

knn-hepmass-accuracy-vs-k-sklearn

September 2, 2021

```
[1]: import pandas as pd

[2]: full_hepm_data = pd.read_csv('all_train.csv.gz')

[5]: full_hepm_data.to_hdf('hepmass.hdf5', 'hepmass')

[6]: full_hepm_data = pd.read_hdf('hepmass.hdf5', 'hepmass')

[7]: n_samples, n_features = full_hepm_data.shape

[8]: n_samples, n_features

[8]: (7000000, 29)

[9]: train_samples, test_samples = 100000, 50000

[10]: train_data, train_labels = full_hepm_data.iloc[0:train_samples, 1:],  
    ↪full_hepm_data.iloc[0:train_samples, 0]

[11]: train_data.shape, train_labels.shape

[11]: ((100000, 28), (100000,))

[12]: test_data, test_labels = full_hepm_data.iloc[train_samples:(train_samples +  
    ↪test_samples), 1:], full_hepm_data.iloc[train_samples:(train_samples +  
    ↪test_samples), 0]

[13]: test_data.shape, test_labels.shape

[13]: ((50000, 28), (50000,))

[14]: from sklearn.neighbors import KNeighborsClassifier  
    from sklearn.metrics import accuracy_score

[15]: njobs = 4
```

```
[20]: def handle_k(k: int) -> float:
        classifier = KNeighborsClassifier(n_neighbors = k, algorithm = 'brute',
        ↪n_jobs = njobs)
        classifier.fit(train_data, train_labels)
        result = classifier.predict(test_data)
        return accuracy_score(test_labels, result)
```

```
[21]: ks = list(range(1, 10)) + list(range(10, 100, 10))
        rs = []

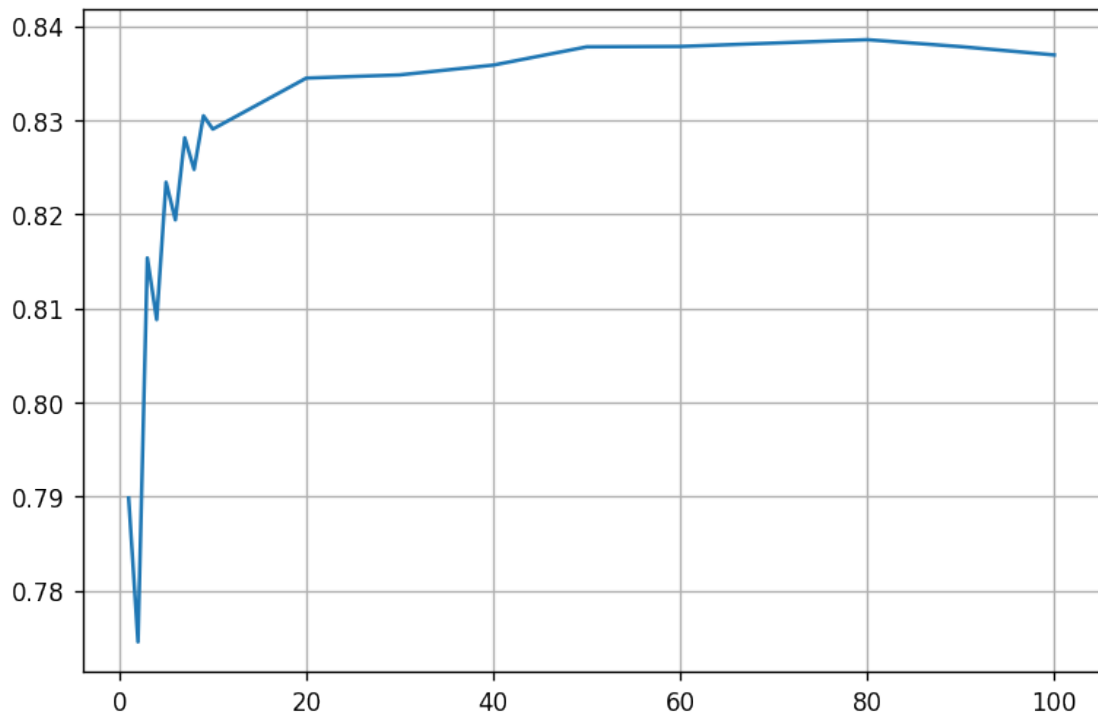
        for k in ks:
            res = handle_k(k)
            print(f'K: {k}, Accuracy: {res}')
            rs.append(res)
```

```
K: 1, Accuracy: 0.78986
K: 2, Accuracy: 0.77456
K: 3, Accuracy: 0.81536
K: 4, Accuracy: 0.8088
K: 5, Accuracy: 0.8234
K: 6, Accuracy: 0.8194
K: 7, Accuracy: 0.82812
K: 8, Accuracy: 0.82476
K: 9, Accuracy: 0.83046
K: 10, Accuracy: 0.82904
K: 20, Accuracy: 0.83446
K: 30, Accuracy: 0.8348
K: 40, Accuracy: 0.83584
K: 50, Accuracy: 0.83778
K: 60, Accuracy: 0.83782
K: 70, Accuracy: 0.83818
K: 80, Accuracy: 0.83854
K: 90, Accuracy: 0.8378
```

```
[22]: ks = ks + [100]
        rs = rs + [handle_k(100)]
```

```
[23]: import matplotlib.pyplot as plt
```

```
[24]: plt.figure(figsize = (7.5, 5), dpi = 120)
        plt.plot(ks, rs)
        plt.grid()
        plt.savefig('hepmass-accuracy-knn-vs-k.pdf')
```



```
[25]: print(rs[-1])
```

```
0.83692
```

```
[ ]:
```