

It's What You Say and How You Say It: Exploring Textual and Audio Features for Podcast Data

Anonymous ACL submission

Abstract

Podcasts are relatively new media in the form of spoken documents or conversations with a wide range of topics, genres, and styles. With a massive increase in the number of podcasts and their listener base, it is beneficial to understand podcasts better, to derive insights into questions such as what makes certain podcasts more popular than others or which tags help in characterizing a podcast. In this work, we provide a comprehensive analysis of hand-crafted features from two modalities, i.e., text and audio. We explore multiple feature combinations considering podcast popularity prediction and multi-label tag assignment as proxy downstream tasks. In our experiments, we use document embeddings, affective features, named entities, tags, and topics as the *textual features*, while multi-band modulation and traditional speech processing features constitute the *audio features*. We find the audio feature *prosody* and textual affective features, *sentiment* and *emotions* are significant for both the downstream tasks. We observe that the combination of textual and audio features helps in improving performance in the popularity prediction task.

1 Introduction

Podcasts have emerged as an exciting medium for entertainment, advertising, news, and information dissemination. According to Nielsen (Nielsen, September 2021), the total number of podcast titles is reaching 2 million with a steady increase in listeners across all demographics. The Interactive Advertising Bureau (IAB) believes that the US-Podcast revenue will see a big jump from \$842 million in 2020 to \$2 billion by 2023. Content creator apps like Anchor (Anchor, 2022), and Riverside (Riverside.fm, 2022) provide an easy framework to record, edit, and publish a podcast on the media platforms. With a steady increase in the listener base and podcast content, several open problems in handling and accessing this information have

emerged. Jones et al. (2021) highlights some of the unique challenges and future directions in the domain of podcast information access. They point out that the existing technologies addressing tasks such as *Search*, *Recommendation*, *Summarization*, and *User experience* are inadequate to handle the multi-genre, multi-style and multi-format composition of podcasts. As this opens up an exciting landscape for future research, we believe that a comprehensive feature analysis of podcasts can serve as an important groundwork in tackling some of these problems. These features need to address *what* podcasts contain and *how* they are delivered.

Podcasts represent inherently heterogeneous data consisting of music and speech in different spoken and written styles and formats. With this work, we investigate the efficacy of individual traditional features and their combinations in understanding *whats* and *hows* in the context of podcast data. Textual and audio features such as tags, sentiment, emotions, and modulation-based features are traditionally considered as a basis of downstream applications. We formulate an application-oriented framework to evaluate and understand the interplay among these features. We consider two separate applications, i.e., *Podcast Popularity Prediction*, and *Podcast Tag Assignment*, each targeting a specific set of features and modality. For example, we hypothesize that textual features may be more informative in the tag assignment task, whereas audio features may be more relevant in the popularity prediction task, as audio captures the style and listening experience important for popularity and tags are more dependent on the content. We also hypothesize that combining features from different modalities can be more informative in characterizing a podcast.

To summarize, in this work, we evaluate various hand-crafted features, from both podcast audio and its textual transcript, that are necessary to understand podcasts. We briefly describe each of these

042
043
044
045
046
047
048
049
050
051
052
053
054
055
056
057
058
059
060
061
062
063
064
065
066
067
068
069
070
071
072
073
074
075
076
077
078
079
080
081
082

083 features and how they contribute to representing
084 podcasts (Section 2). We consider the aforemen-
085 tioned proxy downstream tasks to study and eval-
086 uate these features. We compile and curate the
087 dataset introduced in (Yang et al., 2019) (Section
088 3). We experiment with multiple combinations of
089 text and audio features to evaluate their effect and
090 usefulness in the context of the tasks mentioned
091 above (Section 4). We further provide an in-depth
092 analysis on these features and their efficacy in the
093 context of existing experiment setup (Section 5).
094 Our contribution in this work is a detailed feature
095 analysis useful for future research addressing pod-
096 cast data.

097 2 Features

098 As a podcast is essentially an audio experience with
099 focused textual content, we consider both text and
100 audio modalities in this study. Specifically, we
101 focus on widely used hand-crafted features for our
102 experiments. Next, we present a brief overview of
103 features and their characteristics in the context of a
104 podcast.

105 2.1 Audio Features

106 Audio features we consider need to represent *how*
107 a podcast is delivered. For audio modeling, we
108 use short-time modulation features derived through
109 multi-band filtering and energy operators and the
110 other traditional speech processing features such
111 as Mel-frequency Cepstrum Coefficients (MFCC),
112 prosody features, and emotions to extract salient
113 properties.

114 **Modulation:** Short-time features like energy and
115 absolute amplitude may assign less importance
116 score to *content-rich* podcasts spoken with low
117 amplitude. For example, the spectrogram and
118 magnitude-based features typically used for pod-
119 cast analysis may ignore the low energy and high-
120 frequency sound (Dimitriadis et al., 2005). We
121 believe such sound may contribute to predicting
122 the popularity of a podcast. Though the magnitude
123 information correlates well with the understanding
124 of speech, the spectral content can also convey com-
125 parable speech intelligibility and also provides in-
126 formation about the speaker characteristics that we
127 believe can distinguish popular podcasts (Boashash,
128 1992). Here, we introduce the features that equally
129 consider the energy, amplitude, and spectral con-
130 tent. For any multi-component audio signal like

131 the podcast, we need to break it down into its Am-
132 plitude, and Frequency modulation (AM-FM) com-
133 ponents (Boashash, 1992), since there could be
134 multiple frequencies varying as a function of time.
135 (Zlatintsi et al., 2012) use these AM-FM compo-
136 nents to assign a measure of interest (importance)
137 to audio frames for audio event detection and sum-
138 marization task. For this work, we believe that
139 these AM-FM components can prove beneficial in
140 predicting the popularity of a podcast since they
141 can capture the dynamic nature and preserve the
142 subtle harmonic structures present in audio (Dimit-
143 triadis et al., 2005).

144 Here, we compute the AM-FM features by multi-
145 band filtering the audio signal using 40 Gabor fil-
146 ters (Evangelopoulos and Maragos, 2006). A non-
147 linear energy tracking operator, the Teager Energy
148 Operator (TEO), estimates the squared product of
149 the instant amplitude and frequency for every multi-
150 band filtered signal. To extract these individual
151 features, we rely on the Energy Separation Algo-
152 rithm (ESA) (Evangelopoulos and Maragos, 2006;
153 Maragos et al., 1993). The ESA tracks a filter
154 that records the Maximum average Teager Energy
155 (MTE) and also computes their corresponding val-
156 ues of Mean Instantaneous Amplitude (MIA) and
157 Mean Instantaneous Frequency (MIF). We extract
158 1198-dimensional MTE, MIA, and MIF representa-
159 tions using 12s snippets of the leading 10 minutes
160 of every podcast (Yang et al., 2019). To reduce
161 the computational requirements, we resample the
162 podcast at 16kHz and use a window size of 25ms
163 and a window shift of 10ms.

164 **MFCC:** MFCC is based on the human auditory
165 system and employs a nonlinear scale to corre-
166 late with the human perception of the frequency
167 contents of a sound. We consider 39-dimensional
168 delta and double-delta MFCC representation us-
169 ing the similar parameters adopted for extracting
170 modulation-based features.

171 **Prosody:** Prosody (non-verbal) features effec-
172 tively capture the speaker characteristics, their
173 speaking style, emotional state, and approximately
174 identify the listener’s interest towards a section
175 of an audio (Adell et al., 2005). We use PRAAT
176 to extract 15-dimensional prosody features on the
177 entire podcast (Boersma and Van Heuven, 2001).
178 The prosody features include $F0$ -median, $F0$ -
179 mean, $F0$ -standard deviation, $F0$ -minimum, $F0$ -
180 maximum, number of pulses, number of periods,

their mean and standard deviation, number of unvoiced frames, number and degree of voice breaks, mean autocorrelation, mean noise-to-harmonic and mean harmonic-to-noise ratio.

Affective Audio Features: The emotional state in speech and audio is one of the most important paralinguistic messages captured during human interactions. We use a Wav2Vec2 fine-tuned model on the IEMOCAP database to extract 4-dimensional *audio emotions* features (neutral, happy, sad, and angry) of an entire podcast (Baeovski et al., 2020; Yang et al., 2021).

2.2 Textual Features

According to a Nielsen report (Insights, 2020), podcast engagement has steadily seen growth in the number of heavy as well as light listeners owing to quality content being offered to the listeners. As an exercise to understand the characteristics of such content, we evaluate various textual features. We consider a set of textual features that can help understand *what* a podcast talks about.

Tags: It is observed that every podcast is assigned a set of tags representing its categorization such as *Arts, Society, Sports, Business*, etc. This assignment is often user-defined and captures the scope of the podcast. These tags are a mix of fine-grained and coarse-grained categories. We use an average 100-dimensional Glove embedding vector (Pennington et al., 2014) representing the assigned set of tags to the podcast.

Topics: Even though there exist a set of tags describing a podcast, topics give a detailed list of concepts covered. The topics along with the pre-assigned tags give a much wider representation of the podcast content. We use the unsupervised topic detection algorithm Top2Vec (Angelov, 2020) to generate a list of topics for a given podcast. The topic words are represented using an average 100-dimensional Glove embedding vector for a given podcast.

Affective Textual Features: Emotions play an important role in human cognition, including perception, attention, learning, and reasoning (Tyng et al., 2017). Affective content is more engaging to users than neutral content (Xu et al., 2014). Accordingly, we consider two perspectives of affective features given below.

1. **Sentiment:** Sentiment of a podcast refers to the inclination of the podcast content towards positive or negative polarity. We use the NLTK sentiment analyzer to extract the sentiment of the content. We create a three-dimensional sentiment representation vector for each podcast with scores corresponding to negative, neutral, and positive sentiment that sum up to 1.
2. **Text Emotions:** We consider six basic human emotions, i.e., anger, sad, happy, fear, disgust, and neutral, assigned to sentences in podcast transcripts. We use a zero-shot sentence classification setup with task-aware sentence representations (Halder et al., 2020)¹ to predict the emotion and the corresponding confidence for a sentence, which we further utilize to get an emotion representation for a given podcast.

For each sentence of the podcast, the model predicts one of the above six emotions denoted by $e_1, e_2, e_3, e_4, e_5, e_6$ and gives their respective confidence scores denoted by $s_1, s_2, s_3, s_4, s_5, s_6$. We further denote the probability of an emotion e_i by $P(e_i) = c_i/n$, where c_i is the corresponding count of emotion e_i and n denotes the count of all sentences with detected emotions in the podcast. The confidence score of a model m for a given emotion e_i is denoted by s_i where $P(m|e_i) = s_i$. For the final representation of our emotion probability vector, we calculate posterior probability of each emotion e_i for a given model m as

$$P(e_i|m) = \frac{P(e_i).P(m|e_i)}{\sum_{i=1}^6 P(e_i).P(m|e_i)}$$

Named Entities: Named entities like ‘*person*’ and ‘*organization*’ can play a significant role in attracting audience attention. These are generally not covered in tags or topics features explained above. We consider named entities in the podcast title as one of the features. We predict named entities and the corresponding confidence score using the TARS (Halder et al., 2020) zero-shot sentence classification setup. We predict named entities with four dimensions: ‘*Person*’, ‘*Organization*’, ‘*Location*’, and ‘*Others*’ represented by n_1, n_2, n_3 , and n_4 respectively. To obtain a four-dimensional named

¹https://github.com/flairNLP/flair/blob/master/resources/docs/TUTORIAL_10_TRAINING_ZERO_SHOT_MODEL.md

entities representation for a podcast, we use a similar formulation used in the text emotion feature representation.

Document Embeddings: Pre-trained language models like BERT (Devlin et al., 2019) and DistilBERT (Sanh et al., 2019) are often used to extract content representations. Specifically, we use contextual embeddings to get vector representations for the text content. Considering the longer length of textual transcripts, we use Longformer embeddings (Beltagy et al., 2020) to extract podcast content representations.

3 Dataset

Yang et al. (2019) introduced a dataset of 6,511 English language podcasts for the podcast popularity prediction task from various categories like *Arts*, *Society*, *Sports*, *Business*, etc. They scraped the publicly available iTunes podcast directory to get 46,358 episodes from 18,433 channels active from July 2016 to July 2017. They considered the ranking of channels on the iTunes chart as the basis of channel (and episodes) popularity. Top-200 podcast channels were considered popular, and all the episodes from the popular channels were labeled as *popular*. The remaining episodes were labeled as *unpopular*. At most, one episode published in the most recent two weeks from a channel was considered in the dataset². Joshi et al. (2020) use textual transcripts from this dataset to predict podcast popularity.

For our experiments, we enhance this data further first by scraping the audio files of the podcasts using the links provided along with the dataset. We downloaded 3,526 audio files out of the total 6,511 podcast episodes, as the remaining podcast links were broken or information was missing. In our experiments, we use the text transcripts as provided in the dataset. Since the transcript had no punctuation, we predict the punctuation in the transcripts using a bidirectional RNN and attention-based punctuation restoration technique (Tilk and Alumäe, 2016). Table 1 shows the distribution of the *popular* and *unpopular* episodes in the dataset. The percentage of *popular* and *unpopular* episodes are roughly the same in both modalities, though we have fewer audio files as compared to the text transcripts.

For the analysis of text and audio features on multi-label tag assignment, we extract tags from

	Text	Audio
<i>Popular</i>	837 (12.86%)	454 (12.87%)
<i>Unpopular</i>	5674 (87.14%)	3072 (87.13%)
Total	6511	3526

Table 1: Distribution of the *popular* and *unpopular* podcasts in the dataset.

the RSS feed of the podcasts. Out of 3,526 podcasts with both the text and audio data, the RSS feed (and tags) is available for 3,306 podcasts. Table 2 shows the distribution of tags in 3,306 podcasts. Originally, a podcast can be assigned one or more tags from 105 fine-grained tags. Since we have a small and highly imbalanced dataset, we manually merge these 105 fine-grained tags into 19 coarse-grained tags as described by Apple podcasts³. This allows us to map fine-grained tags such as {*‘Music’*, *‘Music Commentary’*, *‘Music History’*, and *‘Music Interviews’*} under one coarse-grained tag {*‘Music’*}. Table 3 shows the dataset distribution after the dataset was split into 80:20 for train and test set.

Tag	#Podcasts	#Popular Podcasts	%Sentiment Difference
True Crime	18	13	3.14
Fiction	28	6	4.74
Government	41	10	5.90
History	71	9	4.12
Kids & Family	100	23	9.37
Science	126	25	6.21
Music	147	28	9.97
Technology	206	31	9.17
TV & Film	283	28	9.80
Religion & Spirituality	287	38	8.45
Comedy	313	50	8.39
Arts	359	38	9.19
Education	359	65	8.48
Sports	361	20	8.72
News	370	55	6.24
Health & Fitness	371	59	8.63
Leisure	376	38	9.05
Business	486	60	9.39
Society & Culture	514	51	8.12

Table 2: Distribution of tags and their popularity in the dataset. Here, we consider 3,306 podcasts with available data from both text and audio modalities. In %Sentiment Difference column, we show the difference in average %positive and %negative sentiment for each tag.

²<https://github.com/yelongqi/podcast-data-modeling>

³<https://podcasts.apple.com/us/genre/podcasts/id26>

Split	#Podcasts	#Tags	#Tags per Podcast
Train	2671	19	1.46
Test	635	19	1.43

Table 3: Dataset distribution for the task of multi-label tag prediction.

4 Experimental Setup

The podcast data is observed to be multi-modal and heterogeneous, motivating us to experiment with combinations of diverse feature sets. Our experiments consist of studies with individual features as well as their combinations. We conduct our experiments using a 20 core CPU with 64GB RAM. To analyze and evaluate the efficacy of various features mentioned in Section 2 to characterize and understand podcast data, we use two formulations, i.e. *podcast popularity prediction* and *podcast tag assignment*.

4.1 Podcast Popularity Prediction

In this task, we seek to investigate the factors influencing the popularity of a podcast. We posit that multi-modal features capture important information to predict podcast popularity. We experiment with hand-crafted multi-modal features and analyze their efficiency using a fine-tuned XGBoost and bagging classifier-based architecture. We perform the grid search on hyperparameter space resulting in the best found combination for XGBoost as: $\gamma = 0.2$, $\text{maximum depth} = 14$, $\text{estimators} = 120$, $\text{reg_alpha} = 0.8$, and $\text{reg_lambda} = 1.2$. Figure 1 shows our binary podcast popularity classification model. We concatenate the features and use a 5-fold cross-validation as an initial step. We reduce the feature dimension using PCA with 256 components in each fold. Since the data is highly imbalanced (as shown in Table 1), we upsample the minority class (i.e., *popular*) using the SMOTE algorithm (Chawla et al., 2002). Next, we train the XGBoost classifier on the upsampled training data. We then use the hyperparameter tuned bagging classifier with 80 estimators to further address the challenge of data imbalance and synthetic data, as it trains the base XGBoost classifier on the random subset of the original dataset and aggregates the predictions. We use the test data (with PCA and without upsampling) to predict the popularity of the podcasts. We report the results using the macro-F1 score.

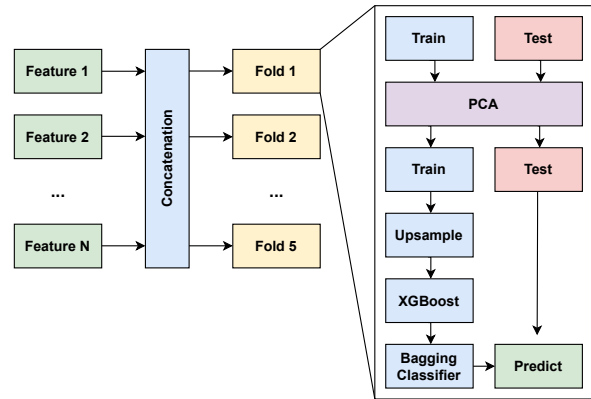


Figure 1: XGBoost and bagging based architecture for podcast popularity prediction. For every fold, the classification setup involves PCA, upsampling, XGBoost, and a bagging classifier.

4.2 Podcast Tag Assignment

In the novel multi-label podcast tag assignment task, we seek to automatically assign appropriate tags to a given podcast. We consider tags available from the RSS feed of the podcast as the ground truth, since these tags are assigned by the hosts of the podcast. As mentioned in Section 3, due to the small and imbalanced nature of the dataset, we manually merge the available 105 fine-grained tags into 19 coarse-grained tags. A podcast can have one or more coarse-grained tags associated with it, similar to a multi-class, multi-label setup. We want to observe the effect of various text and audio features in only text, only audio, and text+audio modalities. One approach to solve multi-label classification is using binary relevance, where we transform the problem into multiple binary models with a one-vs-rest setup. However, since we aim to do feature analysis in different modalities, we create a unified model for multi-label tag assignment by using a simple two-layer perceptron network. For this, we concatenate different combinations of text and audio features and feed them through a fully connected layer, followed by a dropout layer with a dropout probability of 0.3 and another fully connected layer. We use binary cross-entropy loss function with class_weights to handle dataset imbalance. We calculate class_weights for each tag l as N/N_l , where N is the total number of podcasts in the training dataset, and N_l denotes the numbers of podcasts with tag l . Since the dataset is highly imbalanced, we use the weighted-F1 score as the primary evaluation metric. Similar to (Kar et al., 2018), we also evaluate the models on tag recall

Feature(s)	Macro-F1
All Text	0.545
(-) Tags	0.500
(-) Text Emotions	0.537
(-) Sentiment	0.539
(-) Topics	0.547
(-) Named Entities	0.549
All Audio	0.802
(-) MIF	0.603
(-) Prosody	0.784
(-) MFCC	0.801
(-) MTE	0.801
(-) MIA	0.801
(-) Audio Emotions	0.802
All Audio + All Text	0.807

Table 4: Podcast popularity prediction results using multi-modal features. The tags and MIF report the individual best representation, whereas the combination of all audio and text features proves to be more informative in characterizing a podcast.

(TR) and unique tags learned (TL) by the model over weighted-F1 scores. Tag recall is the average recall per tag. Tag recall is calculated as follows:

$$TR = \frac{\sum_{i=1}^T R_i}{|T|}$$

Here, R_i is the recall of the i^{th} tag, and $|T|$ is the total number of tags. We train every setup for a maximum of 75 epochs with early stopping criteria. For all experiments, we use a learning rate of 0.0001, batch size of 4, and Adam optimizer.

5 Results and Analysis

In this section, we report our results and analysis on the performance of hand-crafted features in different modalities for the two aforementioned downstream tasks.

5.1 Podcast Popularity Prediction

In Table 4, we present the ablation study of macro-F1 scores using the text and audio features. The tags and MIF are the most significant features in predicting podcast popularity. Their exclusion results in a performance drop of 9% and 33% w.r.t. the all text and audio features, respectively. This may be due to MIF’s ability to preserve the spectral content information (which otherwise is ignored in only energy-based computations) for the popularity prediction task. The affective features (senti-

Feature(s)	Macro-F1
Modulation + Prosody + Tags	0.820
(-) Modulation	0.520
(-) Prosody	0.803
(-) Tags	0.811
Modulation + Prosody + All Text	0.810
(-) Modulation	0.506
(-) Prosody	0.791
(-) Topics	0.798
(-) Sentiment	0.801
(-) Named Entities	0.802
(-) Tags	0.809
(-) Text Emotions	0.812

Table 5: Podcast popularity prediction results using the top two combinations of multi-modal features with the highest macro-F1 score. The modulation is the most dominant representation of all multi-modal features.

ment and text emotions) contribute equally well to the score. However, the topics and named entities negatively contribute to the prediction results. On their exclusion from all text features, we observe a marginal rise of 0.36% (topics) and 0.73% (named entities) in the macro-F1 score. We also observe an insignificant contribution from the audio emotions to the overall performance. Since the combinations of all audio and text representations report the highest macro-F1 score of 0.807, we further experiment with different combinations of these features to identify the best performing multi-modal feature combination.

Table 5 shows the top two multi-modal feature combinations (i.e., modulation + prosody + tags and modulation + prosody + all text) with the highest macro-F1 score and their ablation results. As can be seen from both the combinations, the modulation features are the most significant ones in capturing the intricacies of a podcast for this task. We believe that this performance improvement is solely due to the ability of modulation-based features in modeling the dynamic and non-linear aspects of an audio (Evangelopoulos and Maragos, 2006).

5.2 Podcast Tag Assignment

In Table 6, we present our multi-label tag assignment results for text, audio, and multi-modal frameworks with their ablation analysis. We observe that the text modality significantly outperforms audio and multi-modal frameworks with over 135% increase in weighted-F1 score. Specifically, the ‘topics’ emerge as the most prominent feature since its

Feature(s)	w-F1	TR	TL
All Text	44.02	78.22	19
(-) Topics	35.31	76.34	19
(-) LF Embeddings	41.89	71.85	19
(-) Named Entities	43.37	79.44	19
(-) Sentiment	43.42	77.07	19
(-) Text Emotions	43.83	79.34	19
All Audio	18.79	80.40	18
(-) Prosody	16.26	57.89	14
(-) MIA	17.55	51.20	15
(-) MTE	17.87	40.78	15
(-) Audio Emotions	18.68	71.84	15
(-) MFCC	19.30	63.84	17
(-) MIF	22.54	70.39	17
All Text + Audio	18.62	69.36	17
(-) LF Embeddings	17.57	70.23	16
(-) Prosody	17.71	63.86	15
(-) MTE	17.90	53.89	15
(-) Audio Emotions	17.91	54.09	12
(-) Topics	18.36	71.92	16
(-) Sentiment	18.49	52.59	15
(-) Named Entities	18.60	59.59	17
(-) MIA	18.63	58.46	15
(-) Text Emotions	18.98	65.40	17
(-) MFCC	19.26	55.45	15
(-) MIF	31.14	55.47	16

Table 6: Ablation analysis of text and audio features for multi-label tag assignment task. We use weighted-F1 (w-F1), tag recall (TR) and unique tags learned (TL) as the evaluation metrics. LF Embeddings represent the Longformer embeddings of the podcast transcript.

removal results in a 19.78% drop in the weighted-F1 score. Textual features like ‘topics’ provide a detailed list of concepts covered in the podcast. This is important in understanding *what* the podcast is about and, in turn, assigning relevant tags.

Longformer Embeddings are the next most important textual feature, followed by Named Entities, Sentiments, and Emotions. Overall, all the textual features play an essential role in tag assignment as the weighted-F1 scores drop after removing any of these features.

The audio features as standalone representations do not perform well in tag assignment tasks. Even with all the audio features, the model fails to learn all the tags. After removing MFCC and MIF, the weighted-F1 scores increase by 2.71% and 19.95%, respectively. Even though these features can capture human speech very well for the task of short-

form audio classification (Bergstra et al., 2006), they fail to provide desirable performance in the case of long-form content like podcasts. This falls in line with our hypothesis that textual features competently capture the complex nature of podcasts for the task of tag assignment. However, prosody features are the most dominant in the audio modality as with its exclusion, the model fails to learn five tags altogether while also producing a significant drop of 13.46% and 27.99% in the weighted-F1 score and tag recall, respectively. Similarly, after excluding MTE features, we see a drop of 04.89% and 49.27% in weighted-F1 and average tag recall metrics, respectively. This may be due to MTE features’ ability to retain the signal envelope variations where a speech activity is detected.

The model also does not benefit from multi-modal setup. Similar to the audio modality, the MFCC, and MIF features fail to capture the tags appropriately. We note an increase of 67.23% in the weighted-F1 score with the removal of MIF features. In a multi-modal framework, longformer embeddings perform the best, followed by audio emotions, prosody features, and sentiments. To understand why sentiments perform well for tag assignment in both text and multi-modal scenarios, we take the average percentage of negative, neutral, and positive sentiment across all podcasts under each tag. On average, 82.6% of every podcast has neutral sentiment irrespective of the tags. In Table 2 We report the difference between the average %positive and %negative sentiment per tag. We believe the model uses these differences to learn mappings to corresponding tags. We can see that the model fails to learn {‘True Crime’, ‘Fiction’, ‘Government’, and ‘Science’} tags completely in the absence of sentiments. These amount to four out of five tags with the lowest difference between average positive and negative sentiment.

5.3 Results Summary

For the popularity prediction task, the *tags* and *MIF* features are most effective for text modality and audio modality experiments, respectively. The MIF representation identifies the important events in a podcast that may impact its popularity. We observe that the model benefits from the inclusion of the tags feature. We can identify from Table 2 that some tags such as ‘True Crime’ in general are much more popular over tags like ‘Society & Culture’. The combination of all audio and text

535 features provides a more informative representation
536 in predicting the popularity of a podcast (Table
537 4). In particular, the combination of modulation,
538 prosody, and tags yields the highest macro-F1 score
539 for popularity prediction (Table 5).

540 For tag assignment, topics and Longformer em-
541 beddings are the most prominent features from text
542 modality as they effectively capture the content
543 within the podcast. The combination of all text fea-
544 tures gives us the highest weighted-F1 score. The
545 standalone audio and multi-modal frameworks give
546 a mediocre performance for multi-label tag assign-
547 ment. Prosody and affective text features are found
548 to be essential in multi-modal setups irrespective
549 of the downstream tasks.

550 6 Related Work

551 One of the first works introducing large-scale pod-
552 cast data and relevant tasks is by Clifton et al.
553 (Clifton et al., 2020) from *Spotify*. They com-
554 piled a corpus of 100,000 podcast episodes com-
555 prising nearly 60,000 hours of speech along with
556 transcriptions. Recent work in Alexander et al.
557 (2021) further enriches this dataset with precom-
558 puted audio features based on prosody and MFCCs.
559 They demonstrate how these features can be used
560 in podcast segment categorization based on deliv-
561 ery(e.g., entertaining, subjective, or discussion).
562 We extract and use similar features in our analysis
563 but on a different dataset. Earlier works based on
564 this dataset, such as abstractive summarization in
565 (Zheng et al., 2020), and PodSumm in (Vartakavi
566 and Garg, 2020) consider pre-trained models such
567 as BART, BERT, and T5. These works do not
568 specifically consider hand-crafted audio and tex-
569 tual features, and their efficacy remains relatively
570 unexplored.

571 Another similar dataset of note was compiled
572 by Yang (Yang et al., 2019) consisting of data
573 from nearly 88,728 podcast episodes on Apple
574 iTunes. Along with data, they also introduce an
575 Adversarial Learning-based Podcast Representa-
576 tion (ALPR) that captures non-textual aspects of
577 podcasts. They evaluate these representations in
578 the context of podcast popularity prediction and
579 prediction of seriousness-energy in podcasts report-
580 ing state-of-the-art results. We enhance this data
581 and use their insights to formulate our experimen-
582 tal framework. Joshi et al. (Joshi et al., 2020)
583 consider DistilBERT based embeddings as textual
584 features with the triplet loss to address the popular-

585 ity prediction task for the data introduced in (Yang
586 et al., 2019) with state-of-the-art results. They note
587 that polarity and subjectivity of features remain
588 similar with no marked difference, thus not very
589 informative for the popularity prediction. We seek
590 to investigate this further by considering various
591 hand-crafted features from both podcast audio and
592 its text transcript, along with enhanced data. While
593 the DistilBERT embeddings with triplet loss act as
594 a black box and are difficult to explain, we focus
595 on hand-crafted features for greater explainabil-
596 ity about how different features contribute towards
597 popularity.

598 (Dhanaraj and Logan, 2005) studied audio and
599 text modalities for popularity prediction of songs
600 using Support Vector Machine and boosting clas-
601 sifiers. A lot of work has been done in genre pre-
602 diction on short-form audio content using MFCCs
603 (Mandel et al., 2006)(Bergstra et al., 2006). Re-
604 cently, (Wilkes et al., 2021) performed feature anal-
605 ysis in text, audio, and video modalities for the task
606 of music genre prediction using machine learning
607 classifiers.(Cascante-Bonilla et al., 2019) use au-
608 dio, text, and video modalities from movie trailers,
609 posters, plots, and other metadata to predict movie
610 genre. However, to the best of our knowledge, no
611 one has explored multi-modal feature analysis on
612 podcasts on the task of multi-label tag assignment.
613 Inspired by these works, we seek to understand how
614 features from different modalities perform (sepa-
615 rately or combined) in the context of proxy appli-
616 cations.

617 7 Concluding Remarks

618 Podcasts are spoken-documents ranging across a
619 wide variety of genres, topics, and styles. Owing
620 to the rapid growth in popularity and global
621 reach, there is a definite need to explore and inves-
622 tigate this new engagement medium and relevant
623 research landscape. In this work, we study differ-
624 ent hand-crafted features and their combinations
625 based on podcast audio and its textual transcript
626 to characterize the podcast data. As can be seen,
627 features capturing distinct qualities like speaker
628 style, content affect, and subject content coverage
629 (i.e., prosody, sentiment, emotions, and topics) are
630 significant irrespective of the downstream task. We
631 believe that the analysis can be helpful in several
632 other downstream tasks such as podcast summa-
633 rization, retrieval, and recommendation.

References

Jordi Adell, Antonio Bonafonte, and David Escudero. 2005. Analysis of prosodic features: towards modelling of emotional and pragmatic attributes of speech. *Procesamiento del lenguaje natural*, (35):277–283.

Abigail Alexander, Matthijs Mars, Josh C Tingey, Haoyue Yu, Chris Backhouse, Sravana Reddy, and Jussi Karlgren. 2021. Audio features, precomputed for podcast retrieval and information access experiments. In *International Conference of the Cross-Language Evaluation Forum for European Languages*, pages 3–14. Springer.

Anchor. 2022. The easiest way to make a podcast. *accessed 11th January 2022*.

Dimo Angelov. 2020. Top2vec: Distributed representations of topics. *arXiv preprint arXiv:2008.09470*.

Alexei Baevski, Henry Zhou, Abdelrahman Mohamed, and Michael Auli. 2020. wav2vec 2.0: A framework for self-supervised learning of speech representations. *arXiv preprint arXiv:2006.11477*.

Iz Beltagy, Matthew E Peters, and Arman Cohan. 2020. Longformer: The long-document transformer. *arXiv preprint arXiv:2004.05150*.

James Bergstra, Norman Casagrande, Dumitru Erhan, Douglas Eck, and Balázs Kégl. 2006. Aggregate features and a da b oost for music classification. *Machine learning*, 65(2-3):473–484.

Boualem Boashash. 1992. Estimating and interpreting the instantaneous frequency of a signal. i. fundamentals. *Proceedings of the IEEE*, 80(4):520–538.

Paul Boersma and Vincent Van Heuven. 2001. Speak and unspeak with praat. *Glott International*, 5(9/10):341–347.

Paola Cascante-Bonilla, Kalpathy Sitaraman, Mengjia Luo, and Vicente Ordonez. 2019. Moviescope: Large-scale analysis of movies using multiple modalities. *arXiv preprint arXiv:1908.03180*.

Nitesh V Chawla, Kevin W Bowyer, Lawrence O Hall, and W Philip Kegelmeyer. 2002. Smote: synthetic minority over-sampling technique. *Journal of artificial intelligence research*, 16:321–357.

Ann Clifton, Sravana Reddy, Yongze Yu, Aasish Pappu, Rezvaneh Rezapour, Hamed Bonab, Maria Eskevich, Gareth Jones, Jussi Karlgren, Ben Carterette, et al. 2020. 100,000 podcasts: A spoken english document corpus. In *Proceedings of the 28th International Conference on Computational Linguistics*, pages 5903–5917.

Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. Bert: Pre-training of deep bidirectional transformers for language understanding. In *Proceedings of the 2019 Conference of the*

North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers), pages 4171–4186. 686
687
688
689

Ruth Dhanaraj and Beth Logan. 2005. Automatic prediction of hit songs. In *ISMIR*, pages 488–491. Cite-seer. 690
691
692

Dimitrios Dimitriadis, Petros Maragos, and Alexandros Potamianos. 2005. Robust am-fm features for speech recognition. *IEEE signal processing letters*, 12(9):621–624. 693
694
695
696

Georgios Evangelopoulos and Petros Maragos. 2006. Multiband modulation energy tracking for noisy speech detection. *IEEE Transactions on audio, speech, and language processing*, 14(6):2024–2038. 697
698
699
700

Kishaloy Halder, Alan Akbik, Josip Krapac, and Roland Vollgraf. 2020. Task-aware representation of sentences for generic text classification. In *Proceedings of the 28th International Conference on Computational Linguistics*, pages 3202–3213. 701
702
703
704
705

Nielsen Insights. 2020. Podcast content is growing audio engagement. *accessed 11th January 2022*. 706
707

Rosie Jones, Hamed Zamani, Markus Schedl, Ching-Wei Chen, Sravana Reddy, Ann Clifton, Jussi Karlgren, Helia Hashemi, Aasish Pappu, Zahra Nazari, et al. 2021. Current challenges and future directions in podcast information access. *arXiv preprint arXiv:2106.09227*. 708
709
710
711
712
713

Brihi Joshi, Shravika Mittal, and Aditya Chetan. 2020. Did you “read” the next episode? using textual cues for predicting podcast popularity. In *Proceedings of the 1st Workshop on NLP for Music and Audio (NLP4MusA)*, pages 13–17. 714
715
716
717
718

Sudipta Kar, Suraj Maharjan, and Tamar Solorio. 2018. Folksonomication: Predicting tags for movies from plot synopses using emotion flow encoded neural network. In *Proceedings of the 27th International Conference on Computational Linguistics*, pages 2879–2891. 719
720
721
722
723
724

Michael I Mandel, Graham E Poliner, and Daniel PW Ellis. 2006. Support vector machine active learning for music retrieval. *Multimedia systems*, 12(1):3–13. 725
726
727

Petros Maragos, James F Kaiser, and Thomas F Quatieri. 1993. Energy separation in signal modulations with application to speech analysis. *IEEE transactions on signal processing*, 41(10):3024–3051. 728
729
730
731

Nielsen. September 2021. Insights for podcast advertisers. *Podcasting Today*. 732
733

Jeffrey Pennington, Richard Socher, and Christopher D Manning. 2014. Glove: Global vectors for word representation. In *Proceedings of the 2014 conference on empirical methods in natural language processing (EMNLP)*, pages 1532–1543. 734
735
736
737
738

739 Riverside.fm. 2022. Record podcasts and videos from
740 anywhere. *accessed 11th January 2022*.

741 Victor Sanh, Lysandre Debut, Julien Chaumond, and
742 Thomas Wolf. 2019. Distilbert, a distilled version
743 of bert: smaller, faster, cheaper and lighter. *arXiv*
744 *preprint arXiv:1910.01108*.

745 Ottokar Tilk and Tanel Alumäe. 2016. Bidirectional
746 recurrent neural network with attention mechanism
747 for punctuation restoration. In *Interspeech 2016*.

748 Chai M Tyng, Hafeez U Amin, Mohamad NM Saad,
749 and Aamir S Malik. 2017. The influences of emotion
750 on learning and memory. *Frontiers in psychology*,
751 8:1454.

752 Aneesh Vartakavi and Amanmeet Garg. 2020.
753 Podsumm–podcast audio summarization. *arXiv*
754 *preprint arXiv:2009.10315*.

755 Ben Wilkes, Igor Vatolkin, and Heinrich Müller. 2021.
756 Statistical and visual analysis of audio, text, and im-
757 age features for multi-modal music genre recognition.
758 *Entropy*, 23(11):1502.

759 Min Xu, Jinqiao Wang, Xiangjian He, Jesse S Jin,
760 Suhuai Luo, and Hanqing Lu. 2014. A three-level
761 framework for affective content analysis and its
762 case studies. *Multimedia tools and applications*,
763 70(2):757–779.

764 Longqi Yang, Yu Wang, Drew Dunne, Michael Sobolev,
765 Mor Naaman, and Deborah Estrin. 2019. More than
766 just words: Modeling non-textual characteristics of
767 podcasts. In *Proceedings of the Twelfth ACM Interna-*
768 *tional Conference on Web Search and Data Mining*,
769 pages 276–284.

770 Shu-wen Yang, Po-Han Chi, Yung-Sung Chuang,
771 Cheng-I Jeff Lai, Kushal Lakhota, Yist Y Lin,
772 Andy T Liu, Jiatong Shi, Xuankai Chang, Guan-
773 Ting Lin, et al. 2021. Superb: Speech processing
774 universal performance benchmark. *arXiv preprint*
775 *arXiv:2105.01051*.

776 Chujie Zheng, Harry Jiannan Wang, Kunpeng Zhang,
777 and Ling Fan. 2020. A baseline analysis for pod-
778 cast abstractive summarization. *arXiv preprint*
779 *arXiv:2008.10648*.

780 Athanasia Zlatintsi, Petros Maragos, Alexandros
781 Potamianos, and Georgios Evangelopoulos. 2012. A
782 saliency-based approach to audio event detection and
783 summarization. In *2012 Proceedings of the 20th Eu-*
784 *ropean Signal Processing Conference (EUSIPCO)*,
785 pages 1294–1298. IEEE.