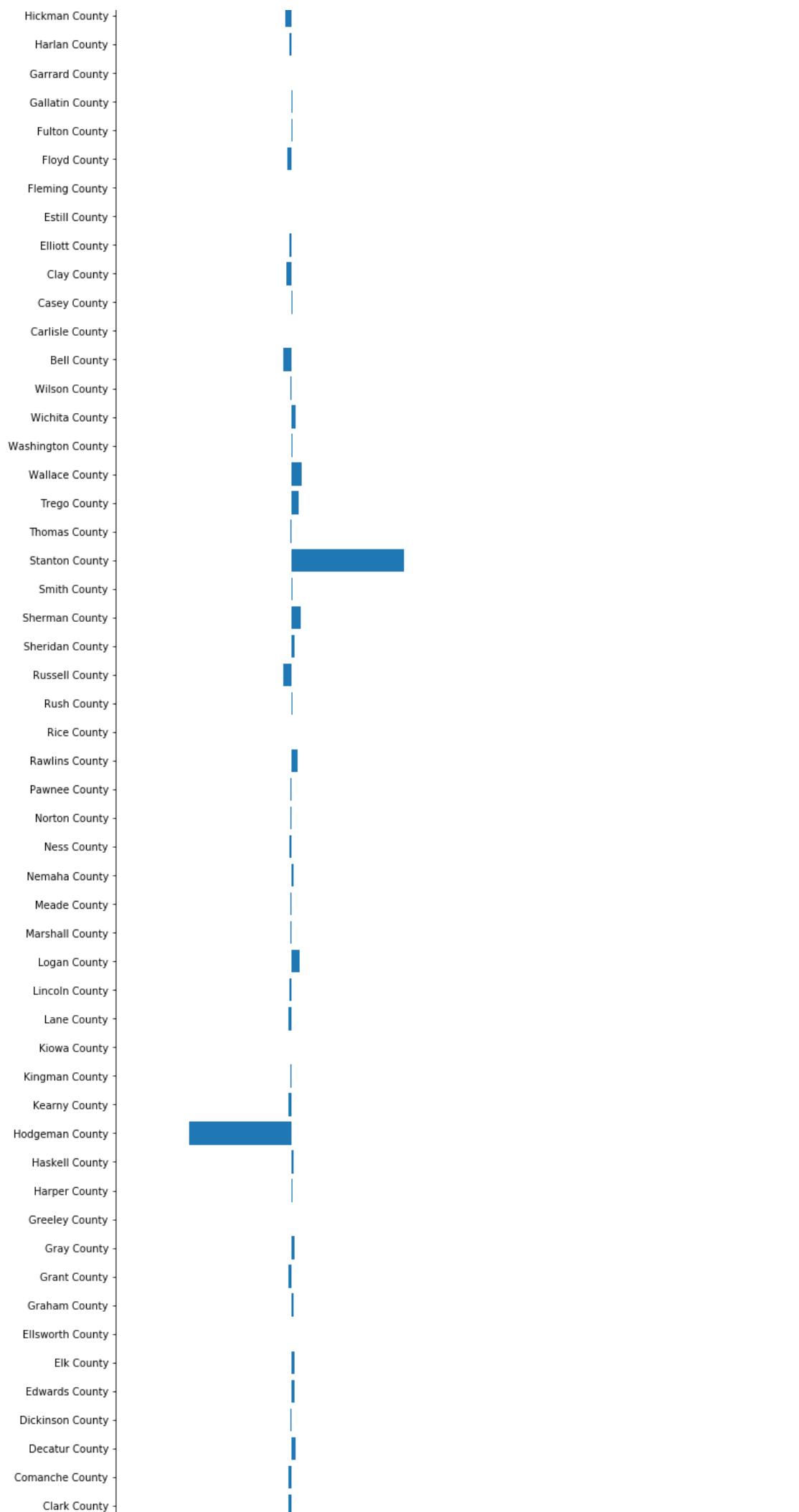


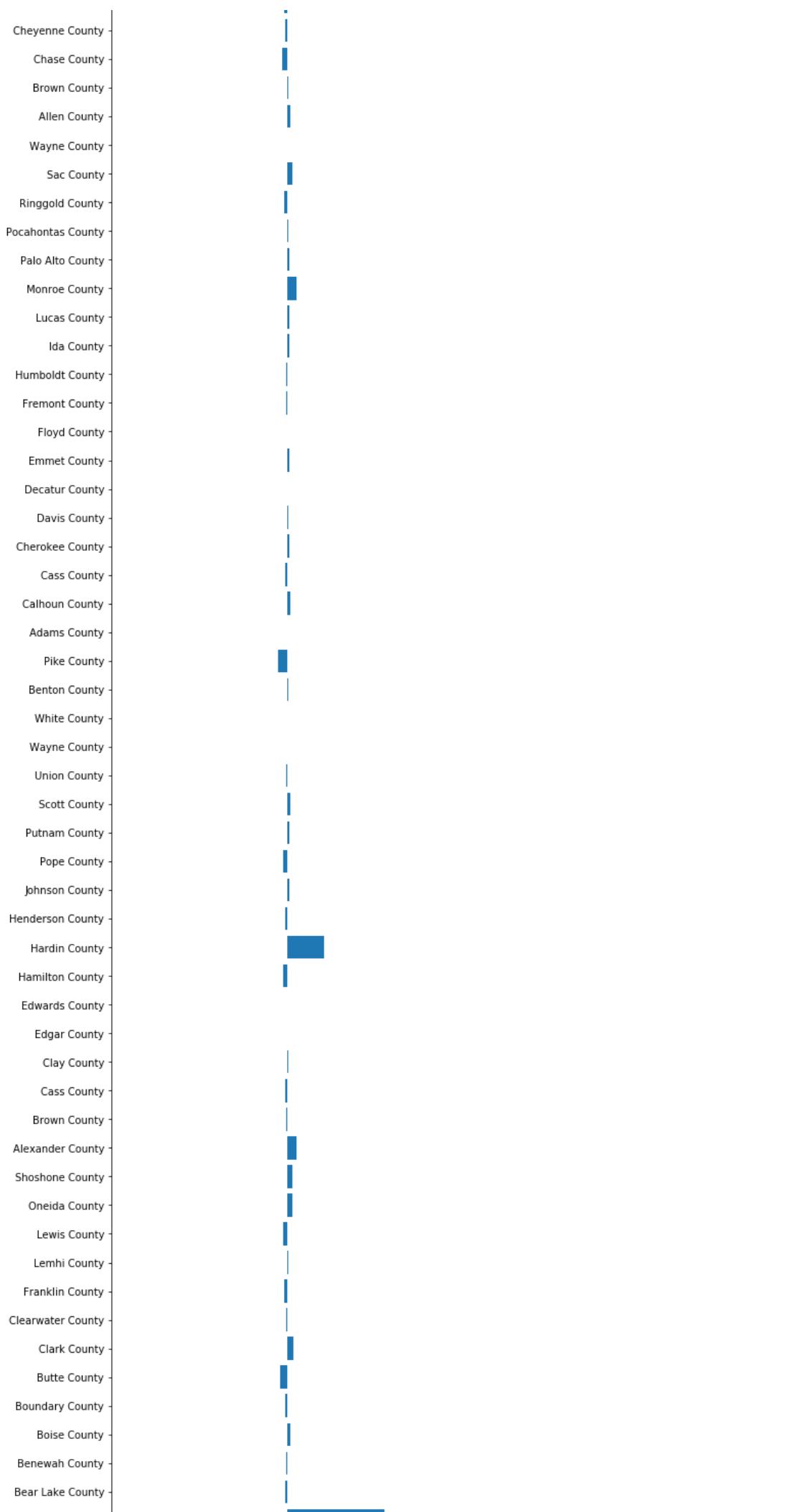


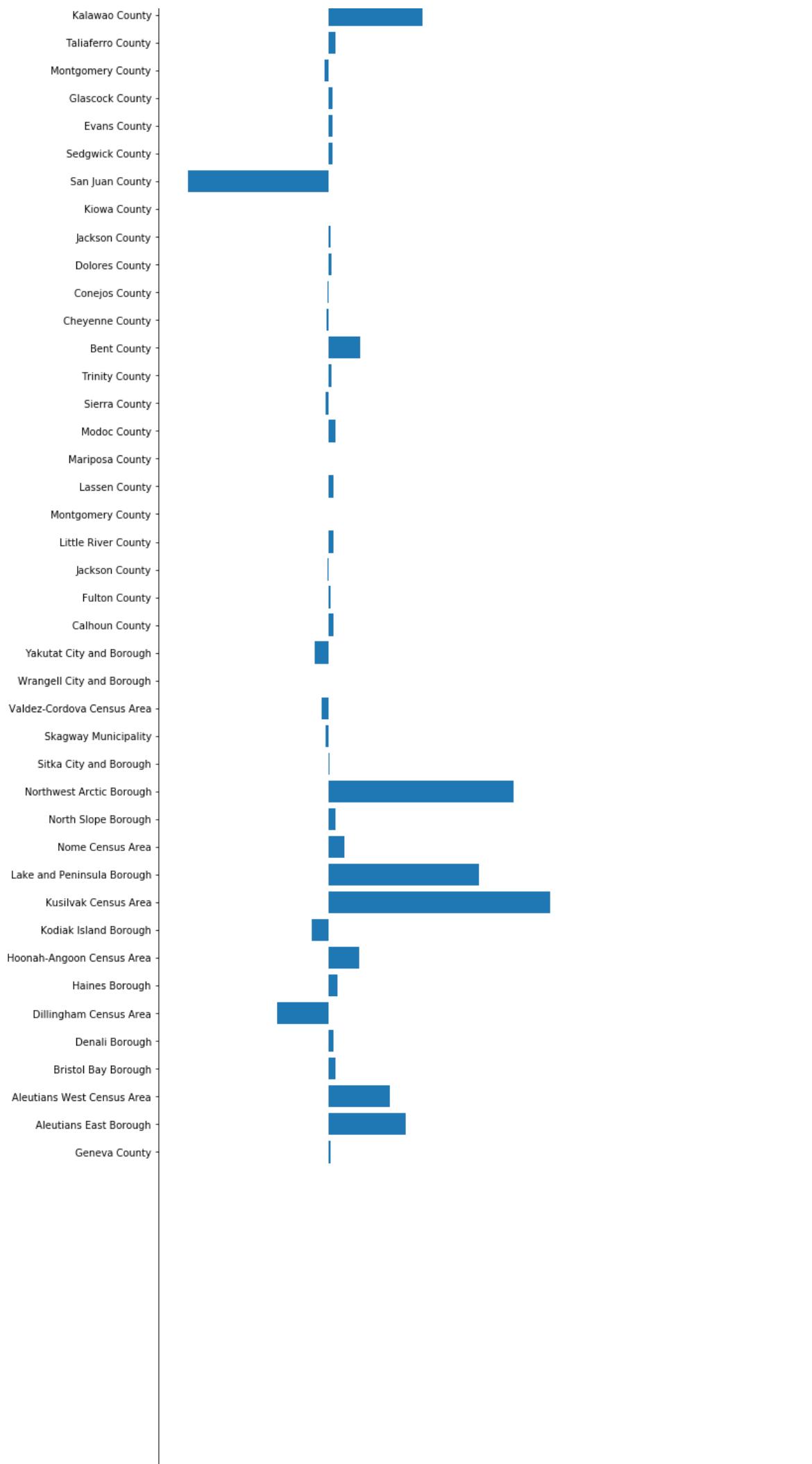


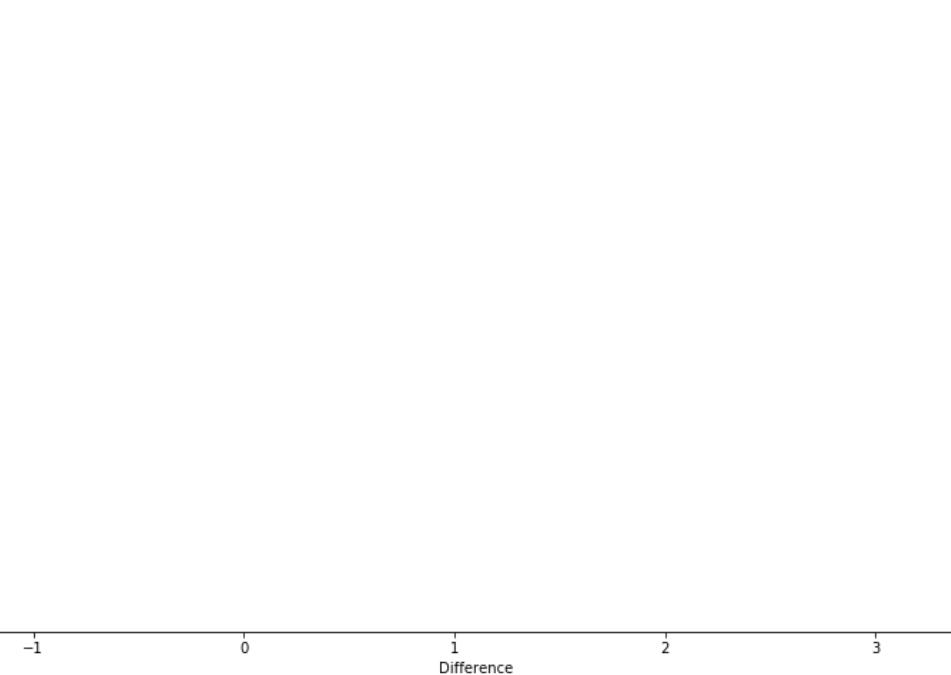












In [57]:

```
1 pred_set[ "DIFF" ] = difference
2 pred_set
```

//anaconda3/lib/python3.7/site-packages/ipykernel_launcher.py:1: SettingWithCopyWarning:

A value is trying to be set on a copy of a slice from a DataFrame.
Try using .loc[row_indexer,col_indexer] = value instead

See the caveats in the documentation: https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-view-versus-a-copy
https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-view-versus-a-copy

"""Entry point for launching an IPython kernel.

Out[57]:

	GEO_ID	COUNTY	STATE	POPULATION	AGE15_19	AGE20_24	AGE25_29	AC
30	0500000US01061	Geneva County	Alabama	26491.0	6.3	5.3	5.4	
67	0500000US02013	Aleutians East Borough	Alaska	3425.0	5.0	6.9	8.1	
68	0500000US02016	Aleutians West Census Area	Alaska	5750.0	4.8	9.2	8.6	
71	0500000US02060	Bristol Bay Borough	Alaska	890.0	5.1	4.3	5.8	
72	0500000US02068	Denali Borough	Alaska	2232.0	5.6	13.4	5.5	
...
3114	0500000US55123	Vernon County	Wisconsin	30516.0	6.6	4.7	4.6	
3117	0500000US55129	Washburn County	Wisconsin	15689.0	5.0	4.0	3.9	
3125	0500000US56003	Big Horn County	Wyoming	11901.0	7.0	4.5	5.0	
3139	0500000US56031	Platte County	Wyoming	8673.0	5.7	4.9	4.3	
3146	0500000US56045	Weston County	Wyoming	7100.0	4.5	6.7	5.1	

529 rows × 104 columns

In [71]:

```
1 # 0.3% is the current national infection rate
2 counties_bigDiff = pred_set.loc[abs(pred_set["DIFF"]) >= 0.3]
3 counties_bigDiff.sort_values(by=['DIFF'], ascending=False)
```

Out[71]:

	GEO_ID	COUNTY	STATE	POPULATION	AGE15_19	AGE20_24	AGE25_29	A
1851	0500000US36047	Kings County	New York	2600747.0	5.4	6.6	9.5	
1830	0500000US36005	Bronx County	New York	1437872.0	6.8	7.9	8.5	
1868	0500000US36081	Queens County	New York	2298513.0	5.1	6.2	8.3	
81	0500000US02158	Kusilvak Census Area	Alaska	8198.0	9.1	8.0	8.8	
86	0500000US02188	Northwest Arctic Borough	Alaska	7734.0	8.1	7.7	9.0	
82	0500000US02164	Lake and Peninsula Borough	Alaska	1375.0	6.7	9.6	7.6	
1607	0500000US30019	Daniels County	Montana	1753.0	5.4	2.1	4.5	
981	0500000US20187	Stanton County	Kansas	2063.0	8.0	3.8	5.5	
2609	0500000US48173	Glasscock County	Texas	1430.0	9.2	11.0	2.1	
548	0500000US15005	Kalawao County	Hawaii	75.0	2.7	0.0	2.7	
1711	0500000US31115	Loup County	Nebraska	585.0	3.6	5.8	5.8	
67	0500000US02013	Aleutians East Borough	Alaska	3425.0	5.0	6.9	8.1	
1610	0500000US30025	Fallon County	Montana	2838.0	4.8	5.9	7.6	
1656	0500000US31005	Arthur County	Nebraska	418.0	4.8	6.2	2.6	
2369	0500000US46017	Buffalo County	South Dakota	2053.0	10.8	8.0	5.9	
1721	0500000US31135	Perkins County	Nebraska	2907.0	5.3	3.5	5.9	
1657	0500000US31007	Banner County	Nebraska	696.0	4.3	3.3	3.4	
1662	0500000US31017	Brown County	Nebraska	2988.0	3.8	3.6	4.4	

	GEO_ID	COUNTY	STATE	POPULATION	AGE15_19	AGE20_24	AGE25_29	A
1661	0500000US31015	Boyd County	Nebraska	2042.0	5.3	3.3	3.5	
2740	0500000US48435	Sutton County	Texas	3865.0	4.3	6.0	7.9	
2143	0500000US40025	Cimarron County	Oklahoma	2189.0	4.3	3.9	8.0	
68	0500000US02016	Aleutians West Census Area	Alaska	5750.0	4.8	9.2	8.6	
73	0500000US02070	Dillingham Census Area	Alaska	4975.0	7.6	7.9	7.9	
929	0500000US20083	Hodgeman County	Kansas	1842.0	5.6	6.2	6.1	
1741	0500000US31175	Valley County	Nebraska	4224.0	5.6	5.0	5.1	
300	0500000US08111	San Juan County	Colorado	544.0	1.8	2.0	11.0	

26 rows × 104 columns

In []:

1	
---	--