

# Random Forest

---

18308045 谷正阳

June 13, 2021

## Contents

<b>1</b>	<b>Importance function</b>	<b>2</b>
1.1	Purity function . . . . .	2
1.2	Information gain . . . . .	2
1.3	Information gain ratio . . . . .	3
1.4	Negative Gini gain . . . . .	3
<b>2</b>	<b>Discretization</b>	<b>4</b>
2.1	The first step: to calculate all split points . . . . .	4
2.2	The second step: to choose the best split point . . . . .	4
2.3	Comparison between Numpy version and PyTorch version . . . . .	5
<b>3</b>	<b>Experiment</b>	<b>5</b>
3.1	Different importance functions . . . . .	5
3.2	Different types of features . . . . .	6
3.3	Different number of features . . . . .	8
3.4	Different number of trees in a forest . . . . .	11
<b>4</b>	<b>Conclusions</b>	<b>12</b>

# 1 Importance function

This function is used to evaluate how good a split of dataset is. The basic idea is that if a label in each subset after a split dominates that subset, the split is good.

## 1.1 Purity function

To evaluate whether a label dominates a subset, we need to build a special function. Since we do binary classification here, we just input the probability of one label in the subset into the function. The return value should be high if the probability is close to 0 or 1. Otherwise, the value should be low. There are two common functions to evaluate the purity of a set, one of which is entropy, the other of which is negative Gini index. Entropy has the form of

$$B(q) = -(q \log(q) + (1 - q) \log(1 - q)), \quad (1)$$

while negative Gini index has the form of

$$Neg\_Gini(p) = -(1 - (p^2 + (1 - p^2))). \quad (2)$$

Based on these two purity functions, we build three importance functions.

## 1.2 Information gain

Information gain is to calculate the gain of entropy after the split. It has the form of

$$Gain(A) = B(\frac{p}{p+n}) - \sum_{v=1}^V \frac{p_v + n_v}{p+n} B(\frac{p_v}{p_v + n_v}). \quad (3)$$

Because after the discretization every attributes' domains are binary, and  $p, p+n$  are the same when comparing splits, the actual equation used here is

$$Neg\_Remainder\_Mul\_N(A) = (p_0 + n_0)(-B(\frac{p_0}{p_0 + n_0})) + (p_1 + n_1)(-B(\frac{p_1}{p_1 + n_1})). \quad (4)$$

The function `NEG_REMAINDER_MUL_N` below takes all the features as inputs and returns an array containing all their *Neg\_Remainder\_Mul\_N*.

```
1 def NEG_B(q):
2     I_q = 1 - q
3     return xlogy(q, q) + xlogy(I_q, I_q)
4
5 def NEG_REMAINDER_MUL_N(S, F):
6     label = S[:, -1:]
```

```

7     S1 = S[:, F]
8
9     p1 = np.sum(S1 & label)
10    N1 = np.sum(S1)
11
12    p0 = np.sum(label) - p1
13    N0 = S.shape[0] - N1
14
15    return N1 * NEG_B(p1 / N1) + N0 * NEG_B(p0 / N0)

```

### 1.3 Information gain ratio

The information gain3 has a weakness, which is that if the domain of an attribute is too big, the second term  $\sum_{v=1}^V \frac{p_v+n_v}{p+n} B(\frac{p_v}{p_v+n_v})$  can be really small. Therefore, the information gain prefer those attributes with larger domain. The information gain ratio has the form of

$$Gain\_ratio(A) = \frac{Gain(A)}{-\sum_{v=1}^V \frac{p_v+n_v}{p+n} \log(\frac{p_v+n_v}{p+n})}, \quad (5)$$

which is actually add a penalty on the information gain3. If the domain of an attribute is large, the  $-\sum_{v=1}^V \frac{p_v+n_v}{p+n} \log(\frac{p_v+n_v}{p+n})$  will be large, so its gain ratio is small.

However, since all the attributes here have binary domains, the information gain ratio is the same as the information gain.

### 1.4 Negative Gini gain

This has a similar form as the information gain3, since it's actually replace the  $\log(p)$  with its approximate substitution  $p - 1$ . This can be calculated faster than information gain since  $p - 1$  is easier to calculate than  $\log(p)$ .

```

1 def NEG_GINI_MINUS_1(p):
2     return p ** 2 + (1 - p) ** 2
3
4 def NEG_GINI_INDEX_MINUS_1_MUL_N(S, F):
5     label = S[:, -1:]
6     S1 = S[:, F]
7
8     p1 = np.sum(S1 & label)
9     N1 = np.sum(S1)
10
11    p0 = np.sum(label) - p1

```

```

12     N0 = S.shape[0] - N1
13
14     return N1 * NEG_GINI_MINUS_1(p1 / N1) + N0 * NEG_GINI_MINUS_1(p0 / N0)

```

## 2 Discretization

Since the data preprocessed in the last part of project are all continuous value. And the decision tree requires those attributes to be discretized. This part is implemented in function read. The function is used to find the best split point to split the domain of an attribute into 2 labels.

### 2.1 The first step: to calculate all split points

To discretize an attribute, firstly, we build an array containing all values appear in the attribute. And we sort it and remove all the duplicate values in it. Now we get a sorted domain of the attributes. Then we can use matrix operation to generate an array containing all split points.

```

1 mids = X[:, f]
2 np.sort(mids)
3 mids = np.unique(mids)
4 mids = (mids[:-1] + mids[1:]) / 2

```

### 2.2 The second step: to choose the best split point

We can get scores of split points using importance function defined above. and choose the split with highest importance.

```

1 x_best = None
2 mid_best = None
3 score_best = -float("inf")
4 for mid in mids:
5     x = (X[:, f:f + 1] >= mid)
6     score = IMPORTANCE(np.concatenate((x, Y), axis=1), np.array([0]))
7     if score_best < score:
8         x_best = x
9         mid_best = mid
10        score_best = score
11 X[:, f:f + 1] = x_best
12 threshold[f] = mid_best

```

## 2.3 Comparison between Numpy version and PyTorch version

Since there are some matrix operation, it seems like the calculation can be speeded up by CUDA which is supported by PyTorch but not supported by NumPy. Therefore, I implement discretization in 2 versions and compare their performance.

```
> python .\performance.py
Numpy CPU:
100%| 256/256 [00:26<00:00, 9.48it/s]
100%| 256/256 [00:25<00:00, 10.01it/s]
100%| 256/256 [00:25<00:00, 10.07it/s]
100%| 256/256 [00:24<00:00, 10.50it/s]
PyTorch CPU:
100%| 256/256 [01:17<00:00, 3.31it/s]
100%| 256/256 [01:17<00:00, 3.32it/s]
100%| 256/256 [01:22<00:00, 3.12it/s]
100%| 256/256 [01:46<00:00, 2.40it/s]
PyTorch GPU:
100%| 256/256 [03:02<00:00, 1.40it/s]
100%| 256/256 [03:30<00:00, 1.22it/s]
100%| 256/256 [04:03<00:00, 1.05it/s]
100%| 256/256 [03:22<00:00, 1.26it/s]
```

The result shows that neither PyTorch with CPU nor PyTorch with GPU is faster than Numpy. The reason may be that there are lots of loop operation which performs worse on PyTorch than on Numpy.

## 3 Experiment

I divide the output into 4 outputs like what I do in the last project. Therefore I need to build 4 random forests in each model.

### 3.1 Different importance functions

I test different importance functions. Since we have discussed that information gain ratio and information gain are the same before, we just test the performance of information gain and negative Gini index gain.

```
XY_list, threshold_list = read('PCA_dataset_1024_512.csv', NEG_REMAINDER_MUL_N)
save_threshold_list(threshold_list, 'threshold_gain_1024_512.txt')
100%| 512/512 [08:45<00:00, 1.03s/it]
100%| 512/512 [08:52<00:00, 1.04s/it]
100%| 512/512 [09:35<00:00, 1.12s/it]
100%| 512/512 [09:32<00:00, 1.12s/it]
```

Figure 1: Discretization using information gain

```
XY_list, threshold_list = read('PCA_dataset_1024_512.csv', NEG_GINI_INDEX_MINUS_1_MUL_N)
save_threshold_list(threshold_list, 'threshold_gini_1024_512.txt')
```

```
512/512 [06:05<00:00, 1.40it/s]
512/512 [06:32<00:00, 1.30it/s]
512/512 [05:58<00:00, 1.43it/s]
512/512 [05:48<00:00, 1.47it/s]
```

Figure 2: Discretization using negative Gini index gain

```
for B in (16, 32, 64):
    XY_list, _ = read('PCA_dataset_1024_512.csv', 'threshold_gain_1024_512.txt')
    counts, counts_sep, tot = test(XY_list, NEG_REMAINDER_MUL_N, B)
    print('\nB:', B)
    print('CV ID\tMultioutput Accuracy\tI Accuracy\t\tN Accuracy\t\tT Accuracy\t\tJ Accuracy')
    for i in range(5):
        print(str(i) + '\t' + str(counts[i] / tot[i]), end='\t')
        for j in range(4):
            print(counts_sep[i][j] / tot[i], end='\t')
        print()
```

```
20/20 [09:04<00:00, 27.21s/it]
```

B: 16	CV ID	Multioutput Accuracy	I Accuracy	N Accuracy	T Accuracy	J Accuracy
0	0	0.25302593659942363	0.768299711815562	0.8737752161383285	0.631700288184438	0.6132564841498559
1	1	0.24380403458213257	0.7538904899135447	0.8564841498559078	0.6172910662824207	0.6409221902017291
2	2	0.26570605187319885	0.7677233429394813	0.8806916426512968	0.6380403458213256	0.6414985590778098
3	3	0.2829971181556196	0.792507204610951	0.8536023054755043	0.6438040345821325	0.6426512968299711
4	4	0.2497116493656286	0.7687427912341407	0.8512110726643599	0.6464821222606689	0.6072664359861591

Figure 3: Attributes selection using information gain

```
for B in (16, 32, 64):
    XY_list, _ = read('PCA_dataset_1024_512.csv', 'threshold_gini_1024_512.txt')
    counts, counts_sep, tot = test(XY_list, NEG_GINI_INDEX_MINUS_1_MUL_N, B)
    print('\nB:', B)
    print('CV ID\tMultioutput Accuracy\tI Accuracy\t\tN Accuracy\t\tT Accuracy\t\tJ Accuracy')
    for i in range(5):
        print(str(i) + '\t' + str(counts[i] / tot[i]), end='\t')
        for j in range(4):
            print(counts_sep[i][j] / tot[i], end='\t')
        print()
```

```
20/20 [08:29<00:00, 25.46s/it]
```

B: 16	CV ID	Multioutput Accuracy	I Accuracy	N Accuracy	T Accuracy	J Accuracy
0	0	0.23919308357348704	0.7596541786743516	0.8668587896253602	0.6328530259365994	0.6161383285302594
1	1	0.28472622478386167	0.7832853025936599	0.861671469740634	0.6628242074927954	0.6357348703170029
2	2	0.2501440922190202	0.7844380403458213	0.8593659942363112	0.6449567723342939	0.5936599423631124
3	3	0.25302593659942363	0.7585014409221902	0.8599423631123919	0.6299711815561959	0.6265129682997118
4	4	0.26239907727797	0.7698961937716263	0.8662053056516724	0.6378316032295271	0.6199538638985006

Figure 4: Attributes selection using negative Gini index gain

Their accuracy is almost the same and the time consumption of Gini is a little less than that of gain.

### 3.2 Different types of features

I test several kinds of features as follows.

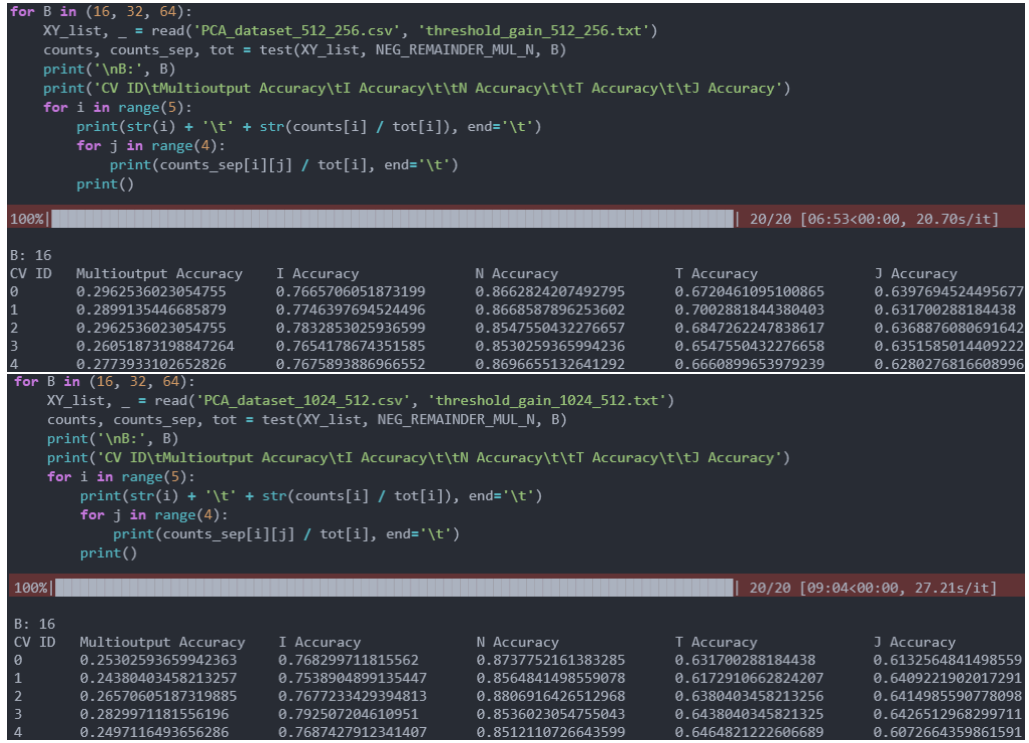


Figure 5: Tf-Idf processed by feature selection and PCA

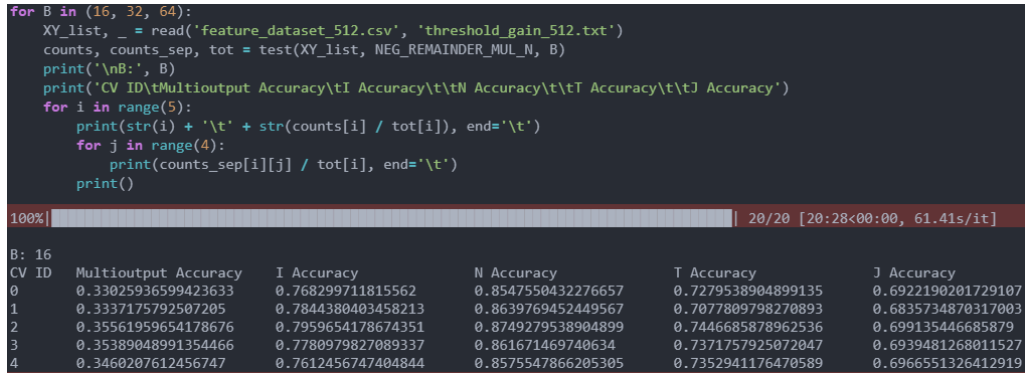


Figure 6: Tf-Idf processed by feature selection

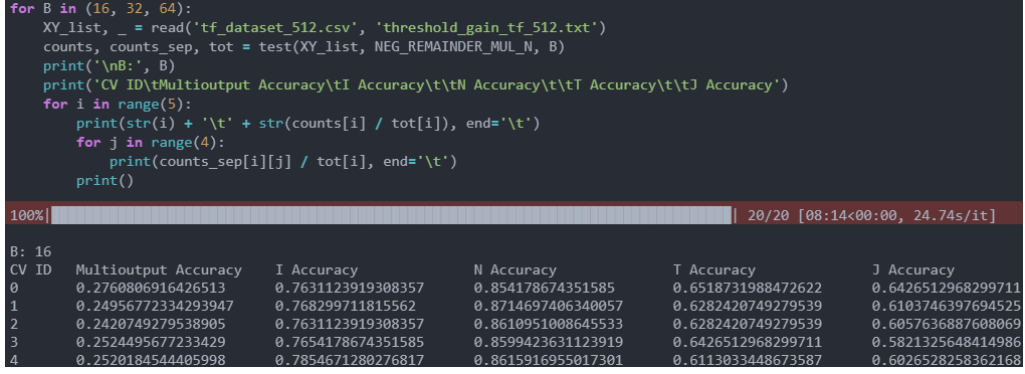


Figure 7: Tf processed by feature selection

It turns out that TF-Idf processed by feature selection has the best performance.

### 3.3 Different number of features

It seems like the larger the number of features is, the better the decision tree predicts. However, if the number of features is too large, those trivial features may greatly affect the random forest. That is because the random forest needs to randomly select attributes before choosing the best attribute. Too many trivial features may cause that the features randomly selected are all trivial.

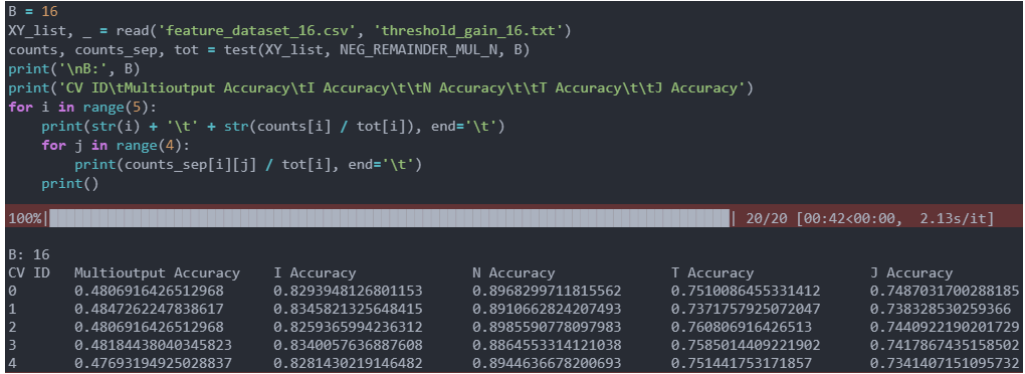


Figure 8: 16 features



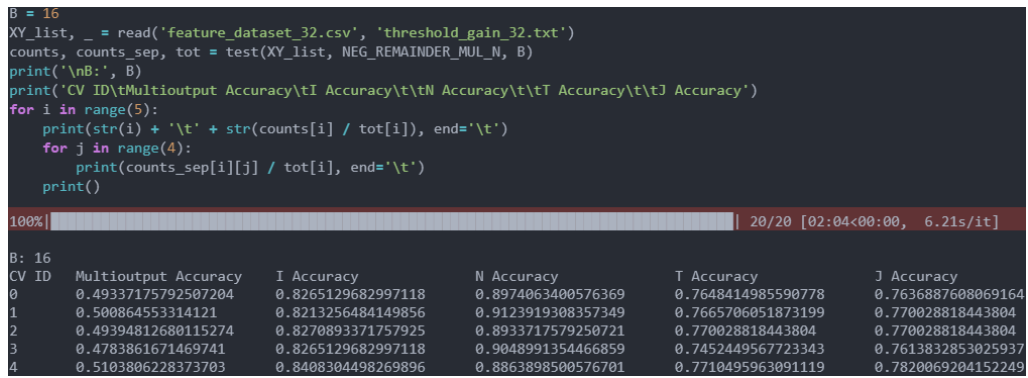


Figure 9: 32 features

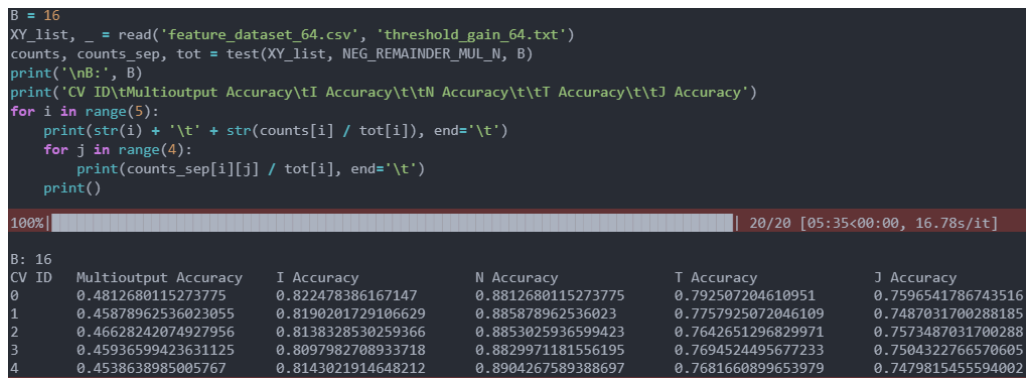


Figure 10: 64 features

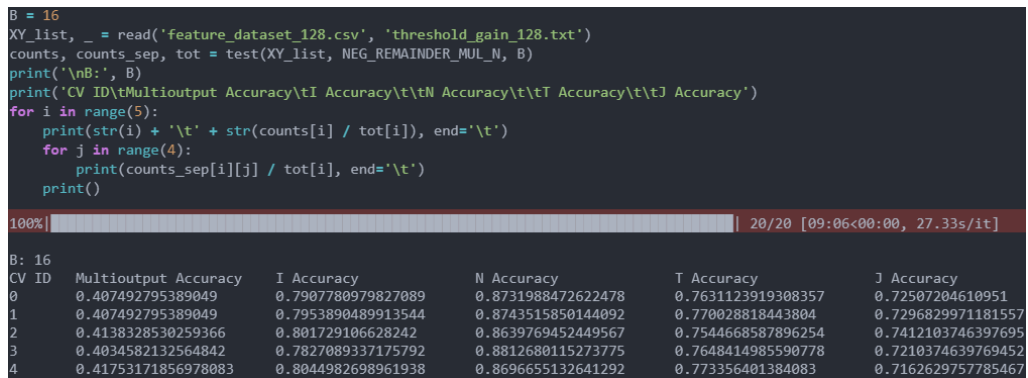


Figure 11: 128 features

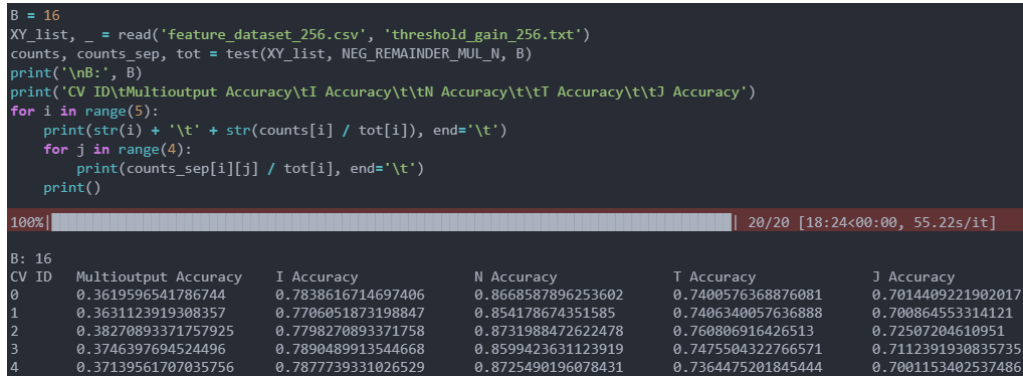


Figure 12: 256 features

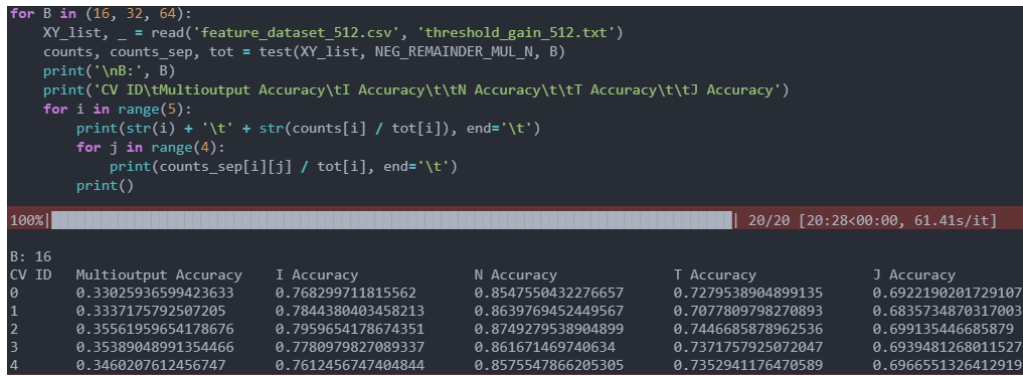


Figure 13: 512 features

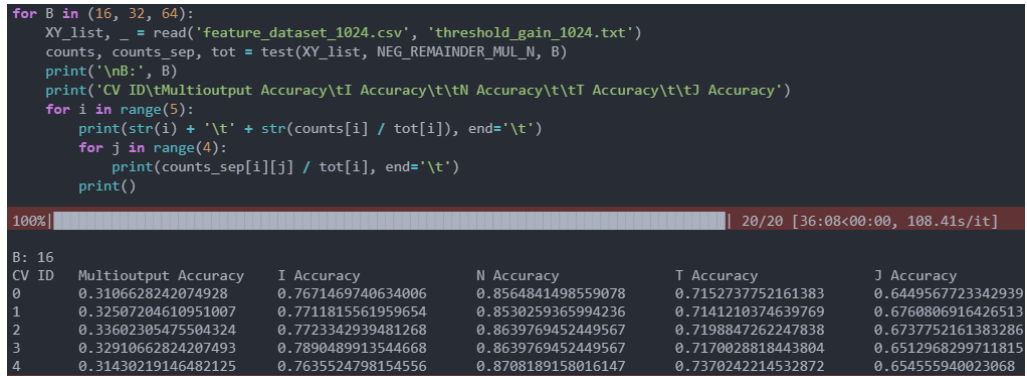


Figure 14: 1024 features

The random forest with 32 features have the best performance.

### 3.4 Different number of trees in a forest

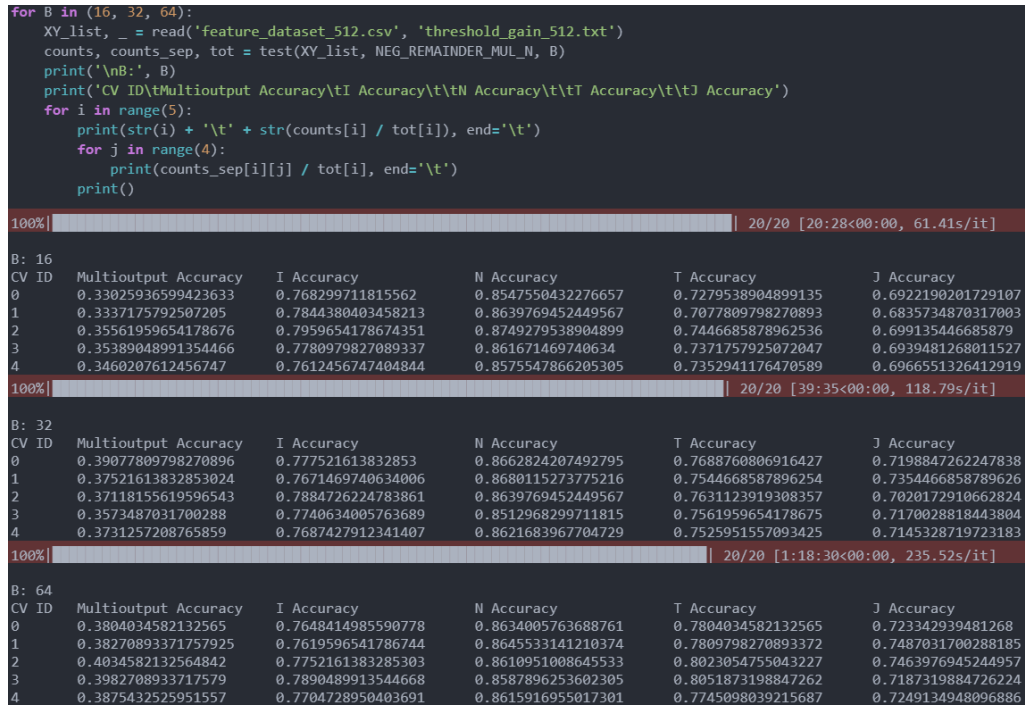


Figure 15: 512 features

