

using Local Outlier Factor (LOF)

The anomaly score of each sample is called Local Outlier Factor. It measures the local deviation of density of a given sample with respect to its neighbors. It is local in that the anomaly score depends on how isolated the object is with respect to the surrounding neighborhood. More precisely, locality is given by k-nearest neighbors, whose distance is used to estimate the local density. By comparing the local density of a sample to the local densities of its neighbors, one can identify samples that have a substantially lower density than their neighbors. These are considered outliers.

In [1]:

```
1 import numpy as np
2 import matplotlib.pyplot as plt
3 from sklearn.neighbors import LocalOutlierFactor
4
```

In [2]:

```
1 np.random.seed(42) # everyone gets the same result
```

In [3]:

```
1 # Generate train data
2 X_inliers = 0.3 * np.random.randn(100, 2)
3 X_inliers = np.r_[X_inliers + 2, X_inliers - 2]
```

In [4]:

```
1 # Generate some outliers
2 X_outliers = np.random.uniform(low=-4, high=4, size=(20, 2))
3 X = np.r_[X_inliers, X_outliers]
```

In [5]:

```
1 n_outliers = len(X_outliers)
2 ground_truth = np.ones(len(X), dtype=int)
3 ground_truth[-n_outliers:] = -1
```

In [7]:

```
1 # fit the model for outlier detection (default)
2 clf = LocalOutlierFactor(n_neighbors=20)
3 # use fit_predict to compute the predicted labels of the training samples
4 # (when LOF is used for outlier detection, the estimator has no predict,
5 # decision_function and score_samples methods).
6 y_pred = clf.fit_predict(X)
7 n_errors = (y_pred != ground_truth).sum()
8 print('The number of errors are ', n_errors)
9 X_scores = clf.negative_outlier_factor_
```

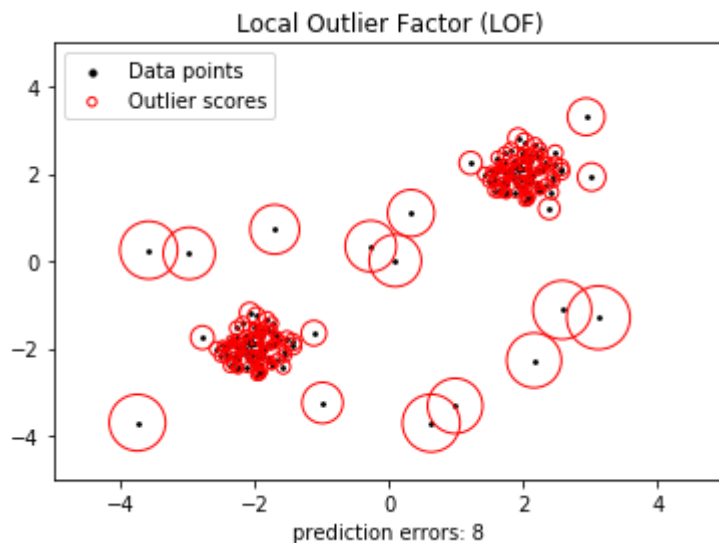
The number of errors are 8

In [8]:

```

1 plt.title("Local Outlier Factor (LOF)")
2 plt.scatter(X[:, 0], X[:, 1], color='k', s=3., label='Data points')
3 # plot circles with radius proportional to the outlier scores
4 radius = (X_scores.max() - X_scores) / (X_scores.max() - X_scores.min())
5 plt.scatter(X[:, 0], X[:, 1], s=1000 * radius, edgecolors='r',
6             facecolors='none', label='Outlier scores')
7 plt.axis('tight')
8 plt.xlim((-5, 5))
9 plt.ylim((-5, 5))
10 plt.xlabel("prediction errors: %d" % (n_errors))
11 legend = plt.legend(loc='upper left')
12 legend.legendHandles[0]._sizes = [10]
13 legend.legendHandles[1]._sizes = [20]
14 plt.show()

```



using isolation forest One efficient way of performing outlier detection in high-dimensional datasets is to use random forests. The ensemble.IsolationForest 'isolates' observations by randomly selecting a feature and then randomly selecting a split value between the maximum and minimum values of the selected feature.

Since recursive partitioning can be represented by a tree structure, the number of splittings required to isolate a sample is equivalent to the path length from the root node to the terminating node.

This path length, averaged over a forest of such random trees, is a measure of normality and our decision function.

Random partitioning produces noticeably shorter paths for anomalies. Hence, when a forest of random trees collectively produce shorter path lengths for particular samples, they are highly likely to be anomalies.

In [9]:

```

1 from sklearn.ensemble import IsolationForest

```

In [10]:

```
1 rng = np.random.RandomState(42)
```

In [11]:

```
1 # Generate train data
2 X = 0.3 * rng.randn(100, 2)
3 X_train = np.r_[X + 2, X - 2]
```

In [12]:

```
1 # Generate some regular novel observations
2 X = 0.3 * rng.randn(20, 2)
3 X_test = np.r_[X + 2, X - 2]
```

In [13]:

```
1 # Generate some abnormal novel observations
2 X_outliers = rng.uniform(low=-4, high=4, size=(20, 2))
```

In [14]:

```
1 # fit the model
2 clf = IsolationForest(max_samples=100, random_state=rng)
3 clf.fit(X_train)
4 y_pred_train = clf.predict(X_train)
5 y_pred_test = clf.predict(X_test)
6 y_pred_outliers = clf.predict(X_outliers)
```

In [29]:

```
1 print(y_pred_test)
2 print(y_pred_outliers)
```

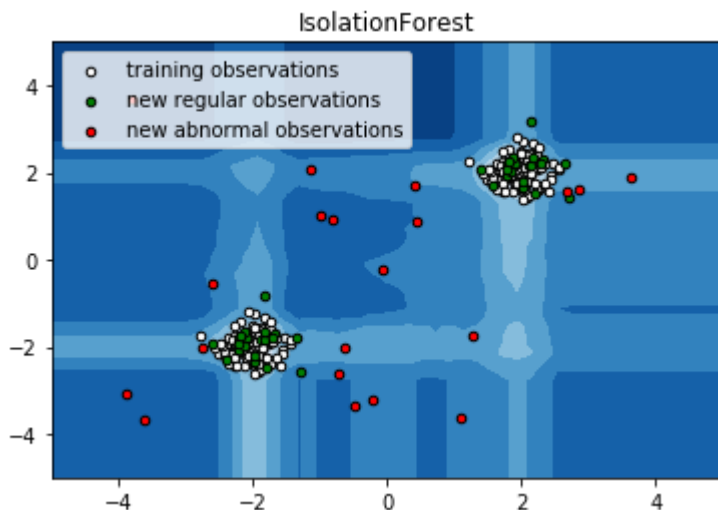
```
[ 1 -1 -1  1 -1 -1 -1  1  1  1 -1 -1  1  1  1  1 -1 -1  1  1 -1 -1
1
-1  1 -1  1  1  1 -1 -1  1  1  1  1  1 -1 -1  1]
[-1 -1 -1 -1 -1 -1 -1 -1 -1 -1 -1 -1 -1 -1 -1 -1 -1 -1 -1]
```

In [15]:

```

1  # plot the line, the samples, and the nearest vectors to the plane
2  xx, yy = np.meshgrid(np.linspace(-5, 5, 50), np.linspace(-5, 5, 50))
3  Z = clf.decision_function(np.c_[xx.ravel(), yy.ravel()])
4  Z = Z.reshape(xx.shape)
5
6  plt.title("IsolationForest")
7  plt.contourf(xx, yy, Z, cmap=plt.cm.Blues_r)
8
9  b1 = plt.scatter(X_train[:, 0], X_train[:, 1], c='white',
10                  s=20, edgecolor='k')
11  b2 = plt.scatter(X_test[:, 0], X_test[:, 1], c='green',
12                  s=20, edgecolor='k')
13  c = plt.scatter(X_outliers[:, 0], X_outliers[:, 1], c='red',
14                  s=20, edgecolor='k')
15  plt.axis('tight')
16  plt.xlim((-5, 5))
17  plt.ylim((-5, 5))
18  plt.legend([b1, b2, c],
19             ["training observations",
20              "new regular observations", "new abnormal observations"],
21             loc="upper left")
22  plt.show()

```



Comparison

In [16]:

```

1  import time
2
3  import numpy as np
4  import matplotlib
5  import matplotlib.pyplot as plt
6
7  from sklearn import svm
8  from sklearn.datasets import make_moons, make_blobs
9  from sklearn.covariance import EllipticEnvelope
10 from sklearn.ensemble import IsolationForest
11 from sklearn.neighbors import LocalOutlierFactor

```

In [17]:

```
1 matplotlib.rcParams['contour.negative_linestyle'] = 'solid'
```

In [18]:

```
1 # Example settings
2 n_samples = 300
3 outliers_fraction = 0.15
4 n_outliers = int(outliers_fraction * n_samples)
5 n_inliers = n_samples - n_outliers
```

In [19]:

```
1 # define outlier/anomaly detection methods to be compared
2 anomaly_algorithms = [
3     ("Isolation Forest", IsolationForest(contamination=outliers_fraction,
4                                           random_state=42)),
5     ("Local Outlier Factor", LocalOutlierFactor(
6         n_neighbors=35, contamination=outliers_fraction))]
```

In [20]:

```
1 # Define datasets
2 blobs_params = dict(random_state=0, n_samples=n_inliers, n_features=2)
3 datasets = [
4     make_blobs(centers=[[0, 0], [0, 0]], cluster_std=0.5,
5                 **blobs_params)[0],
6     make_blobs(centers=[[2, 2], [-2, -2]], cluster_std=[0.5, 0.5],
7                 **blobs_params)[0],
8     make_blobs(centers=[[2, 2], [-2, -2]], cluster_std=[1.5, .3],
9                 **blobs_params)[0],
10    4. * (make_moons(n_samples=n_samples, noise=.05, random_state=0)[0] -
11          np.array([0.5, 0.25])),
12    14. * (np.random.RandomState(42).rand(n_samples, 2) - 0.5)]
```

In [21]:

```
1 # Compare given classifiers under given settings
2 xx, yy = np.meshgrid(np.linspace(-7, 7, 150),
3                      np.linspace(-7, 7, 150))
```

In [22]:

```
1 plt.figure(figsize=(len(anomaly_algorithms) * 2 + 3, 12.5))
2 plt.subplots_adjust(left=.02, right=.98, bottom=.001, top=.96, wspace=.05,
3                     hspace=.01)
4
5 plot_num = 1
6 rng = np.random.RandomState(42)
```

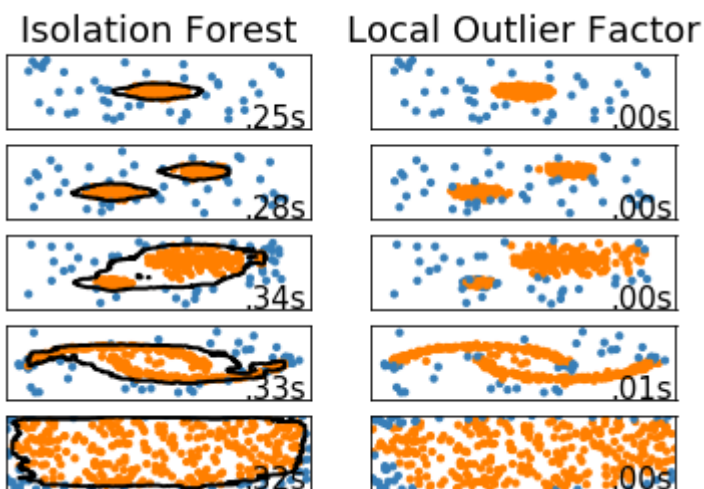
<Figure size 504x900 with 0 Axes>

In [23]:

```

1  for i_dataset, X in enumerate(datasets):
2      # Add outliers
3      X = np.concatenate([X, rng.uniform(low=-6, high=6,
4                                          size=(n_outliers, 2))], axis=0)
5
6      for name, algorithm in anomaly_algorithms:
7          t0 = time.time()
8          algorithm.fit(X)
9          t1 = time.time()
10         plt.subplot(len(datasets), len(anomaly_algorithms), plot_num)
11         if i_dataset == 0:
12             plt.title(name, size=18)
13
14         # fit the data and tag outliers
15         if name == "Local Outlier Factor":
16             y_pred = algorithm.fit_predict(X)
17         else:
18             y_pred = algorithm.fit(X).predict(X)
19
20         # plot the levels lines and the points
21         if name != "Local Outlier Factor": # LOF does not implement predict
22             Z = algorithm.predict(np.c_[xx.ravel(), yy.ravel()])
23             Z = Z.reshape(xx.shape)
24             plt.contour(xx, yy, Z, levels=[0], linewidths=2, colors='black')
25
26         colors = np.array(['#377eb8', '#ff7f00'])
27         plt.scatter(X[:, 0], X[:, 1], s=10, color=colors[(y_pred + 1) // 2])
28
29         plt.xlim(-7, 7)
30         plt.ylim(-7, 7)
31         plt.xticks(())
32         plt.yticks(())
33         plt.text(.99, .01, ('%.2fs' % (t1 - t0)).lstrip('0'),
34                  transform=plt.gca().transAxes, size=15,
35                  horizontalalignment='right')
36         plot_num += 1
37
38 plt.show()

```



In []:

1	
---	--