

MUSIC TRANSCRIPTION WITH MULTILABEL CLASSIFICATION

Nora Huang

Department of Computer Science
University of Victoria
Victoria, B.C, Canada
norah@uvic.ca

Aazim Lakhani

Department of Computer Science
University of Victoria
Victoria, B.C, Canada
aazimlakhani@uvic.ca

Parul Parul

Department of Computer Science
University of Victoria
Victoria, B.C, Canada
pparul@uvic.ca

ABSTRACT

Automatic music transcription is an active research topic in the field of Music Information Retrieval (MIR). The goal of music transcription is to determine the score-like representation of the input audio music. The technologies of it have been used in many mobile and web applications. However, the accuracy of the current transcription approach is still below that of a human expert. The majority of current research is focused on how to separate the pitches from the audio, then estimate them. In this project, time and frequency domain features for piano music files were extracted. The data was filtered through normalization, and feature subsets with maximum entropy and correlation were selected for use in well known classifiers such as Decision Trees, Random Tree, Nearest Neighbor, and SVM. This process enables the project model to detect musical notes which could be used in automatic music transcription.

1. INTRODUCTION

Automatic music transcription is the process of translating an audio music signal into a musical score sheet [?]. There are many applications of this technology, such as helping a student learn how to play music for well known soundtracks using a specific instrument. It could do so by building playable score sheets they could use for practice. It could also allow a musician or composer to record their improvised performance [?] and share their work for others to reproduce. There are various approaches to achieve these goals [?, ?, ?, ?, ?].

An ideal piece of music for transcription has clear, simple, and easily differentiable sounds. However in reality, song compositions are complex pieces of work with many different instruments and human voices. Addressing this complexity is a significant challenge. This project focuses on the problem of monophonic and polyphonic piano music for up to 4 notes. Monophonic includes a single melodic line with no accompaniment, while polyphonic may have multiple melodic lines mostly independent or in

imitation of one another. In other words, this project focuses exclusively on music composed with up to 4 notes played simultaneously.

This project builds a unique model to recognize keys used to played music notes to support music transcription. First, windows are applied to the input music files to segment the audio into frames around 11ms in length. In total, it retrieves 512 digital samples per frame. In total, 63 features are extracted from each frame. These acoustic features are post-processed by removing silent frames or values and labelling each non-silent frame, then analysed with either Weka [?] and Meka [?] to reduce the set into several reasonably sized feature subsets. The selected features are those having high entropy or co-relation, which are then passed through well known classifiers to prepare a model.

Weka is used to analyse monophonic music which can simplified as the single-label classification problem. Conversely, Meka is a multi-label classification extension of Weka to analyse polyphonic music, equivalent to a multilabel classification problem. Polyphonic music exposes the challenge of detecting the right keys used to play a chord when several others notes from potentially many different instruments are played simultaneously. This project reduces the scope of the problem by focusing on single-instrument audio tracks, leaving only the problem of detecting each note in a chord. In total, Multilabel classification [?] is applied using 4 classifiers: IBK [?], J48 [?], SMO [?], and RandomTree [?].

The overview of the project model is described in Figure 1. The classifiers are trained and tested based on this model.

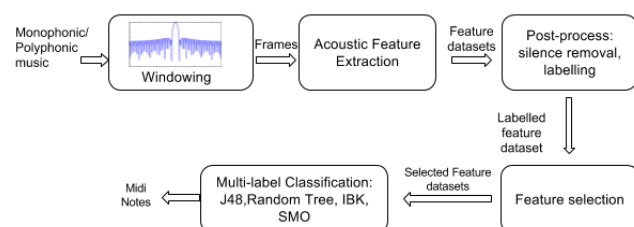


Figure 1. System Flow Diagram

The structure of the report is as follows: first, the dataset for the project is described. The tool and methods used for feature extraction are then described, along with the meaning of each set of features. Next, feature selection de-



scribes how the feature subsets were derived, analysed, and grouped. The classification section details the classifiers used and the parameters of each for every feature subset, including details on how they were trained and tested. The results of all the experiments are then evaluated and visualized for each dataset. Finally, the results are concluded and future work is outlined.

2. DATASET

The model described in Section 1, Figure 1 is evaluated using the MIDI (Musical Instrument Digital Interface) Aligned Piano Sounds (MAPS) dataset [?]. Due to copyright licensing, there are a limited number of datasets available for research. MAPS is a widely used piano database for multipitch estimation and automatic transcription of music. It provides recordings with CD quality (16-bit, 44 kHz sampled stereo audio) sounds. It is freely available under the Creative Commons License. It consists of audio files with corresponding annotations for isolated sounds, chords, and complete pieces of piano music under various different recording conditions, including various closed rooms and close ambient takes. MAPS can be divided into 4 sets:

ISOL (isolated notes and musical excerpts)
RAND (Chords with random pitch notes.)
UCHO (Usual chords from Western music)
MUS (Pieces of piano music)

For monophonic music, we use 2 second notes from the ISOL dataset. This is defined as P1 for the remainder of this report. For polyphonic music, we use polyphonic sounds with 2 to 4 notes from the RAND set. They are similarly defined as P2, P3, and P4 for the remainder of this report.

3. FEATURE EXTRACTION

In the signal processing domain, feature extraction is one of the most widely used methods to analyse a signal. There are several parts to one audio signal offering varying amounts of useful information. This section explores some of these features and selects some for use based on how much they contribute to the goal of automatic music transcription.

3.1 Feature Extraction Tool

Marsyas (Music Analysis, Retrieval and Synthesis for Audio Signals) is an open source software framework for audio processing with specific emphasis on MIR applications [?]. Bextract is one of the most powerful executables provided by Marsyas. It can be used for complete feature extraction and classification experiments with multiple files. However its classification feature is limited and did not fulfill the requirement of this project. As such, only the feature extraction functionality is used.

3.2 Feature Extraction Strategy

Using Bextract provided, the features for each frame are extracted. By default, the frame size of Bextract is 512 samples which matches the dataset sampled from MAPS. In total, it extracts 63 features from each frame. The results are formatted and stored in .arff format file for later processing. Bextract further handles the labeling for each item by using a single label - it does not support multi-labelling. As such, for the polyphonic music requiring multi-labeling, they are manually labeled with multiple values corresponding to the notes it consists of. This is done by post-processed the output .arff files with a script uniquely created for this project. This enables the data to be run directly on Meka. The script derives the labels for the output from the corresponding .txt annotations provided with the sound from MAPS.

3.3 Acoustic Features

From each frame, 63 acoustic features were extracted. They are combinations of the following 5 parameters:

Time-domain Zero-Crossings is the point where the sign of a mathematical function changes (e.g., from positive to negative), represented by a crossing of the axis (zero value) in the graph of the function [?]. It is a commonly used term in electronics, mathematics, sound, and image processing. **Spectral Centroid** is a measure used in digital signal processing to characterise a spectrum. It indicates where the "center of mass" of the spectrum is. Perceptually, it has a robust connection with the impression of "brightness" of a sound.

Rolloff is the steepness of a transmission function with frequency.

Flux is a measure of how quickly the power spectrum of a signal is changing, calculated by comparing the power spectrum for one frame against the power spectrum from the previous frame. More precisely, it is usually calculated as the 2-norm (also known as the Euclidean distance) between the two normalised spectra.

Mel-Frequency Cepstral Coefficients are coefficients that collectively make up an mel-frequency cepstrum (MFC). In sound processing, the MFC is a representation of the short-term power spectrum of a sound, based on a linear cosine transform of a log power spectrum on a non-linear mel scale of frequency.

4. FEATURE SELECTION

The feature set produced by Bextract, while effective, is a large feature set making human understanding of the data difficult. Feature selection was needed in order to reduce the usable attribute set from 63 into a more reasonable size. This section outlines the approach and result of feature set selection.

4.1 Filter-based Selection

There are multiple techniques that may be used to reduce a feature set, including using domain knowledge of the data

to manually reduce it, tailoring the set for a specific classifier, or using a systematic filter-based approach [?]. Due to the complexity of the notes and their acoustic features, and the desire to have a general feature set for use in any classifier, the systematic approach was used to decide which features were most useful.

Three approaches were used to create 5 feature sets for each PN dataset: 1) learner, 2) entropy, and 3) pearson correlation coefficient. The accuracy results of each feature set was compared against the original for 4 different classifiers in 5-fold cross validation. In order to create a standardized set selection across the datasets (i.e., they all have the same thresholds and features), the P1 dataset was selected as the baseline. That is, the feature subsets generated for P1 were assumed to be sufficient for P2 - P4.

4.2 Results

WEKAs select attributes toolset [?] was used to generate data for each approach. Two feature subsets were created from the entropy results; any attribute that provided an information gain greater or equal to 0.05 was used in the first entropy feature data set. A second set was created using a threshold of 0.05. Similarly, the pearson result was split into two feature subset using threshold values of 0.03 and 0.02. The threshold values were chosen by observation as they provided either approximately half the original feature set or slightly less which was agreed to be a reasonable attribute set size. Table7.4 summarizes the results.

Table 1. Feature subset selection attribute size results

Filter Algorithm	Threshold	Attributes Size	Reduced %
None	N/A	63	–
Learner	N/A	28	66
Entropy	0.05	14	77
	0.03	31	51
Pearson	0.03	27	57
	0.02	36	43

5. PREPROCESSING WITH NORMALIZATION

Normalization scales all numeric variables in the range [0,1] to ensure the outputs are scale invariant. In other words, it makes sure big values dont dominate smaller ones. One possible formula is given below:

$$\chi_{new} = \frac{\chi - \chi_{min}}{\chi_{max} - \chi_{min}}$$

Normalization was used on the P1 original dataset (i.e., with all 63 attributes) to determine if this technique could improve the accuracy of the results. As seen in Table5, the results of normalization did not significantly improve the accuracy of 4 different classifiers. This result meant that normalization was not applied, and was not attempted for the remaining P2-P4 datasets.

Table 2. Accuracy Results Compare (%)

	IBK	J48	SMO	Random tree
Without Normalization	97.98	88.98	54.75	87.65
With Normalization	97.98	88.98	54.76	87.67

6. MULTILABEL CLASSIFICATION

In MIR, multi-label classification are usually used to identify genres or emotion recognition, but are not used in transcription. Current research into automatic transcription tends to focus only on either the main melody, or on the approach by separating the multiple melodic lines and processing them individually. This project explores multilabel classification for polyphonic transcription by treating each note from the multiple melodic lines as a label, where each label is one of 88 classes corresponding to 88 piano keys (i.e., 11 octaves of 8 notes).

Using the feature attribute subsets derived in Section 4, four classifiers were used to experiment with results on all datasets from P1 to P4. The four classifiers were: 1) IBK, 2) J48, 3) SMO, and 4) RandomTree. Training and testing was executed in 5 fold cross-validation on each dataset. This section presents the raw data results and some visualization of the data.

7. EVALUATION

This section presents the results of running the training and testing classifiers on the P1 dataset in Weka, and the same approach for P2 - P4 in Meka. All datasets use the feature subsets derived in Section 4, and each dataset is separated into its own subsection.

7.1 P1 Monophonic Dataset Evaluation

This section presents the results of the first dataset, provides visualizations of the J48 tree, and analyses the J48 confusion matrix. J48 was chosen for further analysis as it was quick to reproduce on any machine while not sacrificing accuracy.

Feature subset results

Table 7.1 summarizes the result of each classifier on each feature subset. Anecdotally, SMO took the longest amount of time to run (approximately 5-10 minutes) and was assumed to have the worst performance, validated by its low correct percentage classified results. IBK took a moderate amount of time to run while providing the highest accuracy. RandomTree was slightly faster than J48 but produced consistent lower accuracies.

J48 Analysis

It was found that the initialization of all J48 trees for each feature subset was identical. That is, each set of attributes derived the same root node and first layer of their J48 classifier. However, to emphasize the differences between the

Table 3. Accuracy results of different feature subsets on different classifiers (%)

Filter (Threshold)	IBK	J48	SMO	Random tree
Original	97.98	88.98	54.75	87.65
Learner	97.37	88.54	46.11	87.63
Entropy (0.05)	91.48	87.08	27.88	86.87
Entropy (0.03)	98.61	88.59	51.86	87.83
Pearson (0.03)	99.41	87.27	48.96	87.26
Pearson (0.02)	99.14	88.10	48.58	87.86

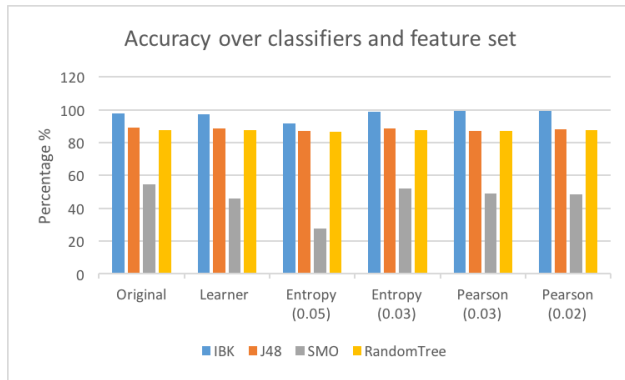


Figure 2. Accuracy for each dataset and classifiers for P1 datasets

feature set trees, Figures 5 through 10 depict the general structures of each.

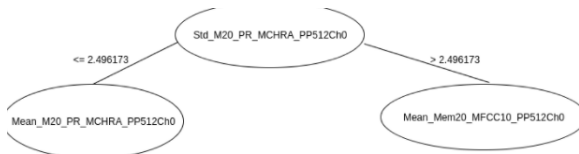


Figure 3. Root Node and First layer for all P1 instances

The confusion matrix of each feature set was reviewed. The intersection of notes with the highest misclassified percentage value is summarized in Table 7.1.

It is not surprising that the feature sets have similar misclassified note intersections. It is reasonable to expect that similar sounding notes would be consistently misclassified no matter which subset and classifier were used.

7.2 P2 Dataset

This dataset has 2 musical notes being played simultaneously. Table 6 summarizes the percentage output of the accuracy for each feature set. The accuracy is considered as note1 being correctly classified, note2 being correctly classified, all notes are exactly being classified. The train and

Table 4. P1 Feature subset highest confusion results

Filter (Threshold)	Intersection	Notes
Original	(m28, m27)	55 / 3420
Learner	(m28, m27)	66 / 3420
Entropy (0.05)	(m23, m24)	73 / 3074
Entropy (0.03)	(m28, m27)	75 / 3420
Pearson (0.03)	(m22, m21)	72 / 3370
Pearson (0.02)	(m28, m27)	61 / 3420

test time are also summarized in this section. Even though we spread the test over different computers and the running time varies on different machines, we can still see the trend of it. For P2-P4 dataset SMO is almost impossible without high computation power. Attempting this classifier on the original dataset it took overnight without finishing. It was not a feasible classifier to use for audio transcription. As such it was skipped for P2-P4 due to poor performance.

In table 6, it is clear that IBK has highest accuracy but takes the longest amount of time compared to the other two classifiers. The entropy dataset with threshold of 0.05 consistently has the worst performance, while the remainder all have similar accuracy results.

Figures 4 and 5 are visual representations of tables 5 and 6. It is clear to see the similar performances of all the results.

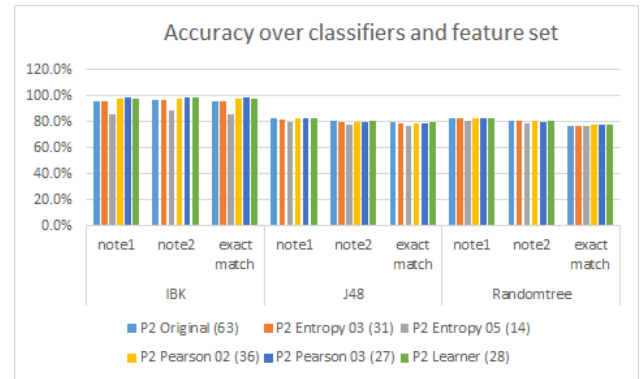


Figure 4. Accuracy for each dataset and classifiers for P2 datasets

7.3 P3 Dataset

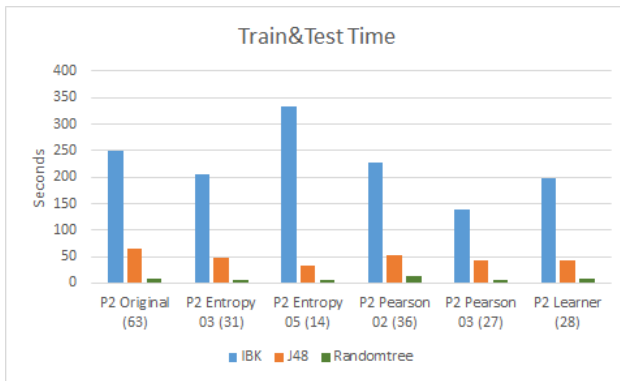
This dataset has 3 musical notes being played simultaneously. Table 7 summarizes the accuracy for the dataset against the 3 classifiers, and table 8 summarizes the run times.

Table 5. Percent Accuracy results of different feature subsets on different classifiers for P2

P2 Filter (threshold)	IBK			J48			Random Tree		
	note1	note2	exact match	note1	note2	exact match	note1	note2	exact match
Original	96.3	96.9	96.1	83.1	80.9	79.7	82.7	80.9	76.9
Learner	98.4	98.6	98.3	82.9	80.7	79.5	82.7	81.0	78.3
Entropy (0.05)	86.2	88.8	85.5	80.2	78.0	76.4	80.9	78.4	76.9
Entropy (0.03)	96.3	96.9	96.1	82.3	80.1	78.8	82.7	80.9	76.9
Pearson (0.03)	99.0	99.1	99.0	82.5	80.1	78.8	82.6	79.7	78.1
Pearson (0.02)	97.6	98.0	97.5	82.5	80.2	79.0	82.9	80.6	78.2

Table 6. Training and Testing times for P2

Random Tree Time(S)	J48 Time(S)	IBK Time(S)	P2 Filter(Threshold)
8.3	64.96	248.6	Original
8.4	43.5	197.0	Learner
6.0	33.7	334.4	Entropy(0.05)
5.5	46.8	205.2	Entropy(0.03)
5.2	43.4	139.3	Pearson 03(0.03)
12.3	51.9	226.5	Pearson 02(0.02)

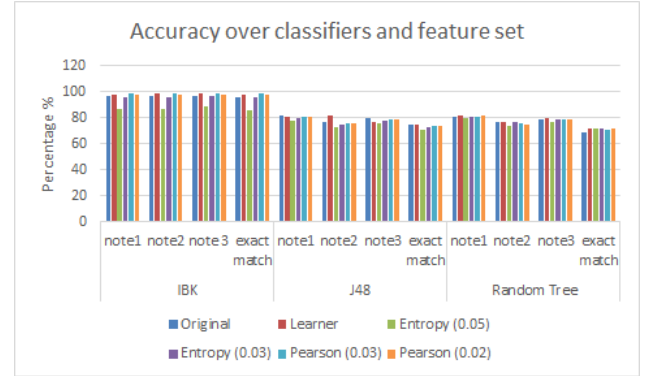
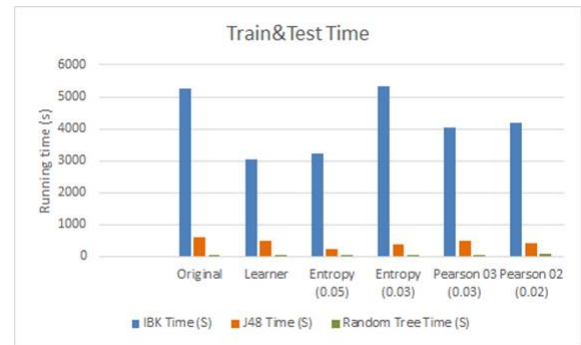
**Figure 5.** Train and Test Time for each dataset and classifiers(P2)

7.4 P4 Dataset

This dataset has 4 musical notes being played simultaneously. Table 9 summarizes the accuracy results. Table 10 summarizes the test time.

7.5 Exact Match Accuracy Comparison for PN Dataset

In this section we compare the exact match accuracy over different monophonic and polyphonic notes data set with various feature sets. It is not surprising that, the higher number of notes a dataset has, the lower accuracy the classifier provides. Interestingly, this misclassification margin is somewhat consistent across the datasets. Overall, IBK offers the highest accuracy while RandomTree is the lowest. Table 11 summarizes these findings.

**Figure 6.** Accuracy for each dataset and classifiers for P3 datasets**Figure 7.** Train and Test Time for each dataset and classifiers(P3)

8. CONCLUSION

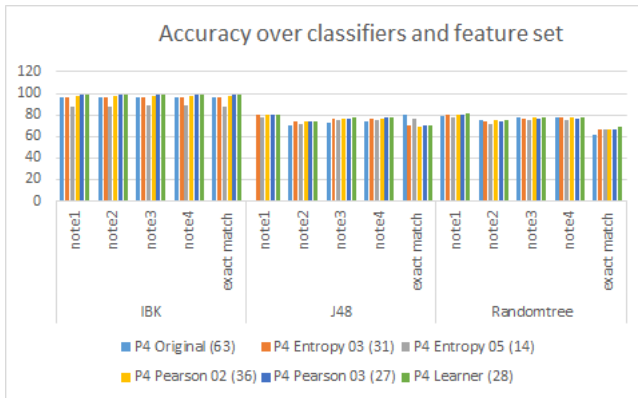
This exploratory project used the MAPS dataset for various audio files, ranging from single notes to 4 notes played simultaneously. To extract the features from the data source, Bextract was used, deriving 63 feature attributes for analysis. WEKA was used to develop 5 reduce feature sets using filter-based selection approach. 4 classifiers, IBK, SMO, J48, and Random tree for P1 were found. Their results were presented and compared for accuracy depending on which feature subset and which classifier were used. Using MEKA, IBK, J48, and RandomTree classifiers were used on datasets P2 - P4. Similarly, these accuracy results were presented and compared. In summary, this project determined IBK to be the best classifier for determining notes from an input audio file. In conclusion, this project successfully provided the framework for future work on

Table 7. Percent Accuracy results of different feature subsets on different classifiers for P3

P4 Filter (threshold)	IBK				J48				Random Tree			
	note1	note2	note 3	exact match	note1	note2	note3	exact match	note1	note2	note3	exact match
Original	96.4	96.5	97.0	96.3	82.0	76.7	79.7	75.1	80.4	76.8	78.6	68.6
Learner	98.4	98.5	98.6	98.4	80.7	81.7	76.5	74.8	81.7	76.7	79.8	72.2
Entropy (0.05)	86.8	87.1	88.7	86.2	78.1	72.6	75.8	70.7	79.5	73.7	77.1	71.7
Entropy (0.03)	95.9	96.0	96.5	95.7	80.2	74.7	77.9	72.9	80.9	76.5	79.2	71.9
Pearson (0.03)	99.0	99.0	99.1	98.9	80.9	75.5	78.6	73.8	80.4	75.6	78.5	71.3
Pearson (0.02)	97.7	97.8	98.0	97.6	81.0	75.6	78.7	73.9	81.4	75.3	78.9	72.3

Table 8. Training and Testing times for P3

Random Tree Time(S)	J48 Time(S)	IBK Time(S)	P4 Filter(Threshold)
60.8	597.8	5245.9	Original
64.5	510.2	3051.5	Learner
49.5	242.3	3228.5	Entropy(0.05)
70.6	392.7	5334.4	Entropy(0.03)
64.5	480.3	4055.7	Pearson 03(0.03)
81.1	405.3	4194.2	Pearson 02(0.02)

**Figure 8.** Accuracy for each dataset and classifiers for P4 datasets

transcribing a piano audio file into a human-readable score sheet.

9. FUTURE WORK

The results of this project support the ultimate goal of MIR to audio transcription to scoresheets by use of classifiers on sampled audio files. However, these results are limited to the ideal world where the number of notes and complexity of the file is already established. The most useful results will only be produced when the problem of unknown number of notes is solved. Future work should focus on any number of polyphonic music transcription using multilabel classification. Adding a blank class to the current 88 key classes should be one thing that we can try. We can suppose there are 10 notes being played simultaneously (by humans 10 fingers). If actually there are only 4 notes, then 6 out of ten will be classified as blank class while the 4 notes can be classified as those in 88 key classes.

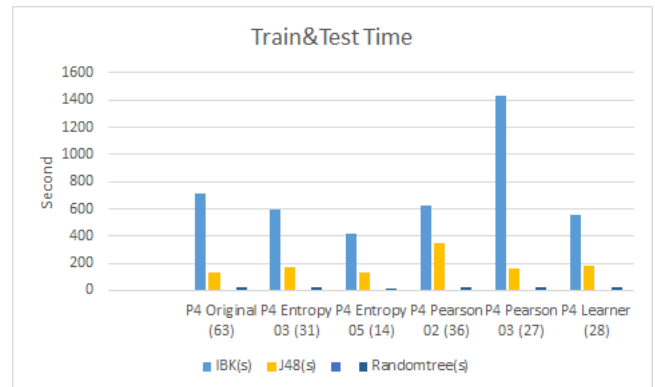
**Figure 9.** Train and Test Time for each dataset and classifiers(P4)

Table 9. Percent Accuracy results of different feature subsets on different classifiers for P4

P4 Filter (threshold)	IBK					J48					Random Tree				
	note1	note2	note3	note4	exact match	note1	note2	note3	note4	exact match	note1	note2	note3	note4	exact match
Original	96.6	96.7	96.7	96.9	96.5	76.9	70.7	73.3	74.0	79.7	79.4	75.5	77.1	77.8	61.9
Learner	98.5	98.5	98.5	98.6	98.4	80.7	74.5	77.3	77.9	70.5	81.3	75.1	77.6	77.9	69.0
Entropy (0.05)	88.2	88.2	88.4	89.4	87.7	78.3	72.0	74.7	75.4	76.4	78.3	72.0	74.7	75.4	66.8
Entropy (0.03)	96.1	96.1	96.2	96.5	96.0	80.0	73.9	76.6	77.0	69.7	79.8	73.9	76.3	77.3	66.9
Pearson (0.03)	99.0	99.0	99.0	99.0	98.9	80.1	73.7	76.5	77.1	69.6	80.4	73.7	76.6	76.4	66.9
Pearson (0.02)	97.8	97.8	97.9	98.0	97.7	79.7	73.5	76.2	76.8	69.4	80.6	74.9	77.2	78.3	66.9

Table 10. Training and Testing times for P4

Random Tree Time(S)	J48 Time(S)	IBK Time(S)	P4 Filter(Threshold)
25.1	135.6	712.3	Original
19.7	183.0	559.5	Learner
17.4	133.6	416.8	Entropy(0.05)
20.7	171.2	598.7	Entropy(0.03)
20.8	161.3	1427.2	Pearson 03(0.03)
25.4	345.5	627.5	Pearson 02(0.02)

Table 11. Summary of exact match percentages across all datasets

P4 Filter (threshold)	IBK				J48				Random Tree			
	P1	P2	P3	P4	P1	P2	P3	P4	P1	P2	P3	P4
Original	98.0	96.1	96.3	96.7	89.0	79.7	75.1	79.7	87.7	76.9	68.6	61.9
Learner	97.4	98.3	98.4	98.4	88.5	79.5	74.8	70.5	87.6	78.3	72.2	69
Entropy (0.05)	91.5	85.5	86.2	87.7	87.1	76.4	70.7	76.4	86.9	76.9	71.7	66.8
Entropy (0.03)	98.6	96.1	95.7	96	88.6	78.8	72.9	69.7	87.8	76.9	71.9	66.9
Pearson (0.03)	99.4	99	98.9	98.9	87.3	78.8	73.8	69.6	87.6	78.1	71.3	66.9
Pearson (0.02)	99.1	97.5	97.6	97.7	88.1	79	73.9	69.4	87.9	78.2	72.3	66.9