

# COMPSCI 589 Final Project

William Cai, Maxwell Tang

9 May 2025

See `README.md` for running instructions.

## 1 MNIST Dataset

Maxwell's implementations were used for the entirety of the MNIST dataset.

For the MNIST dataset, we chose to use the multilayer perceptron (MLP), k-nearest neighbors (KNN), and random forest (RF) architectures to predict on our dataset. We chose this because decision trees are a subset of random forests and naive bayes is hard to adapt to numerical data.

First, we swept over the secondary hyperparameters in each algorithm. For MLP, we used a training rate of 0.1, a regularization cost of 0.01, and used unbatched gradient descent. As shown in Table 1, you can see that a model shape of [64, 128, 10] achieved the highest accuracy. For RF, we used the set size as a stopping criteria: nodes below the splitting threshold cannot be split. As shown in Table 2, a splitting threshold of 5 achieves the highest accuracy.

As shown in Figure 1,  $k < 5$  is the best. This indicates that the dataset is almost perfectly separated into disjoint clusters. As shown in Figure 2, the f1 score for 0 plateaus after only 100 epochs, whereas the f1 scores for 1, 8, and 9 don't reach a plateau until around 300-400 epochs in. This indicates that 1, 8, and 9 were harder classes to distinguish while 0 was much easier to recognize.

Model Shape	Epoch	Accuracy
[64, 128, 10]	989	0.952145
[64, 64, 10]	977	0.933768
[64, 10]	988	0.920987
[64, 16, 10]	991	0.878687
[64, 16, 16, 10]	998	0.803554
[64, 8, 10]	997	0.722818
[64, 8, 8, 10]	999	0.533731

Table 1: The results for the hyperparameter tuning step on the MLP architecture trained on the MNIST dataset. The best epoch is used for each hyperparameter setting.

Minimum Splittable Size	ntree	Mean Accuracy
5	100	0.977182
10	100	0.975515
20	50	0.967722
30	100	0.954926
40	100	0.953256
50	50	0.943811
2	10	0.523063

Table 2: The results for the hyperparameter tuning step on the random forest architecture trained on the MNIST dataset. The best epoch is used for each hyperparameter setting.

The overall accuracy plateaus after around 300 epochs. As shown in Figure 3, each metric is essentially flat after  $ntree = 20$ .

The best performing model overall was the KNN model with at around 0.99 accuracy. This is way better than the other model architectures, indicating that some specific feature of KNN is really well suited to MNIST. I hypothesize that this is due to the MNIST dataset being very strictly segmented, but in a complex way. Since random forests and MLPs tend to have trouble fitting very detailed data, this explains why MLPs and random forests behaved relatively poorly. Since low-k KNN models fit these sorts of datasets really well, that would explain why KNN had such a high accuracy.

## 2 Rice Dataset

Maxwell’s implementations were used for the entirety of the rice dataset.

For the rice dataset, we chose to use the multilayer perceptron (MLP), k-nearest neighbors (KNN), and random forest (RF) architectures to predict on our dataset. We chose this because decision trees are a subset of random forests and naive bayes is hard to adapt to numerical data.

First, we swept over the secondary hyperparameters in each algorithm. For MLP, we used a training rate of 0.1, a regularization cost of 0.01, and used unbatched gradient descent. As shown in Table 3, you can see that a model shape of [7, 8, 8, 1] achieved the highest accuracy. For RF, we used the set size as a stopping criteria: nodes below the splitting threshold cannot be split. As shown in Table 4, a splitting threshold of 30 achieves the highest accuracy.

As shown in Figure 4, a k value of around 3 or 4 is the best. This indicates that the dataset doesn’t have very strongly mixed classes. As shown in Figure 5, all metrics plateau around 100 epochs in. As shown in Figure 6, the metrics rise steeply and then plateau. The exact values for these phenomena can be found in the referenced graphs and in the code. The best performing model overall was the MLP with an accuracy of 0.93. This is really close the best performance of random forests, which was 0.928, but this was significantly higher than the best performance of the KNN architecture, which was 0.89. This might indicate

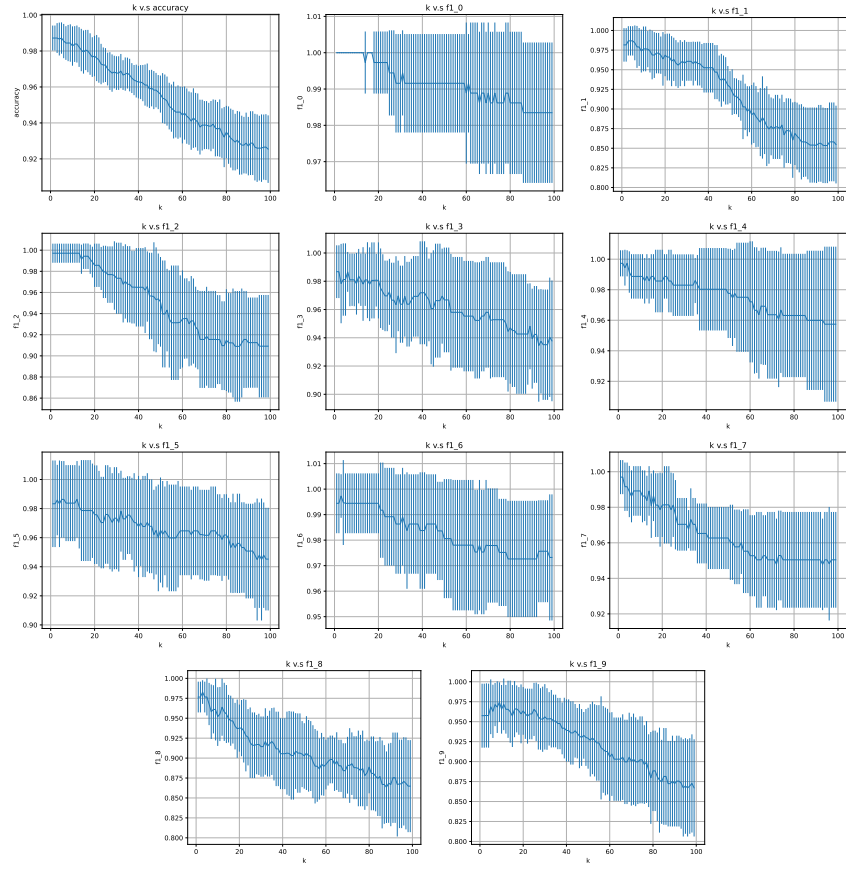


Figure 1: Performance metrics graphed against  $k$  for KNN on the MNIST dataset.

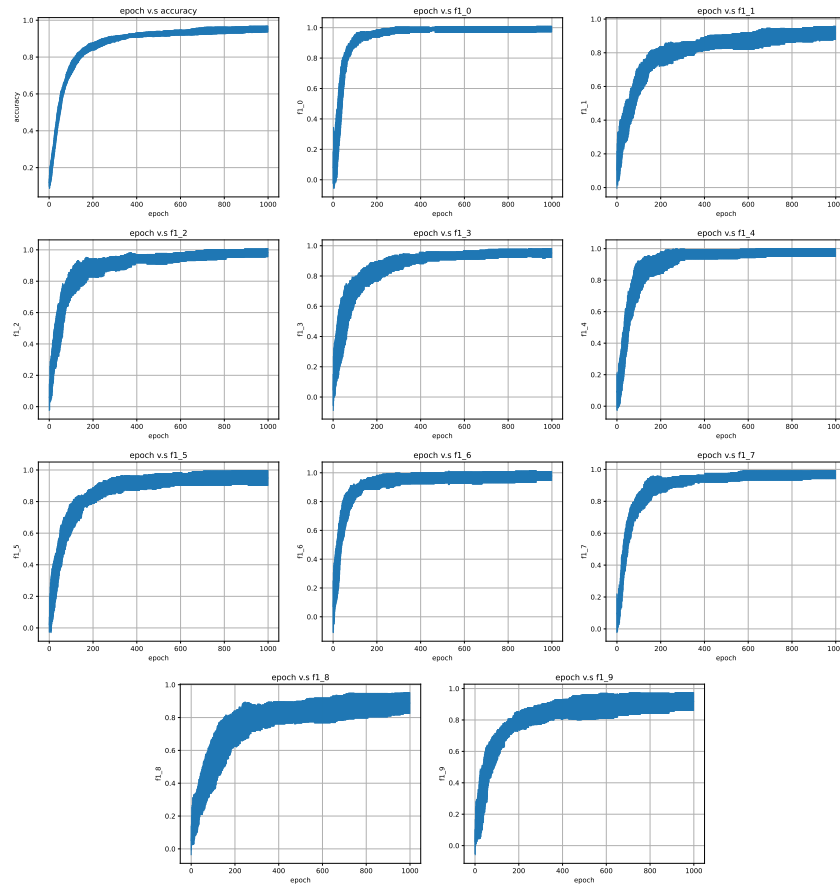


Figure 2: Training graphs for a MLP trained on the MNIST dataset.

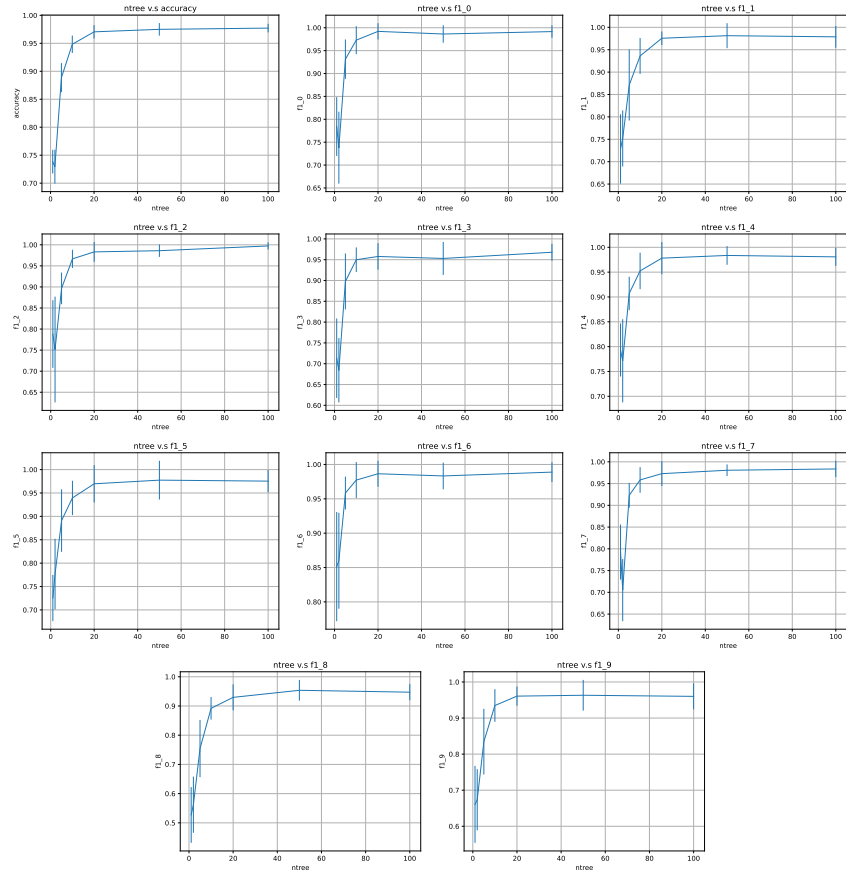


Figure 3: Performance metrics graphed against ntree for random forests trained on the MNIST dataset

Model Shape	Epoch	Mean Accuracy
[7, 8, 8, 1]	765	0.930971
[7, 2, 1]	979	0.929921
[7, 2, 4, 1]	996	0.929659
[7, 1]	670	0.929396
[7, 16, 16, 1]	553	0.928871
[7, 4, 16, 1]	735	0.928871
[7, 2, 16, 1]	656	0.928871
[7, 16, 1]	584	0.928871
[7, 16, 8, 1]	757	0.928084
[7, 8, 2, 1]	969	0.928084
[7, 8, 1]	535	0.928084
[7, 8, 4, 1]	638	0.928084
[7, 16, 2, 1]	760	0.927822
[7, 2, 8, 1]	960	0.927822
[7, 8, 16, 1]	748	0.927822
[7, 16, 4, 1]	732	0.927822
[7, 4, 8, 1]	835	0.927559
[7, 4, 4, 1]	993	0.927297
[7, 4, 1]	781	0.927034
[7, 4, 2, 1]	770	0.925722
[7, 2, 2, 1]	999	0.887139

Table 3: The results for the hyperparameter tuning step on the MLP architecture trained on the rice dataset. The best epoch is used for each hyperparameter setting.

Minimum Splittable Size	ntree	Mean Accuracy
30	100	0.928215
20	100	0.927979
40	100	0.927428
50	100	0.927008
10	50	0.926772
2	50	0.926772
5	50	0.926772

Table 4: The results for the hyperparameter tuning step on the Random Forest architecture trained on the rice dataset. The best ntree value is used for each hyperparameter setting.

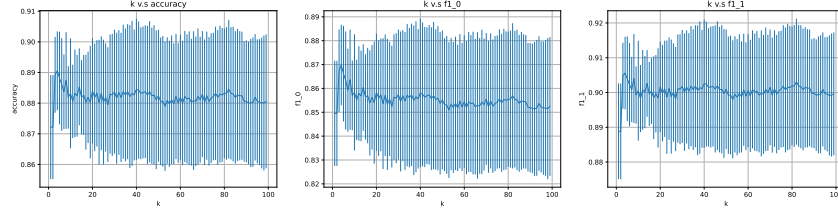


Figure 4: Performance metrics graphed against k for KNN on the rice dataset.

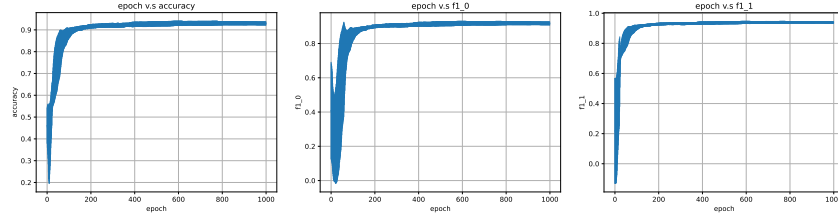


Figure 5: Training graphs for a MLP trained on the rice dataset.

that the dataset has differently behaved regions, where some areas have more outliers than others, making KNN models difficult to fit to the dataset.

### 3 Credit Approval Dataset

This section is implemented by William Cai, including the code and the writeup.

#### 3.1 Nature of the Dataset and Algorithms

The Credit Approval Dataset is composed of 653 credit card applications. It has 6 numerical features and 9 categorical features. The goal is to predict whether a given credit card application should be approved or not.

Because the dataset contains high-dimensional data, traditional data anal-

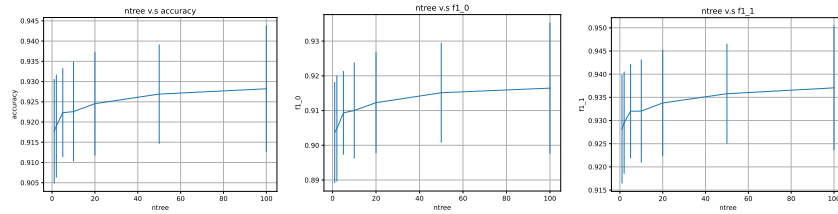


Figure 6: Performance metrics graphed against ntree for random forests trained on the rice dataset

ysis like regressions are ill-suited for the tasks. We should explore multiple powerful machine learning algorithms to address this problem. Because most tools are not good at handling categorical features, we can use one-hot encoding to preprocess the categorical features. This would further expand the dimension of the dataset.

The  $k$ -NN algorithm is good at finding the similarity between instances by calculating their Euclidean distances in the high-dimensional space. This method is effective, since we suspect that all features about the voice attributes are of similar importance. At least, when we do not know about the relative importance of features, treating them equally is a good baseline starting point. The drawback of  $k$ -NN is that the dataset being super high in dimension means many "useless" features that mislead  $k$ -NN to a dead end. But this is not our case:  $6 + 9$  features are moderate, even after one-hot encoding expansion, and they are all clearly related to credit approval or not approval by the dataset provider.

The random forest algorithm could be suitable for this problem. It is built on many decision trees and uses bootstrapping and feature sampling to get robust results. If there are some features that are more important in classifying the disease, the decision tree will find such features based on the larger information gain via entropy or the Gini index. It compliments  $k$ -NN's treating features equally. This allows it to be resistant to some noise in the data. The drawback is that it could be hard to interpret than a simple decision tree, and it is often computationally expensive. We can alleviate this by multiprocessing to use all CPU cores.

We also explore neural networks for this dataset. With adequately sophisticated neural network configuration and non-linear activation functions like Sigmoid or ReLU at each neuron, and enough training data points, a NN can learn and potentially approximate any curves, and discover the hidden rules if there are any. This opens new ways for us to let the neural network learn from data by adjusting weights and biases in an automated fashion, without us to guess how to pick features. This powerful method comes with potential over curve-fitting, so we thus deploy regularization to mitigate such risks. Most complex NNs need intensive parallel computing, which may require RTX A6000, A100, H100, or H200. Fortunately, our dataset is small or moderate, and our CPU is quite capable of doing it.

### 3.1.1 Dataset Preprocessing

Since this dataset has 9 categorical features, in addition to 6 numerical features, we need to use one-hot encoding on these categorical features.

The unique values in those features are:

```
{'attr1_cat': {'a', 'b'},  
'attr4_cat': {'l', 'u', 'y'},  
'attr5_cat': {'g', 'gg', 'p'},  
'attr6_cat': {'aa',
```



```

'c',
'cc',
'd',
'e',
'ff',
'i',
'j',
'k',
'm',
'q',
'r',
'w',
'x'},
'attr7_cat': {'bb', 'dd', 'ff', 'h', 'j', 'n', 'o', 'v', 'z'},
'attr9_cat': {'f', 't'},
'attr10_cat': {'f', 't'},
'attr12_cat': {'f', 't'},
'attr13_cat': {'g', 'p', 's'}})

```

For example, `attr13_cat` has three unique values: `'g'`, `'p'`, `'s'`. The one-hot encodings for this feature are `[1, 0, 0]` for `'g'`, `[0, 1, 0]` for `'p'`, and `[0, 0, 1]` for `'s'`. We therefore replace column `attr13_cat` with three columns: `attr13_cat_g`, `attr13_cat_p`, and `attr13_cat_s`.

We call this encoded file `credit_approval_adj.csv`. Now, the features in the dataset are all numerical, conducive to our data analysis algorithms below.

## 3.2 Algorithms

### 3.2.1 $k$ -NN

The features are first scaled to `[0, 1]` using our own `MinMaxScaler` class before feeding the data to  $k$ -NN. This is a straightforward but important step to avoid meaningless or misleading scaling or unit distortions among features in calculating the distance between two instances, i.e., how similar and close they are.

Our  $k$ -NN code supports many options as shown below.

```

$ python knn_cv.py --help
usage: knn_cv.py [-h] [-p PATH] [-n ROUNDS] [-kmin KMIN] [-kmax KMAX]
                  [-exclude_self EXCLUDE_SELF] [-skip_norm SKIP_NORMALIZATION]
                  [-header HEADER] [-kfold NUM_KFOLDS] [-random_state RANDOM_STATE]

```

optional arguments:

```

-h, --help            show this help message and exit
-p PATH, --path PATH  CSV path of data file, default wdbc_wo_header.csv
-n ROUNDS, --rounds ROUNDS
                        Rounds for KNN run, default 20

```

```

-kmin KMIN, --kmin KMIN
    Min K in KNN, default 1
-kmax KMAX, --kmax KMAX
    Max K for KNN, default 52
-exclude_self EXCLUDE_SELF, --exclude_self EXCLUDE_SELF
    Exclude self in traing, default False
-skip_norm SKIP_NORMALIZATION, --skip_normalization SKIP_NORMALIZATION
    Skip normalization, default False
-header HEADER, --header HEADER
    CSV file header or not, default None
-kfold NUM_KFOLDS, --num_kfolds NUM_KFOLDS
    number of folds for stratified K-Fold, default 5
-random_state RANDOM_STATE, --random_state RANDOM_STATE
    random seed like 42, -5, default None

```

The following command runs  $k$ -NN with our own stratified  $k$ -fold with shuffling.

```

$ time python knn_cv.py --header true --exclude_self false \
    --num_kfolds 10 -p credit_approval_adj.csv

```

Stratified Results:

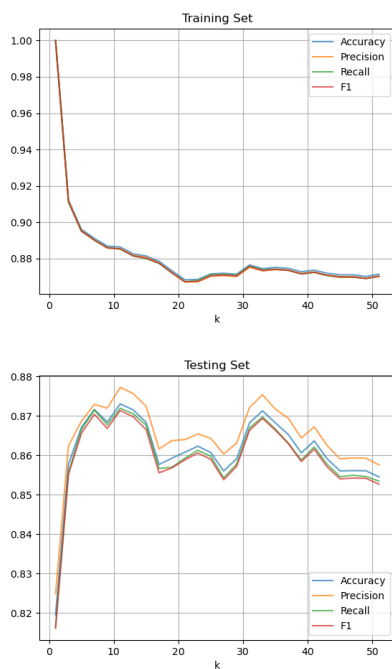
Mean Train Metrics:

	Accuracy	Precision	Recall	F1
1	1.000000	1.000000	1.000000	1.000000
3	0.912202	0.911839	0.910916	0.911331
5	0.896040	0.895191	0.895044	0.895107
7	0.891103	0.890180	0.890112	0.890135
9	0.886853	0.885822	0.885937	0.885871
11	0.886342	0.885441	0.885214	0.885304
13	0.882597	0.881749	0.881370	0.881508
15	0.881404	0.880668	0.879958	0.880257
17	0.878517	0.877540	0.877383	0.877427
19	0.873241	0.872151	0.872109	0.872120
21	0.868301	0.867091	0.867269	0.867174
23	0.868473	0.867140	0.867876	0.867451
25	0.871535	0.870216	0.870964	0.870540
27	0.871876	0.870564	0.871277	0.870878
29	0.871365	0.870025	0.870842	0.870382
31	0.876468	0.875220	0.875925	0.875504
33	0.874424	0.873182	0.873767	0.873424
35	0.875105	0.873969	0.874166	0.874042
37	0.874595	0.873533	0.873507	0.873492
39	0.872723	0.871688	0.871538	0.871582

41	0.873575	0.872528	0.872382	0.872437
43	0.871873	0.870855	0.870601	0.870702
45	0.871022	0.870008	0.869726	0.869842
47	0.871021	0.870039	0.869757	0.869845
49	0.870173	0.869143	0.868982	0.869011
51	0.871364	0.870297	0.870200	0.870218

Mean Test Metrics:

	Accuracy	Precision	Recall	F1
1	0.819488	0.824862	0.816351	0.816155
3	0.857770	0.862433	0.855628	0.855216
5	0.867096	0.868678	0.866533	0.865591
7	0.871618	0.872966	0.871550	0.870431
9	0.868470	0.871992	0.867745	0.866837
11	0.873062	0.877274	0.871951	0.871372
13	0.871547	0.875667	0.870562	0.869821
15	0.868423	0.872482	0.867687	0.866561
17	0.857699	0.861628	0.856682	0.855569
19	0.859285	0.863745	0.856981	0.856832
21	0.860799	0.864012	0.859279	0.858857
23	0.862338	0.865450	0.861281	0.860605
25	0.860728	0.864342	0.859813	0.859001
27	0.856041	0.860347	0.854305	0.853799
29	0.859142	0.863121	0.857753	0.857224
31	0.868280	0.872170	0.866977	0.866384
33	0.871310	0.875383	0.869755	0.869383
35	0.868233	0.871711	0.866642	0.866369
37	0.865274	0.869420	0.863071	0.862982
39	0.860634	0.864413	0.858807	0.858445
41	0.863664	0.867254	0.862140	0.861568
43	0.859119	0.862477	0.857696	0.857051
45	0.856041	0.859118	0.854583	0.854023
47	0.856136	0.859350	0.854958	0.854257
49	0.856088	0.859276	0.854622	0.854159
51	0.854526	0.857605	0.853472	0.852701



As we can see in the above chart, the training set achieves accuracy 100

Let me exclude the data point itself in the classification by setting `--exclude_self true`.

```
$ python knn_cv.py --header true --exclude_self true \
  --num_kfolds 10 -p credit_approval_adj.csv
```

Stratified Results:

Mean Train Metrics:

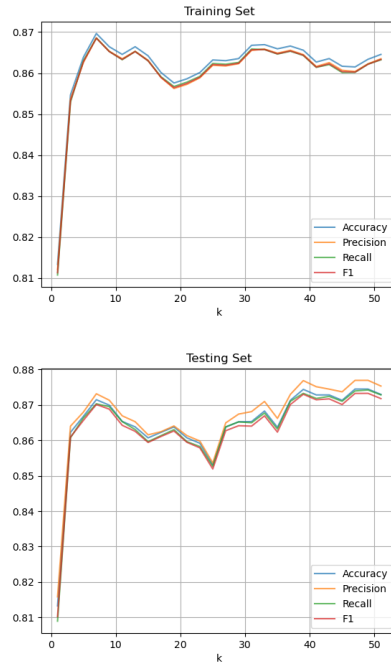
	Accuracy	Precision	Recall	F1
1	0.813338	0.812113	0.810749	0.811312
3	0.854688	0.853702	0.852928	0.853262
5	0.863875	0.862587	0.863030	0.862767
7	0.869662	0.868591	0.868548	0.868516
9	0.866428	0.865320	0.865268	0.865251
11	0.864556	0.863498	0.863267	0.863329
13	0.866429	0.865313	0.865270	0.865254
15	0.864214	0.863138	0.862988	0.863006
17	0.860135	0.858920	0.859097	0.858951
19	0.857581	0.856233	0.856728	0.856435
21	0.858603	0.857239	0.857791	0.857475
23	0.860133	0.858840	0.859191	0.858977
25	0.863196	0.861898	0.862343	0.862083

27	0.863027	0.861745	0.862159	0.861908
29	0.863536	0.862292	0.862591	0.862398
31	0.866773	0.865540	0.865873	0.865671
33	0.866941	0.865859	0.865738	0.865755
35	0.865919	0.864867	0.864642	0.864707
37	0.866602	0.865563	0.865331	0.865398
39	0.865579	0.864518	0.864300	0.864366
41	0.862688	0.861616	0.861399	0.861455
43	0.863538	0.862558	0.862080	0.862264
45	0.861669	0.860689	0.860146	0.860366
47	0.861495	0.860396	0.860179	0.860244
49	0.863366	0.862263	0.862179	0.862159
51	0.864555	0.863498	0.863235	0.863323

Mean Test Metrics:

	Accuracy	Precision	Recall	F1
1	0.813260	0.815834	0.808964	0.810260
3	0.862265	0.864002	0.860729	0.860907
5	0.866787	0.867982	0.866320	0.865640
** 7	0.871474	0.873149	0.870328	0.870131
9	0.869937	0.871333	0.869495	0.868801
11	0.865344	0.866925	0.865289	0.864250
13	0.863805	0.865266	0.863049	0.862544
15	0.860703	0.861553	0.859636	0.859405
17	0.862265	0.862438	0.861342	0.861099
19	0.863781	0.864069	0.863027	0.862644
21	0.860703	0.861345	0.859636	0.859454
23	0.859188	0.859809	0.858265	0.857887
25	0.853080	0.853640	0.852670	0.851925
27	0.863733	0.864932	0.863820	0.862728
29	0.865248	0.867416	0.865209	0.864131
31	0.865296	0.868116	0.864971	0.864019
33	0.868279	0.870985	0.867692	0.866916
35	0.863686	0.866188	0.863190	0.862329
37	0.871309	0.872981	0.871025	0.870122
39	0.874387	0.876868	0.873269	0.873002
41	0.872824	0.875145	0.871840	0.871453
43	0.872824	0.874460	0.872414	0.871711
45	0.871309	0.873698	0.871025	0.870110
47	0.874481	0.876941	0.873940	0.873251
49	0.874481	0.876948	0.874217	0.873274
51	0.872966	0.875330	0.872828	0.871788

(\*\* indicates the optimal tuning parameters)



We found the Testing set resemble the training set in accuracy, precision, recall, and F1. This further boost the confidence of reliable results.

As we can see,  $k$ -NN improves performs from  $k=1$  tp 7, as the model seek closeness from kneighboring data points. All accuracy, precision, recall, and F1 peaks at  $k=7$ , indicating maximum signal when  $k=7$ . As more data points included, the signal is diluted, and performance decline, until  $k=20$ . Performances slowly climbs afterwards for both training and testing sets. It never reach previous peaks on training set.

How to pick the best  $k$  values? I would like to pick  $k= 5, 7$ , and  $9$ . From charts, all accuracy, precision, recall, and F1 are about 0.87, for training set. Such results are confined in testing set. For credit approval application, false positive means big financial loss from persons who can not pay back. So high precision is very beneficial. The highest F1 near 0.87 shows robustness of these result.

Ideally, I would peak the optimal  $k = 7$ , exactly at the performance peak, with surrounding performance at  $k = 5$  and  $9$  also stable and good.

### 3.3 Random Forest

```
$ python random_forest.py --help
```

```
usage: random_forest.py [-h] [-p PATH] [-ntrees NUM_TREES] [-dmin MIN_DEPTH]
                        [-dmax MAX_DEPTH] [-smin MIN_SAMPLES] [-igmin MIN_INFO_GAIN]
                        [-skip_norm SKIP_NORMALIZATION] [-bst ISBST] [-gini USE_GINI]
                        [-early_stop EARLY_STOP_THRESHOLD] [--random_state RANDOM_STATE]
```

optional arguments:

```
-h, --help            show this help message and exit
-p PATH, --path PATH  csv file path, like wdbc.csv
-ntrees NUM_TREES, --num_trees NUM_TREES
                        number of trees if it is greater than zero, 0 to loop over [1 5 10 20 30 40 50], less than 0 to loop over [1 5 10 20 30 40 50 75 100 125 150 200 250 300], default 0
-dmin MIN_DEPTH, --min_depth MIN_DEPTH
                        Min depth of tree, default 2
-dmax MAX_DEPTH, --max_depth MAX_DEPTH
                        Max depth of tree, default 8
-smin MIN_SAMPLES, --min_samples MIN_SAMPLES
                        Min samples for split, default 2
-igmin MIN_INFO_GAIN, --min_info_gain MIN_INFO_GAIN
                        Minimum information gain, default 0.001
-skip_norm SKIP_NORMALIZATION, --skip_normalization SKIP_NORMALIZATION
                        skip simulation, default True
-bst ISBST, --isBST ISBST
                        Use numerical splitting, default False
-gini USE_GINI, --use_gini USE_GINI
                        Use Gini coefficient, default False
-early_stop EARLY_STOP_THRESHOLD, --early_stop_threshold EARLY_STOP_THRESHOLD
                        Early majority stop threshold like 0.85, default 1.00
--random_state RANDOM_STATE
                        random seed like 42, -5, default None
```

Use use -ntrees -1 to let the random forest algorithms to loop over [1 5 10 20 30 40 50 75 100 125 150 200 250 300] trees.

To handle the large computational demand, multiprocessing are used to use multicores (overcome python GIL lock). multiprocessing queues are used to feed workers with tasks, and collect results. For details, please see code.

```
$ python random_forest.py -bst true --random_state 42 -p credit_approval_adj.csv
```

Random Forest Performance:

	Accuracy	Precision	Recall	F1
1	0.701682	0.704034	0.697459	0.696212
5	0.774921	0.780442	0.768733	0.769602

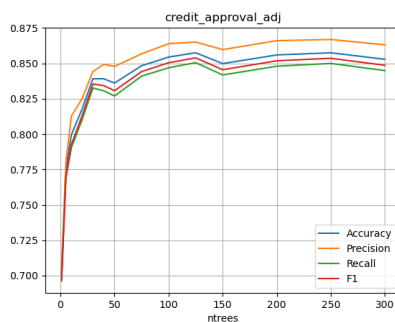
10	0.799349	0.812550	0.790203	0.792452
20	0.817718	0.825202	0.810173	0.812904
30	0.839116	0.844212	0.832614	0.835392
40	0.839116	0.849252	0.830607	0.834262
50	0.836086	0.847841	0.826931	0.830663

All our algorithms are properly handling random state, making results repeatable.

```
$ python random_forest.py -bst true -ntrees -1 --random_state 42 \
    -p credit_approval_adj.csv
```

Random Forest Performance:

	Accuracy	Precision	Recall	F1
1	0.701682	0.704034	0.697459	0.696212
5	0.774921	0.780442	0.768733	0.769602
10	0.799349	0.812550	0.790203	0.792452
20	0.817718	0.825202	0.810173	0.812904
30	0.839116	0.844212	0.832614	0.835392
40	0.839116	0.849252	0.830607	0.834262
50	0.836086	0.847841	0.826931	0.830663
75	0.848300	0.856746	0.840959	0.844188
100	0.854431	0.863912	0.846859	0.850430
125	0.857484	0.865068	0.850527	0.853893
150	0.849827	0.859689	0.841794	0.845493
200	0.855934	0.866037	0.847973	0.851810
250	0.857484	0.866872	0.849954	0.853587
300	0.852869	0.863082	0.844869	0.848623



We also explored the Gini index to see whether it offers an improvement over entropy in picking the best features and values to split.

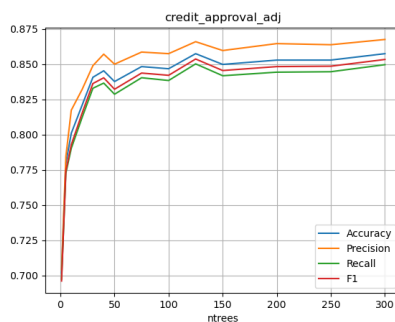
```
$ python random_forest.py -bst true -ntrees -1 -gini true \
    --random_state 42 -p credit_approval_adj.csv
```



# Random Forest Performance:

	Accuracy	Precision	Recall	F1
1	0.701705	0.704163	0.697487	0.696098
5	0.779490	0.784502	0.772624	0.774221
10	0.800970	0.817330	0.790218	0.792797
20	0.820830	0.832013	0.812121	0.815204
30	0.840701	0.848978	0.832905	0.836265
40	0.845305	0.857116	0.836575	0.840332
50	0.837659	0.849999	0.828665	0.832200
75	0.848289	0.858594	0.840377	0.843701
100	0.846762	0.857467	0.838404	0.842091
** 125	0.857461	0.865968	0.850212	0.853666
150	0.849827	0.859689	0.841794	0.845493
200	0.852904	0.864611	0.844324	0.848299
250	0.852916	0.863803	0.844639	0.848562
300	0.857461	0.867535	0.849639	0.853345

(\*\* indicates the optimal tuning parameters)



Both entropy and the Gini index offer very similar performances. They all start with low accuracy, precision, recall, and F1 for low ‘ntrees’ counts. This is expected, as a random forest relies on collecting the wisdom of weak learners. Each decision tree offers better than random performance, but it is far from stellar. As ‘ntrees’ reaches 40, performance is significantly improved, with accuracy 0.841 and F1 0.840. The model continues to improve, synergizing the knowledges of more decision trees, until ntrees around 125. This is the point all metrics peak together, with accuracy 0.857, precision 0.865, recall 0.850, and F1 0.854. Then the performance more or less is flat. This is the point of diminishing return for more computing.

The chart by the Gini index is similar with small variations. This confirmation is welcome in our analysis.

Since the performance is stable and very flat from ‘ntrees’ 100 to 300, we would pick the peak performance near the low end ‘ntrees = 125’ as the optimal hyperparameter for the random forest algorithm for highest performance.

### 3.4 NN

We should caution against overfitting with deep and dense neural network. We only have 584 data samples. Considering there are 46 features in the data after one-hot encoding, if the first layer has 16 neurons, it takes  $46 \times 16 = 736$  parameters, which is more than our data points.

```
$ python nn.py --help
usage: nn.py [-h] [-p PATH] [-input NUM_INPUT] [-output NUM_OUTPUT]
             [-neuron HIDDEN_NEURONS] [-activation ACTIVATION] [-batch BATCH_SIZE]
             [-lr LEARNING_RATE] [-rlambda RLAMBDA]
             [-loss_delta LOSS_DELTA_THRESHOLD] [-k K_EPOCH_SHUFFLE]
             [-epoch EPOCHS] [-kfold NUM_KFOLDS] [-random_state RANDOM_STATE]
             [-proj PROJECT_NAME_PREFIX]
```

optional arguments:

```
-h, --help            show this help message and exit
-p PATH, --path PATH  csv file path, like wdbc.csv
-input NUM_INPUT, --num_input NUM_INPUT
                      number of input, default 5
-output NUM_OUTPUT, --num_output NUM_OUTPUT
                      number of output, default 1
-neuron HIDDEN_NEURONS, --hidden_neurons HIDDEN_NEURONS
                      Hidden layer neurons, default [10]
-activation ACTIVATION, --activation ACTIVATION
                      Non-output activation function, default sigmoid
-batch BATCH_SIZE, --batch_size BATCH_SIZE
                      batch size, default 16
-lr LEARNING_RATE, --learning_rate LEARNING_RATE
                      Learning rate, default 0.01
-rlambda RLAMBDA, --rlambda RLAMBDA
                      Regularization lambda, default 0.01
-loss_delta LOSS_DELTA_THRESHOLD, --loss_delta_threshold LOSS_DELTA_THRESHOLD
                      Early stop loss limit, default 0.001
-k K_EPOCH_SHUFFLE, --k_epoch_shuffle K_EPOCH_SHUFFLE
                      Shuffle training data per k epochs, default -1 (never)
-epoch EPOCHS, --epochs EPOCHS
                      number of epochs, default 1000
-kfold NUM_KFOLDS, --num_kfolds NUM_KFOLDS
                      number of folds for stratified K-Fold, default 5
-random_state RANDOM_STATE, --random_state RANDOM_STATE
                      random seed like 42, -5, default None
-proj PROJECT_NAME_PREFIX, --project_name_prefix PROJECT_NAME_PREFIX
                      Project name prefix, like myproject, default proj
```

We want to build a simple network with strong performance, rather than a complex network of over curve-fitting. So we start explore network with single

hidden layer with low neuron count, then two hidden layers, then three hidden layers. We shall see that a moderate network is adequate.

The tuning of learning rate and regularization rate are key to performing model.

```
$ python nn.py --hidden_neurons "16 16" -lr 0.01 -rlambda 0.01 \
--batch_size 16 --loss_delta_threshold 0.001 \
--k_epoch_shuffle -1 -kfold 5 --random_state 42 \
-p credit_approval_adj.csv
```

—	NN	—	batch size	—	lr	—	rlambda	—	loss	—	Accuracy	—	Precision	—
Recall	F1	—	:	:	—	:	:	—	:	—	:	:	—	:
46, [2], 1	16	—	0.01	—	0.01	—	0.35041	—	0.86066	—	0.86090	—	0.86199	—
0.85989	—	46, [4], 1	64	—	0.04	—	0.04	—	0.36057	—	0.85597	—	0.85635	—
0.85680	—	0.85511	—	46, [8], 1	16	—	0.01	—	0.01	—	0.36485	—	0.86989	—
0.87063	—	0.87044	—	0.86901	—	46, [8], 1	64	—	0.05	—	0.1	—	0.34365	—
0.86537	—	0.86711	—	0.86857	—	0.86491	—	46, [16], 1	64	—	0.01	—	0.05	—
0.35033	—	0.86677	—	0.86722	—	0.86844	—	0.86611	—	46, [64], 1	64	—	0.05	—
0.05	—	0.05	—	0.36635	—	0.86544	—	0.86644	—	0.86581	—	0.86441	—	46, [2, 2], 1
64	—	0.04	—	0.04	—	0.34602	—	0.86539	—	0.86536	—	0.86625	—	0.86461
—	46, [4, 2], 1	64	—	0.04	—	0.04	—	0.37200	—	0.86229	—	0.86224	—	0.86223
—	0.86119	—	46, [4, 4], 1	64	—	0.04	—	0.04	—	0.36472	—	0.86253	—	0.86353
—	0.86199	—	0.86138	—	46, [4, 8], 1	32	—	0.05	—	0.1	—	0.34901	—	0.88234
—	0.88544	—	0.88674	—	0.88218	—	46, [4, 8], 1	64	—	0.05	—	0.1	—	0.34754
—	0.86385	—	0.86599	—	0.86490	—	0.86285	—	46, [8, 8], 1	64	—	0.05	—	0.05
—	0.35577	—	0.86703	—	0.86706	—	0.86748	—	0.86608	—	46, [8, 8], 1	16	—	0.01
—	0.01	—	0.38466	—	0.85460	—	0.85632	—	0.85381	—	0.85327	—	46, [16, 8], 1	16
—	0.04	—	0.04	—	0.35071	—	0.87467	—	0.87804	—	0.87884	—	0.87446	—
—	46, [16, 8], 1	16	—	0.06	—	0.04	—	0.35512	—	0.87145	—	0.87665	—	0.87558
—	0.87109	—	46, [16, 8], 1	32	—	0.06	—	0.04	—	0.35298	—	0.86994	—	0.87336
—	0.87335	—	0.86959	—	46, [16, 8], 1	64	—	0.06	—	0.02	—	0.44303	—	0.84996
—	0.85180	—	0.84806	—	0.84801	—	46, [16, 16], 1	64	—	0.04	—	0.04	—	0.37055
—	0.86075	—	0.86086	—	0.86114	—	0.85976	—	46, [32, 32], 1	64	—	0.04	—	0.04
—	0.38846	—	0.85931	—	0.86097	—	0.85849	—	0.85788	—	46, [16, 16], 1	64	—	0.03
—	0.02	—	0.51221	—	0.84536	—	0.84673	—	0.84395	—	0.84350	—	46, [16, 16], 1	64
—	0.06	—	0.02	—	0.44872	—	0.84844	—	0.84795	—	0.84850	—	0.84742	—
—	46, [32, 32], 1	64	—	0.04	—	0.04	—	0.38846	—	0.85931	—	0.86097	—	0.85849
—	0.85788	—	46, [4, 4, 4],											

For single layer, as we can see the performance metrics are similar for neurons 2, 4, 8, 16 and 32. They all have low loss around 0.35, accuracy 0.86, precision 0.86, recall 0.86, and F1 0.86. This hint simple network with low neurons. It

might suggest more than 8 neurons are unnecessary complexity, yet less than 2 is unnecessary limiting.

From the above table, we can see that network with two hidden layers are in general not much better performing than the single hidden layer. Some more complexed models performs a lot worse.

Three hidden layers does not help performance. It seems brittle: with small changes and the performances tanks. This is like due to overfitting.

We are mostly interesting in simple network with two hidden layers, for these reasons: it is as performing as single layer but offer a lot more options to tune. In particular, we pick hidden layer [8, 2] network to further tune for crank up performance.

### 3.4.1 Tune hyperparameters

Since the dataset is credit approval, F1 and accuracy is of ultimate importance, because it could mean that the business could avoid unnecessary loss for bad approvals, or missed profit for deny qualified applicants. For the same or close accuracy and F1 values, we would like the ones with lower loss because the higher confidence it implied.

NN	batch size	lr	rlambda	loss	Accuracy	Precision	Recall	F1
46, [4, 8], 1	32	0.05	0.1	0.34780	0.88083	0.88410	0.88535	0.88068
46, [4, 8], 1	32	0.05	0.05	0.34602	0.88083	0.88367	0.88479	0.88059
46, [4, 8], 1	32	0.025	0.05	0.34648	0.86996	0.87093	0.87282	0.86955
46, [4, 8], 1	32	0.04	0.05	0.34622	0.87467	0.87684	0.87831	0.87436
46, [4, 8], 1	32	0.05	0.025	0.38195	0.85453	0.85756	0.85606	0.85352
46, [4, 8], 1	32	0.06	0.04	0.35022	0.87609	0.87824	0.87928	0.87577
46, [4, 8], 1	32	0.06	0.12	0.35361	0.87614	0.88020	0.88104	0.87601
* 46, [4, 8], 1	32	0.04	0.075	0.34435	0.88083	0.88347	0.88479	0.88060
46, [4, 8], 1	32	0.035	0.08	0.34396	0.87770	0.88006	0.88193	0.87746
46, [4, 8], 1	32	0.045	0.09	0.34624	0.88083	0.88382	0.88507	0.88063
46, [4, 8], 1	32	0.04	0.09	0.34531	0.88083	0.88347	0.88479	0.88060
46, [4, 8], 1	32	0.04	0.07	0.34423	0.87927	0.88157	0.88307	0.87903

(\*\* indicates the optimal tuning parameters)

```
$ python nn.py --hidden_neurons "4 8" -lr 0.04 -rlambda 0.075 \
--batch_size 32 --loss_delta_threshold 0.0001 \
--k_epoch_shuffle -1 -kfold 10 --random_state 42 -p parkinsons.csv
```

Stratified K-Fold losses:

[0.3301494868319709, 0.24885913737392373, 0.4531295853945854, 0.29652771447233056, 0.3964249  
mean loss: 0.3443485917243486

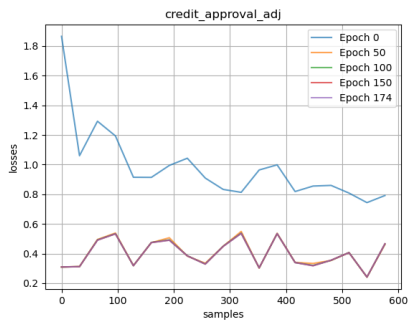
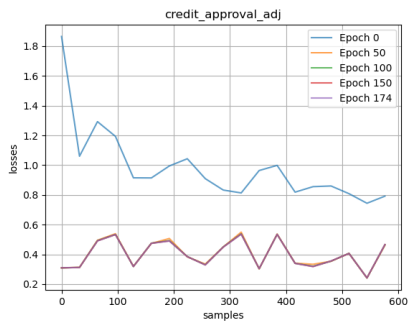
Stratified K-Fold Performances:

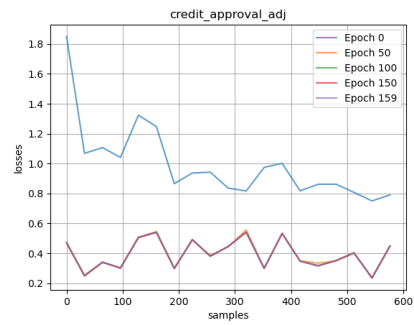
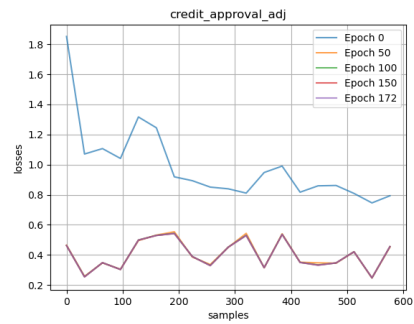
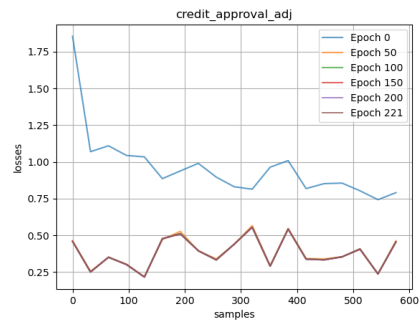
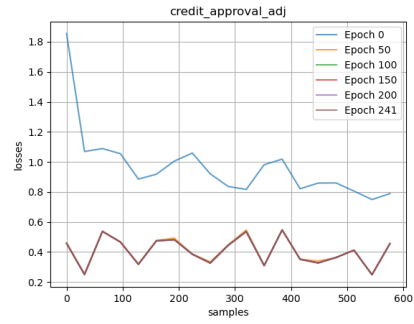
	Accuracy	Precision	Recall	F1
0	0.878788	0.880515	0.883333	0.878676
1	0.939394	0.938419	0.941667	0.939171

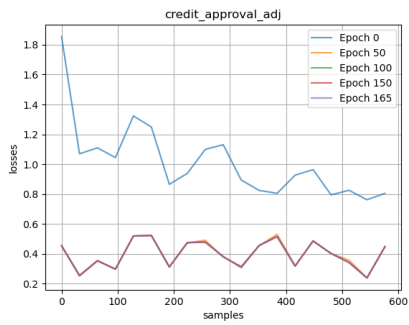
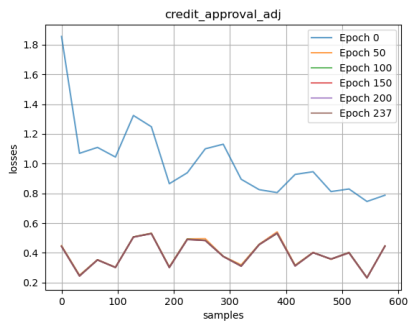
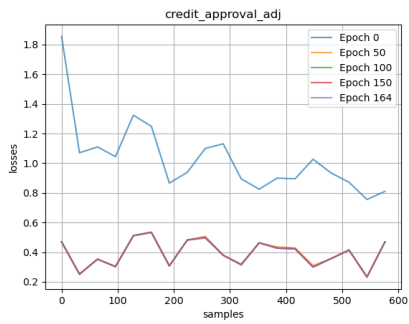
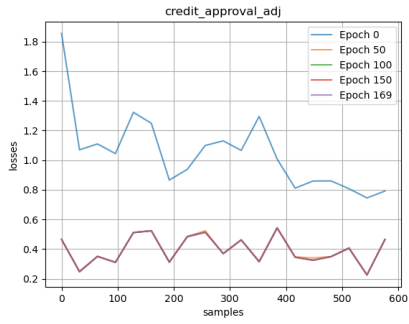
2	0.818182	0.817096	0.819444	0.817512
3	0.893939	0.893939	0.897222	0.893720
4	0.833333	0.837327	0.838889	0.833295
5	0.848485	0.847426	0.850000	0.847926
6	0.846154	0.860521	0.857759	0.846117
7	0.890625	0.890625	0.894089	0.890384
8	0.875000	0.882759	0.882759	0.875000
9	0.984375	0.986111	0.982759	0.984186

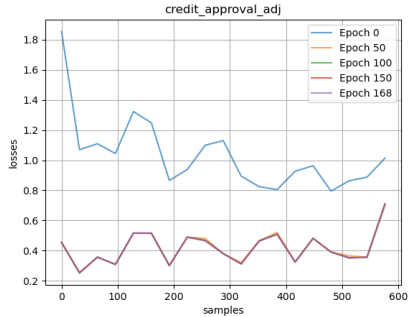
Stratified K-Fold Mean Performances:

Accuracy	0.880828
Precision	0.883474
Recall	0.884792
F1	0.880599









### 3.5 Summary

All KNN, random forest, and NN delivered solid results:

KNN: accuracy 0.8715 F1 0.8701  
RD: accuracy 0.8575 F1 0.8537  
NN: accuracy 0.8808 F1 0.8806

All accuracy  $> 0.85$ , all F1  $> 0.85$ . The consistency boost our confidence in the results.

RF show remarkable resistance to noise, with stable and consistent high performance for  $ntrees \geq 125$ .

NN is the top performer here. We know NN shines on large data set. But we only have 653 data samples here. This shows the power of NN and proper hyperparameter tuning.

KNN performers well here too. This might due to the clean data set in credit approval dataset, despite those data tends to be noisy. All features are more or less equally important. And, it has moderate demension. This is where KNN shines.

## 4 Parkinsons Dataset

This section is implemented by William Cai, including the code and the writeup.

### 4.1 Nature of the Dataset and Algorithms

The Oxford Parkinson's Diseases Detection Dataset is composed of a range of biomedical voice measurements of people with or without Parkinson's disease. It has 22 numerical features. The goal is to detect whether a person is healthy or with Parkinson's disease.

Because the dataset contains high-dimensional data, traditional data analysis like regressions are ill-suited for the tasks. We should explore multiple powerful machine learning algorithms to address this problem.



The  $k$ -NN algorithm is good at finding the similarity between instances by calculating their Euclidean distances in a high-dimensional space. This method is effective, since we suspect that all features about the voice attributes are of similar importance. At least, when we do not know about the relative importance of features, treating them equally is a good baseline starting point. The drawback of  $k$ -NN is that the dataset being super high in dimension means that there are many "useless" features that mislead  $k$ -NN to a dead end. But this is not our case: 22 features are moderate, and they are all clearly relating to Parkinson's (or the absence thereof), by researchers.

The random forest algorithm could be suitable for this problem. It is built on many decision trees and uses bootstrapping and feature sampling to get robust results. If there are some features that are more important in classifying the disease, the decision tree will find out such features by obtaining the larger information gain via entropy or the Gini index. It compliments  $k$ -NN's treating features equally. This allows it to be resistant to some noise in the data. The drawback is that it could be hard to interpret than a simple decision tree, and it is often computationally expensive. We can alleviate this by multiprocessing to use all CPU cores.

We also explore neural networks for this dataset. With adequately sophisticated neural network configuration and non-linear activation functions like Sigmoid or ReLU at each neuron, and enough training data points, a NN can learn and potentially approximate any curves, and discover the hidden rules if there are any. This opens new ways for us to let the neural network learn from data by adjusting weights and biases in an automated fashion, without us to guess how to pick features. This powerful method comes with potential over curve-fitting, so we thus deploy regularization to mitigate such risks. Most complex NNs need intensive parallel computing, which may require RTX A6000, A100, H100, or H200. Fortunately, our dataset is small or moderate, and our CPU is quite capable of doing it.

## 4.2 Algorithms

### 4.2.1 $k$ -NN

The features are first scaled to  $[0, 1]$  using our own MinMaxScaler class before feeding them to the  $k$ -NN. This is a straightforward but important step to avoid meaningless or misleading scaling or unit distortions among features in calculating the distance between two instances, i.e., how similar and close they are.

Our  $k$ -NN code support many options as shown below.

```
$ python knn_cv.py --help
usage: knn_cv.py [-h] [-p PATH] [-n ROUNDS] [-kmin KMIN] [-kmax KMAX]
                  [-exclude_self EXCLUDE_SELF] [-skip_norm SKIP_NORMALIZATION]
                  [-header HEADER] [-kfold NUM_KFOLDS] [-random_state RANDOM_STATE]
```

optional arguments:

```

-h, --help          show this help message and exit
-p PATH, --path PATH CSV path of data file, default wdbc_wo_header.csv
-n ROUNDS, --rounds ROUNDS
                    Rounds for KNN run, default 20
-kmin KMIN, --kmin KMIN
                    Min K in KNN, default 1
-kmax KMAX, --kmax KMAX
                    Max K for KNN, default 52
-exclude_self EXCLUDE_SELF, --exclude_self EXCLUDE_SELF
                    Exclude self in training, default False
-skip_norm SKIP_NORMALIZATION, --skip_normalization SKIP_NORMALIZATION
                    Skip normalization, default False
-header HEADER, --header HEADER
                    CSV file header or not, default None
-kfold NUM_KFOLDS, --num_kfolds NUM_KFOLDS
                    number of folds for stratified K-Fold, default 5
-random_state RANDOM_STATE, --random_state RANDOM_STATE
                    random seed like 42, -5, default None

```

The following command runs  $k$ -NN with our own stratified  $k$ -fold with shuffling.

```
$ python knn_cv.py --header true --exclude_self false --num_kfolds 10 -p parkinsons.csv
```

Stratified Results:

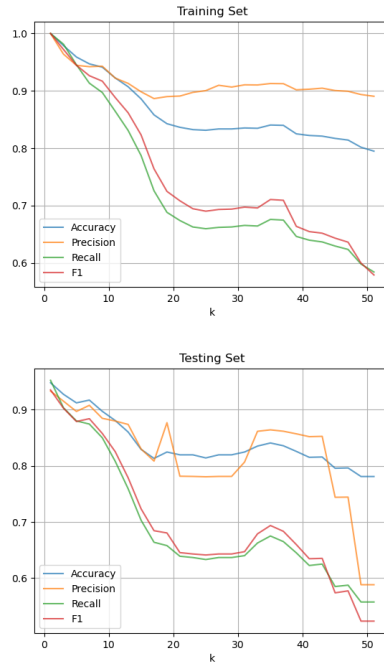
Mean Train Metrics:

	Accuracy	Precision	Recall	F1
1	1.000000	1.000000	1.000000	1.000000
3	0.978925	0.964033	0.981339	0.972133
5	0.958986	0.944525	0.945483	0.944647
7	0.947025	0.942123	0.913403	0.926396
9	0.941321	0.942955	0.897139	0.916935
11	0.922512	0.922249	0.864443	0.888386
13	0.907148	0.913029	0.830931	0.861749
15	0.886050	0.898369	0.787209	0.823180
17	0.858118	0.886512	0.725786	0.764043
19	0.842741	0.890008	0.688295	0.724802
21	0.836469	0.890918	0.673992	0.708410
23	0.832472	0.897354	0.662788	0.694616
25	0.831342	0.900377	0.659757	0.690288
27	0.833621	0.909732	0.662077	0.693267
29	0.833618	0.906606	0.662838	0.693937
31	0.835326	0.910511	0.665539	0.697385
33	0.834767	0.910267	0.664429	0.695919

35	0.840456	0.912796	0.675951	0.710637
37	0.839881	0.912527	0.674762	0.709311
39	0.825076	0.901811	0.646309	0.663870
41	0.822238	0.902895	0.639790	0.654743
43	0.821108	0.904649	0.636734	0.651953
45	0.817140	0.900572	0.629537	0.643290
47	0.814270	0.899253	0.623643	0.636318
49	0.801747	0.893592	0.598194	0.599940
51	0.794906	0.890550	0.584294	0.579130

Mean Test Metrics:

	Accuracy	Precision	Recall	F1
1	0.947807	0.933036	0.952143	0.935038
3	0.927251	0.915000	0.903214	0.902480
5	0.911988	0.896712	0.879881	0.878833
7	0.916959	0.907635	0.874286	0.883803
9	0.896988	0.884510	0.849881	0.857592
11	0.880585	0.879688	0.807976	0.825155
13	0.860058	0.873578	0.759048	0.778530
15	0.829240	0.829874	0.702976	0.723322
17	0.813129	0.808335	0.663810	0.684416
19	0.824503	0.876722	0.657857	0.680398
21	0.819503	0.781419	0.638929	0.645394
23	0.819503	0.781072	0.636667	0.642997
25	0.813947	0.780418	0.633095	0.641183
27	0.819503	0.781072	0.636667	0.642997
29	0.819503	0.781072	0.636667	0.642997
31	0.824503	0.806656	0.640000	0.647072
33	0.835058	0.861395	0.662500	0.678832
35	0.840614	0.863968	0.675000	0.693671
37	0.835614	0.861517	0.665000	0.683322
39	0.825614	0.856873	0.645000	0.659731
41	0.815058	0.851849	0.622500	0.634543
43	0.815614	0.852421	0.625000	0.635135
45	0.795614	0.743896	0.585000	0.573803
47	0.796170	0.744210	0.587500	0.577287
49	0.780906	0.588216	0.557500	0.523440
51	0.780906	0.588216	0.557500	0.523440



As we can see in the above chart, the training set achieves accuracy 100% when  $k = 1$ . This is because the model is comparing a data point itself (the closest distance is zero) for classification, hence the `--exclude_self false` option. This distortion continues for small  $k$  values. This self-included effect gradually fades away for large  $k$ .

Let me exclude the data point itself in the classification by setting `--exclude_self true`.

```
$ python knn_cv.py --header true --exclude_self true --num_kfolds 10 -p parkinsons.csv
```

Stratified Results:

Mean Train Metrics:

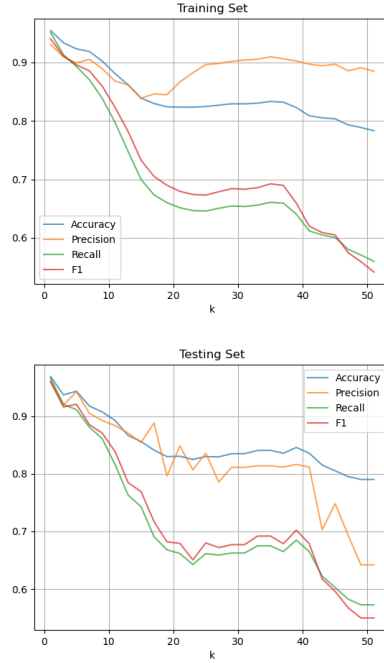
	Accuracy	Precision	Recall	F1
1	0.954421	0.930580	0.951038	0.940060
3	0.933330	0.909115	0.912898	0.910676
5	0.923077	0.898944	0.892860	0.895655
7	0.918522	0.905150	0.870353	0.885373
9	0.901989	0.889059	0.838336	0.859139
11	0.880343	0.867794	0.796666	0.822646
13	0.861534	0.861212	0.747568	0.781513
15	0.838771	0.838938	0.699802	0.732706
17	0.829621	0.845800	0.673396	0.705025
19	0.823920	0.844723	0.660274	0.689796

21	0.823365	0.866482	0.651374	0.679489
23	0.823365	0.882016	0.646667	0.674076
25	0.824501	0.896163	0.645854	0.673138
27	0.826793	0.897953	0.650531	0.678705
29	0.829079	0.901527	0.654399	0.684094
31	0.829063	0.903987	0.653535	0.683088
33	0.830202	0.905246	0.655858	0.685654
35	0.833053	0.909434	0.660888	0.692374
37	0.831891	0.906024	0.659267	0.689606
39	0.822751	0.902289	0.640662	0.659229
41	0.808553	0.896764	0.611883	0.620037
43	0.805121	0.894006	0.604907	0.608932
45	0.803423	0.896831	0.600740	0.605018
47	0.793173	0.885580	0.580673	0.574719
49	0.788604	0.890632	0.570587	0.559031
51	0.782926	0.884918	0.559875	0.541670

Mean Test Metrics:

	Accuracy	Precision	Recall	F1
1	0.969181	0.961131	0.966190	0.959954
** 3	0.937251	0.919702	0.920714	0.916103
5	0.943626	0.942589	0.911429	0.920953
7	0.917807	0.905189	0.881190	0.885579
9	0.907544	0.892630	0.861190	0.870721
11	0.892515	0.883787	0.815595	0.838304
13	0.866696	0.870858	0.763095	0.784677
15	0.856140	0.854596	0.742857	0.769041
17	0.841170	0.888424	0.690833	0.717167
19	0.830058	0.796377	0.668095	0.681895
21	0.830614	0.848583	0.661667	0.679328
23	0.825058	0.806970	0.642500	0.650556
25	0.830058	0.835899	0.661429	0.680047
27	0.829503	0.785805	0.658929	0.671878
29	0.835058	0.811356	0.662500	0.677039
31	0.835058	0.811356	0.662500	0.677039
33	0.840614	0.813930	0.675000	0.691878
35	0.840877	0.814024	0.675000	0.691952
37	0.835877	0.811831	0.665000	0.678710
39	0.845877	0.816475	0.685000	0.702301
41	0.835877	0.811831	0.665000	0.678710
43	0.815322	0.702859	0.622500	0.617668
45	0.805322	0.748254	0.602500	0.595871
47	0.795322	0.693830	0.582500	0.567594
49	0.790322	0.641856	0.572500	0.549667
51	0.790322	0.641856	0.572500	0.549667

(\*\* indicates the optimal tuning parameters)



We found the testing set resembles the training set in accuracy, precision, recall, and F1. This further boosts the confidence of reliable results.

As we can see,  $k$ -NN performs best for  $k < 10$ , then the performance sharply drops for  $k = 10$  20, then stabilizes for  $k = 20$  40, then drops again. It could be that Parkinson's disease marks or absent marks are very notbale, so keeping a small  $k$  is like keeping signal strong, whereas including far away data points dilutes the signal, thus weakening performance.

How to pick the best  $k$  values? I would like to pick  $k = 1, 3$ , or  $5$ . From the charts, all accuracy, precision, recall, and F1 are above 0.90 for the training set. Such results are confirmed in testing set. For medical detection, false negatives are extremely harmful, so recall  $\geq 0.92$  is very beneficial. F1  $\geq 0.91$  shows the robustness of these results.

Due to risk of conincidence, we would favor  $k = 3$ , or  $5$  over  $k = 1$ .

## 4.2.2 Random Forest

```
$ python random_forest.py --help
```

```
usage: random_forest.py [-h] [-p PATH] [-ntrees NUM_TREES] [-dmin MIN_DEPTH]
                        [-dmax MAX_DEPTH] [-smin MIN_SAMPLES] [-igmin MIN_INFO_GAIN]
                        [-skip_norm SKIP_NORMALIZATION] [-bst ISBST] [-gini USE_GINI]
                        [-early_stop EARLY_STOP_THRESHOLD] [--random_state RANDOM_STATE]
```

optional arguments:

```
-h, --help            show this help message and exit
-p PATH, --path PATH  csv file path, like wdbc.csv
-ntrees NUM_TREES, --num_trees NUM_TREES
                        number of trees if it is greater than zero, 0 to loop over [1 5 10 2
                        30 40 50], less than 0 to loop over [1 5 10 20 30 40 50 75 100 125
                        150 200 250 300], default 0
-dmin MIN_DEPTH, --min_depth MIN_DEPTH
                        Min depth of tree, default 2
-dmax MAX_DEPTH, --max_depth MAX_DEPTH
                        Max depth of tree, default 8
-smin MIN_SAMPLES, --min_samples MIN_SAMPLES
                        Min samples for split, default 2
-igmin MIN_INFO_GAIN, --min_info_gain MIN_INFO_GAIN
                        Minimum information gain, default 0.001
-skip_norm SKIP_NORMALIZATION, --skip_normalization SKIP_NORMALIZATION
                        skip simulation, default True
-bst ISBST, --isBST ISBST
                        Use numerical splitting, default False
-gini USE_GINI, --use_gini USE_GINI
                        Use Gini coefficient, default False
-early_stop EARLY_STOP_THRESHOLD, --early_stop_threshold EARLY_STOP_THRESHOLD
                        Early majority stop threshold like 0.85, default 1.00
--random_state RANDOM_STATE
                        random seed like 42, -5, default None
```

Use ‘-ntrees -1’ to let the random forest algorithm loop over [1 5 10 20 30 40 50 75 100 125

To handle the large computational demand, multiprocessing is used to use multicores (overcon

```
$ python random_forest.py -bst true --random_state 42 -p parkinsons.csv
```

Random Forest Performance:

	Accuracy	Precision	Recall	F1
1	0.810688	0.752122	0.756743	0.743992
5	0.861215	0.827278	0.789080	0.803977

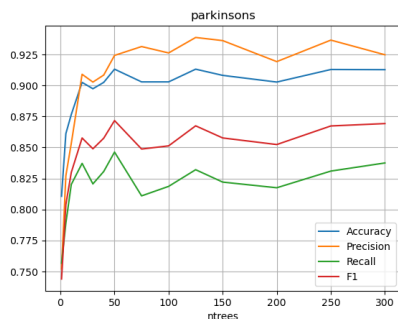
10	0.876606	0.853321	0.820192	0.829893
20	0.902517	0.909036	0.837203	0.857684
30	0.897389	0.902747	0.820651	0.848971
40	0.902517	0.908428	0.830651	0.857514
50	0.913171	0.924270	0.846322	0.871725

All our algorithms are properly handling random state, making results repeatable.

```
$ python random_forest.py -bst true -ntrees -1 --random_state 42 -p parkinsons.csv
```

Random Forest Performance:

	Accuracy	Precision	Recall	F1
1	0.810688	0.752122	0.756743	0.743992
5	0.861215	0.827278	0.789080	0.803977
10	0.876606	0.853321	0.820192	0.829893
20	0.902517	0.909036	0.837203	0.857684
30	0.897389	0.902747	0.820651	0.848971
40	0.902517	0.908428	0.830651	0.857514
50	0.913171	0.924270	0.846322	0.871725
75	0.902908	0.931414	0.810996	0.848773
100	0.902908	0.926258	0.818659	0.851321
125	0.913171	0.938731	0.832107	0.867483
150	0.908171	0.936210	0.822107	0.857708
200	0.902780	0.919270	0.817548	0.852392
250	0.912908	0.936608	0.830996	0.867389
300	0.912780	0.924785	0.837548	0.869305



We also explored the Gini index to see whether it offers improvement over entropy in picking the best features and values to split.

```
$ python random_forest.py -bst true -ntrees -1 -gini true --random_state 42 -p parkinsons.csv
```

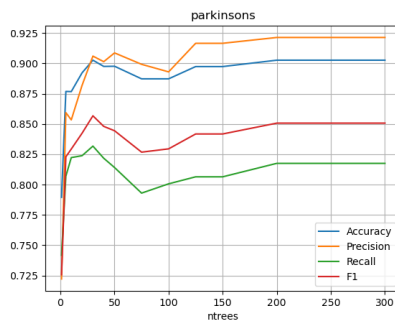
Random Forest Performance:

	Accuracy	Precision	Recall	F1
1	0.789501	0.722170	0.741724	0.725693



5	0.876862	0.859376	0.806973	0.822951
10	0.876734	0.853395	0.822299	0.829350
20	0.892247	0.881943	0.823870	0.842421
30	0.902645	0.906008	0.831762	0.856761
40	0.897517	0.901399	0.821762	0.848041
50	0.897645	0.908642	0.814100	0.844437
75	0.887254	0.899238	0.792989	0.826783
100	0.887254	0.892997	0.800651	0.829498
** 125	0.897382	0.916576	0.806437	0.841780
150	0.897382	0.916576	0.806437	0.841780
200	0.902645	0.921372	0.817548	0.850715
250	0.902645	0.921372	0.817548	0.850715
300	0.902645	0.921372	0.817548	0.850715

(\*\* indicates the optimal tuning parameters)



Both entropy and the Gini index offer very similar performances. They all start with low accuracy, precision, recall, and F1 for low ‘ntrees’ counts. This is expected, as a random forest relies on collection of the wisdom of weak learners. Each decision tree offers better than random performance, but it is far from stellar. As ‘ntrees’ reaches 20, performances are significantly improved, with accuracy 0.90 and F1 0.85. The model continues to improve, synergizing the knowledge of more decision trees, until ‘ntrees’ = around 125. This is the point where all metrics peak together, with accuracy 0.913, precision 0.939, recall 0.832, and F1 0.867. Then the performance more or less is flat. This is the point of diminishing returns for more computing.

The chart by Gini index is similar. It is just smoother, with less variation. This confirmation is welcome in our analysis.

So, we would pick ‘ntrees = 125’ as the optimal hyperparameter for random forest algorithms.

### 4.2.3 Neural Network

```
$ python nn.py --help
usage: nn.py [-h] [-p PATH] [-input NUM_INPUT] [-output NUM_OUTPUT]
             [-neuron HIDDEN_NEURONS] [-activation ACTIVATION] [-batch BATCH_SIZE]
             [-lr LEARNING_RATE] [-rlambda RLAMBDA]
             [-loss_delta LOSS_DELTA_THRESHOLD] [-k K_EPOCH_SHUFFLE]
             [-epoch EPOCHS] [-kfold NUM_KFOLDS] [-random_state RANDOM_STATE]
             [-proj PROJECT_NAME_PREFIX]
```

optional arguments:

```
-h, --help            show this help message and exit
-p PATH, --path PATH  csv file path, like wdbc.csv
-input NUM_INPUT, --num_input NUM_INPUT
                     number of input, default 5
-output NUM_OUTPUT, --num_output NUM_OUTPUT
                     number of output, default 1
-neuron HIDDEN_NEURONS, --hidden_neurons HIDDEN_NEURONS
                     Hidden layer neurons, default [10]
-activation ACTIVATION, --activation ACTIVATION
                     Non-output activation function, default sigmoid
-batch BATCH_SIZE, --batch_size BATCH_SIZE
                     batch size, default 16
-lr LEARNING_RATE, --learning_rate LEARNING_RATE
                     Learning rate, default 0.01
-rlambda RLAMBDA, --rlambda RLAMBDA
                     Regularization lambda, default 0.01
-loss_delta LOSS_DELTA_THRESHOLD, --loss_delta_threshold LOSS_DELTA_THRESHOLD
                     Early stop loss limit, default 0.001
-k K_EPOCH_SHUFFLE, --k_epoch_shuffle K_EPOCH_SHUFFLE
                     Shuffle training data per k epochs, default -1 (never)
-epoch EPOCHS, --epochs EPOCHS
                     number of epochs, default 1000
-kfold NUM_KFOLDS, --num_kfolds NUM_KFOLDS
                     number of folds for stratified K-Fold, default 5
-random_state RANDOM_STATE, --random_state RANDOM_STATE
                     random seed like 42, -5, default None
-proj PROJECT_NAME_PREFIX, --project_name_prefix PROJECT_NAME_PREFIX
                     Project name prefix, like myproject, default proj
```

We want to build a simple network with strong performance, rather than a complex network of over curve-fitting. So we start by exploring a network with a single hidden layer with a low neuron count, then two hidden layers, then three hidden layers. We shall see that a moderate network is adequate.

The tuning of learning rate and regularization rate is key to a performing model.

```
$ python ../hw/nn.py --hidden_neurons "16 16" -lr 0.01 -rlambda 0.01 \
--batch_size 16 --loss_delta_threshold 0.001 \
--k_epoch_shuffle -1 -kfold 5 --random_state 42 -p parkinsons.csv
```

	— NN —	batch size	— lr —	rlambda	— loss —	Accuracy	Precision	—
Recall	—	F1	—	—	—	—	—	—
22, [2], 1	—	16	—	0.01	—	0.01	—	0.33684 — 0.85107 — 0.83281 — 0.75487 —
0.77863	—	22, [4], 1	—	16	—	0.01	—	0.31594 — 0.86173 — 0.84761 —
0.77720	—	0.79734	—	22, [8], 1	—	16	—	0.01 — 0.01 — 0.29868 — 0.87199 —
0.86445	—	0.79065	—	0.81275	—	22, [16], 1	—	16 — 0.01 — 0.01 — 0.32073
— 0.86699	—	0.85735	—	0.78732	—	0.80774	—	22, [32], 1 — 16 — 0.01 —
0.01 — 0.30235	—	0.87199	—	0.86445	—	0.79065	—	0.81275 — 22, [2, 2], 1
— 16 — 0.01 — 0.01	—	0.37030	—	0.85134	—	0.82861	—	0.76264 — 0.78344 —
— 22, [4, 4], 1	—	16	—	0.01	—	0.01	—	0.36277 — 0.84634 — 0.82300 — 0.75165
— 0.77424	—	*	22, [8, 8], 1	—	16	—	0.01	— 0.36429 — 0.84107 —
0.81596	—	0.74820	—	0.76880	—	22, [16, 16], 1	—	16 — 0.01 — 0.01 —
0.33252	—	0.86673	—	0.86217	—	0.77287	—	0.79985 — 22, [32, 32], 1 — 16
— 0.01 — 0.01	—	0.32554	—	0.87199	—	0.87675	—	0.77632 — 0.80660 —
22, [8, 4], 1	—	16	—	0.01	—	0.01	—	0.33749 — 0.85660 — 0.83565 — 0.76609
— 0.78888	—	22, [2, 2, 2], 1	—	16	—	0.01	—	0.01 — 0.39751 — 0.81516
— 0.79722	—	0.69019	—	0.70921	—	22, [4, 4, 4], 1	—	16 — 0.01 — 0.01 —
0.35818	—	0.84634	—	0.82300	—	0.75165	—	0.77424 — 22, [4, 2, 2], 1 — 16
— 0.01 — 0.01	—	0.55770	—	0.75398	—	0.37699	—	0.50000 — 0.42986 — 22,
[4, 4, 2], 1	—	16	—	0.01	—	0.01	—	0.40452 — 0.85686 — 0.89128 — 0.72433 —
0.76005	—	22, [8, 4], 1	—	16	—	0.01	—	0.05 — 0.36773 — 0.85620 — 0.85964
— 0.74398	—	0.77626	—	22, [8, 4], 1	—	16	—	0.05 — 0.05 — 0.36645 —
0.86173	—	0.88341	—	0.73989	—	0.77749	—	22, [8, 4], 1 — 16 — 0.1 — 0.1
— 0.39576	—	0.84107	—	0.91353	—	0.67667	—	0.70438 — 22, [8, 4], 1 —
16 — 0.05 — 0.07	—	0.37681	—	0.85686	—	0.90222	—	0.71667 — 0.75696 —
22, [8, 4], 1	—	16	—	0.03	—	0.2	—	0.55819 — 0.75398 — 0.37699 — 0.50000 —
0.42986	—	22, [8, 4], 1	—	16	—	0.02	—	0.0 — 0.32641 — 0.88726 — 0.88326
— 0.80843	—	0.83154	—	22, [8, 4], 1	—	16	—	0.02 — 0.05 — 0.36456 —
0.86673	—	0.87547	—	0.75854	—	0.79141	—	22, [8], 1 — 16 — 0.01 — 0.01
— 0.32851	—	0.86146	—	0.85513	—	0.76943	—	0.79440 — * 22, [8, 2], 1 —
16 — 0.01 — 0.01	—	0.35864	—	0.84607	—	0.82341	—	0.75153 — 0.77357 —
* 22, [8, 2], 1	—	16	—	0.02	—	0.01	—	0.31955 — 0.86686 — 0.84734 — 0.78720
— 0.80685	—	22, [8, 2], 1	—	16	—	0.04	—	0.02 — 0.34631 — 0.85607 —
0.84009	—	0.75720	—	0.78235	—	22, [8, 2], 1	—	16 — 0.04 — 0.01 — 0.31300
— 0.85621	—	0.82825	—	0.78031	—	0.79290	—	22, [8, 2], 1 — 16 — 0.02 —
0.02 — 0.36044	—	0.85134	—	0.83045	—	0.75498	—	0.77902 — * 22, [8, 2], 1
— 16 — 0.02 — 0.005	—	0.29603	—	0.86173	—	0.84151	—	0.77720 — 0.79624 —
— 22, [8, 2], 1	—	16	—	0.04	—	0.005	—	0.29365 — 0.87700 — 0.84671 — 0.82253
— 0.82731	—	22, [8, 2], 1	—	16	—	0.04	—	0.0025 — 0.29750 — 0.89252 —
0.87518	—	0.83943	—	0.84835	—	22, [8, 2], 1	—	16 — 0.04 — 0.00375 —
0.29504	—	0.88726	—	0.86059	—	0.83598	—	0.84160 — 22, [8, 2], 1 — 16

— 0.08— 0.005 — 0.31504 — 0.86161 — 0.81824 — 0.80475 — 0.80616 — —  
 22, [8, 2], 1 — 16 — 0.04— 0.01 — 0.31300 — 0.85621 — 0.82825 — 0.78031  
 — 0.79290 — — 22, [8, 2], 1 — 16 — 0.04— 0.075 — 0.38195 — 0.85686 —  
 0.90222 — 0.71667 — 0.75696 — — 22, [8, 2], 1 — 16 — 0.04— 0.0075 —  
 0.30041 — 0.86673 — 0.84058 — 0.80153 — 0.80959 — — 22, [8, 2], 1 — 16  
 — 0.03— 0.005 — 0.29198 — 0.87687 — 0.85118 — 0.80820 — 0.82236 — —  
 22, [8, 2], 1 — 16 — 0.03— 0.0025 — 0.28229 — 0.89226 — 0.86751 — 0.83931  
 — 0.84862 — — 22, [8, 2], 1 — 16 — 0.06— 0.0025 — 0.30148 — 0.87713 —  
 0.83884 — 0.82920 — 0.82964 — — 22, [8, 2], 1 — 16 — 0.045— 0.0025 —  
 0.30119 — 0.88726 — 0.86059 — 0.83598 — 0.84160 — — 22, [8, 2], 1 — 16 —  
 0.025— 0.0025 — 0.28486 — 0.87161 — 0.83847 — 0.80475 — 0.81554 — —  
 22, [8, 2], 1 — 16 — 0.035— 0.0025 — 0.28680 — 0.89779 — 0.87810 — 0.85054  
 — 0.85744 — — \*\* 22, [8, 2], 1 — 16 — 0.03— 0.00125 — 0.27804 — 0.89252  
 — 0.86539 — 0.84709 — 0.85054 — — 22, [8, 2], 1 — 16 — 0.06— 0.00125 —  
 0.30785 — 0.87213 — 0.83134 — 0.82586 — 0.82402 — — 22, [8, 2], 1 — 16  
 — 0.045— 0.00125 — 0.30393 — 0.88213 — 0.85116 — 0.83253 — 0.83648 —  
 — 22, [8, 2], 1 — 16 — 0.025— 0.00125 — 0.28118 — 0.87687 — 0.84843 —  
 0.80820 — 0.82193 — — \*\* 22, [8, 2], 1 — 16 — 0.035— 0.00125 — 0.27777 —  
 0.89779 — 0.87810 — 0.85054 — 0.85744 — — \*\* 22, [8, 2], 1 — 16 — 0.035—  
 0.000625 — 0.27635 — 0.89779 — 0.87810 — 0.85054 — 0.85744 — — \*\* 22,  
 [8, 2], 1 — 16 — 0.035— 0.000325 — 0.27597 — 0.89779 — 0.87810 — 0.85054  
 — 0.85744 — — 22, [8, 2], 1 — 16 — 0.04— 0.000325 — 0.27383 — 0.89266 —  
 0.86056 — 0.85475 — 0.85359 —

For a single layer, as we can see the performance metrics are similar for neurons 2, 4, 8, 16 and 32. They all have low loss around 0.31, accuracy 0.86, precision 0.85, recall 0.78, and F1 0.80. This hints a simple network with low neurons. It might suggest that more than 8 neurons mean unnecessary complexity, yet less than 2 is unnecessarily limiting.

From the above table, we can see that networks with two hidden layers are in general not much better performing than networks with a single hidden layer. Some more complex models perform a lot worse.

Three hidden layers do not help performance. It seems brittle: with small changes, performance tanks. This is likely due to overfitting.

We are mostly interested in a simple network with two hidden layers, for these reasons: it is as performing as a single layer but offers a lot more options to tune. In particular, we pick hidden layer [8, 2] network to further tune for cranking up performance.

#### 4.2.4 Tuning hyperparameters

kfold = 10

```
$ python nn.py --hidden_neurons "8 2" -lr 0.05 -rlambda 0.000325 --batch_size 16 --loss_delt
```

NN	batch size	lr	rlambda	loss	accuracy	precision	recall	f1
22, [8, 2], 1	16	0.01	0.01	0.34188	0.85673	0.84666	0.76488	0.78372
22, [8, 2], 1	16	0.04	0.000325	0.24287	0.90251	0.87037	0.87762	0.86774
22, [8, 2], 1	16	0.08	0.000325	0.26106	0.91778	0.89453	0.89679	0.88940
22, [8, 2], 1	16	0.03	0.00016	0.25103	0.89196	0.86160	0.85071	0.84491
22, [8, 2], 1	16	0.04	0.00024	0.24345	0.90251	0.87037	0.87762	0.86774
* 22, [8, 2], 1	16	0.05	0.000325	0.23790	0.91278	0.88762	0.87786	0.88027
22, [8, 2], 1	16	0.05	0.00024	0.23961	0.91278	0.88762	0.87786	0.88027
* 22, [8, 2], 1	16	0.05	0.0004	0.23627	0.91278	0.88762	0.87786	0.88027

(\*\* indicates the optimal tuning parameters)

The optimal result is achieved by:

```
$ python nn.py --hidden_neurons "8 2" -lr 0.05 -rlambda 0.0004 \
--batch_size 16 --loss_delta_threshold 0.001 \
--k_epoch_shuffle -1 -kfold 10 --random_state 42 -p parkinsons.csv
```

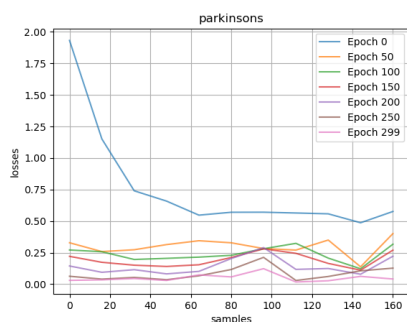
Stratified K-Fold losses:

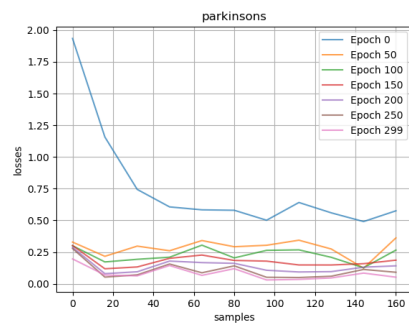
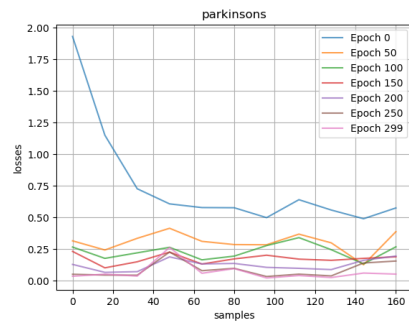
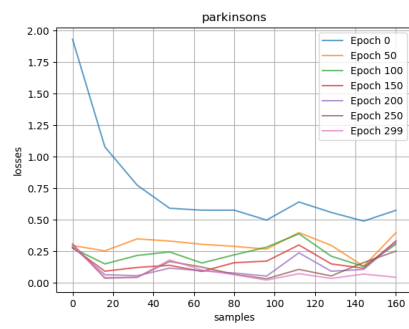
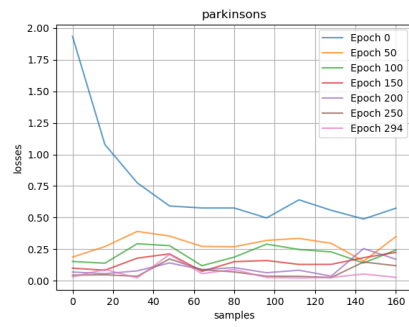
[0.418539881293964, 0.12249701755547995, 0.08748663880663718, 0.11358799573722644, 0.478264]

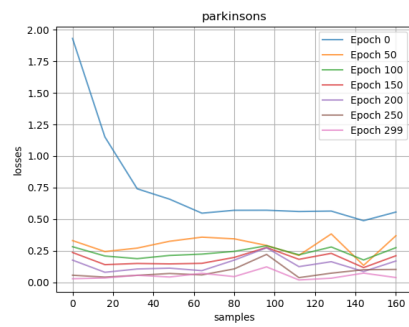
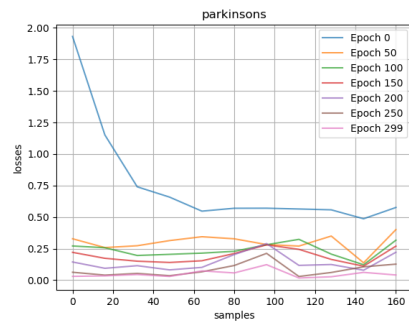
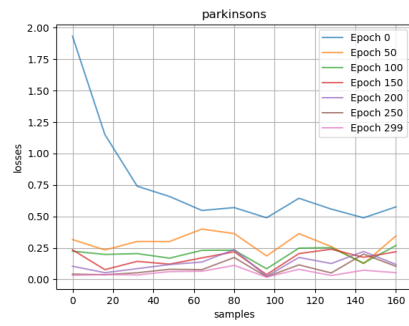
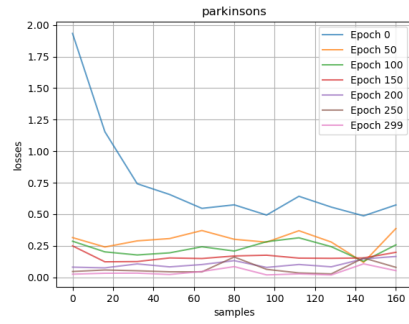
mean loss: 0.2362675003632269

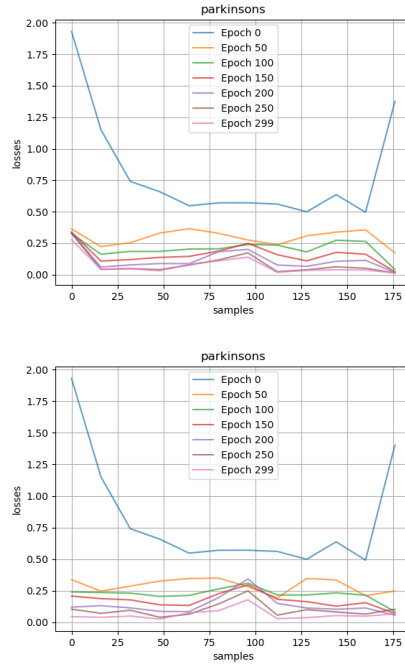
Stratified K-Fold Performances:

	Accuracy	Precision	Recall	F1
0	0.950000	0.968750	0.900000	0.928315
1	0.900000	0.866667	0.866667	0.866667
2	0.950000	0.916667	0.966667	0.937304
3	0.950000	0.968750	0.900000	0.928315
4	0.850000	0.797619	0.833333	0.811912
5	0.850000	0.812500	0.766667	0.784946
6	0.900000	0.866667	0.866667	0.866667
7	1.000000	1.000000	1.000000	1.000000
8	0.888889	0.839286	0.839286	0.839286
9	0.888889	0.839286	0.839286	0.839286









#### 4.2.5 Summary

All KNN, random forest, and NN delivered solid results:

KNN: accuracy 0.9373 F1 0.9161  
RD: accuracy 0.8974 F1 0.8418  
NN: accuracy 0.8128 F1 0.8803

All accuracy  $> 0.89$ , all F1  $> 0.84$ . The consistency boost our confidence in the results.

RF show remarkable resistance to noise, with stable and consistent high performance for  $ntrees \geq 125$ .

KNN is the top performer here. This might due to the clean data set in Parkinsons dataset. All features are more or less equally important. And, it has moderate demension. This is where KNN shines.

NN performs well too. We know NN shines on large data set. But we have less than 200 data samples here. This shows the power of NN and proper hyperparameter tuning.

## 5 Summary of Results

BOLD indicates the highest performance. RF\*: Random Forest

Note that for F1 scores in MNIST, the worst F1 score is used



Algorithm	Credit Approval		Parkinsons		Rice Grains		Hand-Written Digits	
	<b>Accuracy</b>	<b>F1</b>	<b>Accuracy</b>	<b>F1</b>	<b>Accuracy</b>	<b>F1</b>	<b>Accuracy</b>	<b>F1</b>
KNN	0.8715	0.8701	<b>0.9373</b>	<b>0.9161</b>	0.8906	0.9056	0.8906	0.9579
RF*	0.8575	0.8537	0.8974	0.8418	0.9282	0.9370	0.9772	0.9473
NN	<b>0.8808</b>	<b>0.8806</b>	0.9128	0.8803	0.9310	0.9398	0.9521	0.8881