

# PROBLEM STATEMENT



Reduce losses by bidding correctly



Predict popularity based on characteristics of music



Bid valid if popularity predicted is higher or equal to actual

POPULARITY		BID PRICE	EXPECTED REVENUE (in 10k \$)
	VERY LOW	1	2
	LOW	2	4
	AVERAGE	3	6
	HIGH	4	8
	VERY HIGH	5	10

Acounsticness	Mode
Danceability	Release Date
Instrumentalness	Speechiness
Key	Tempo
Energy	Valence
Liveliness	Year
Explicit	Duration min
Loudness	Popularity

# DATASET



NUMBER OF TRAINING EXAMPLES

12227

NUMBER OF FEATURES

16

CATEGORICAL FEATURES

2

NUMERIC FEATURES

14

MISSING VALUES

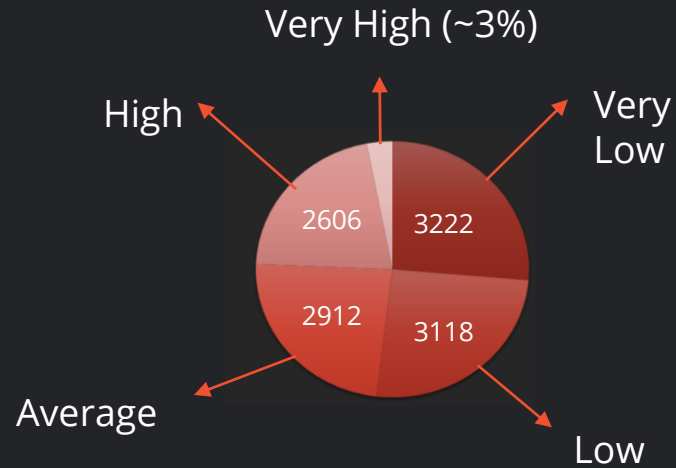
NONE

OUTLIERS

NONE  
REMOVED

DUPLICATE VALUES

NONE

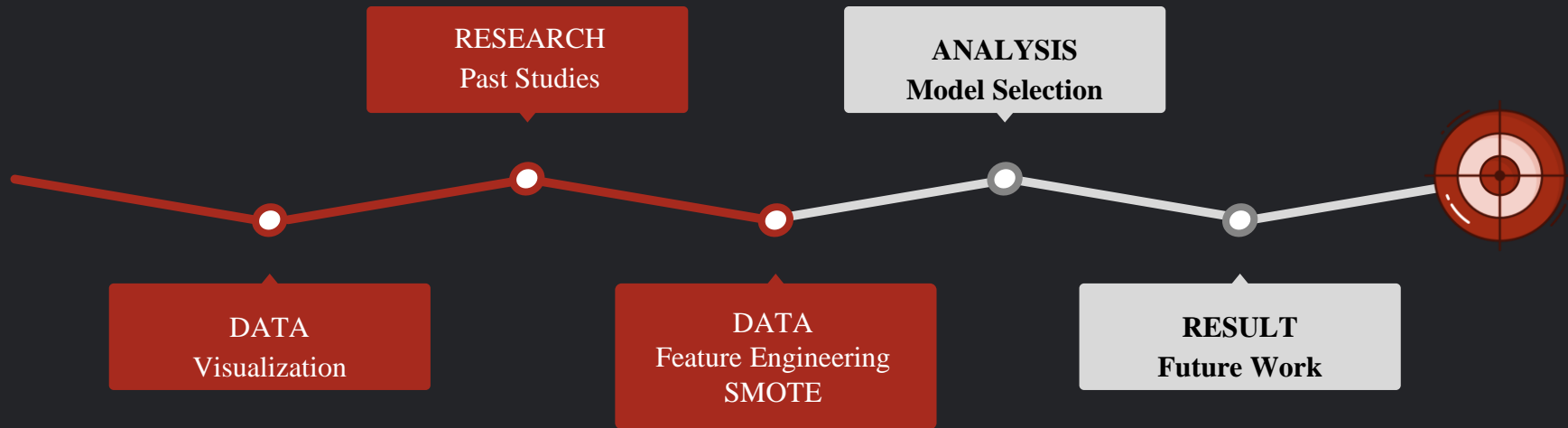


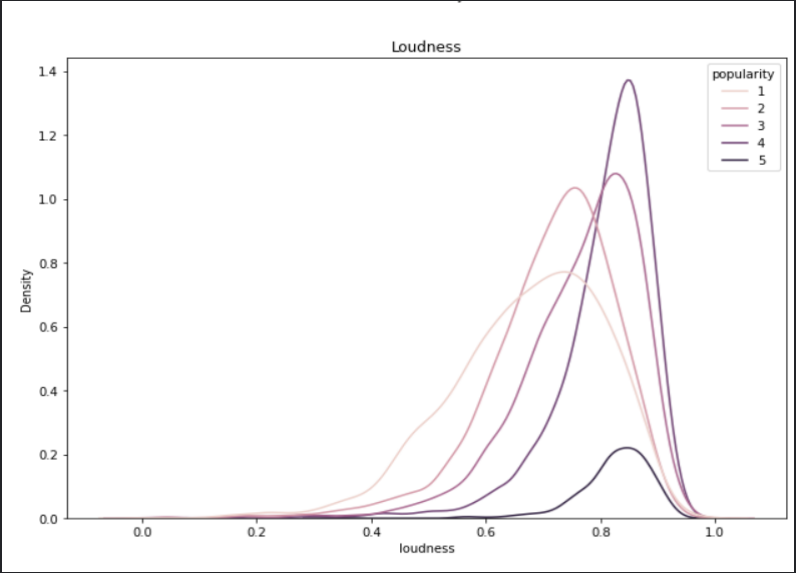
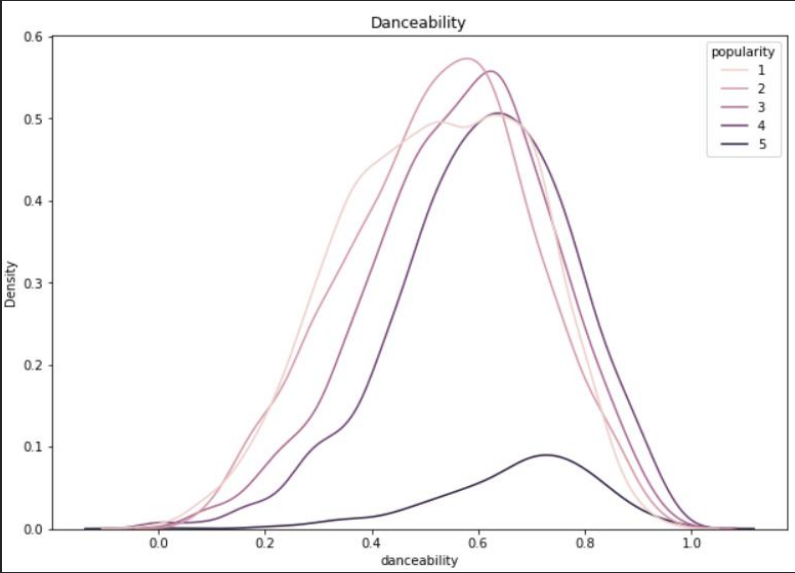
Oversampling (SMOTE)



Class Imbalance

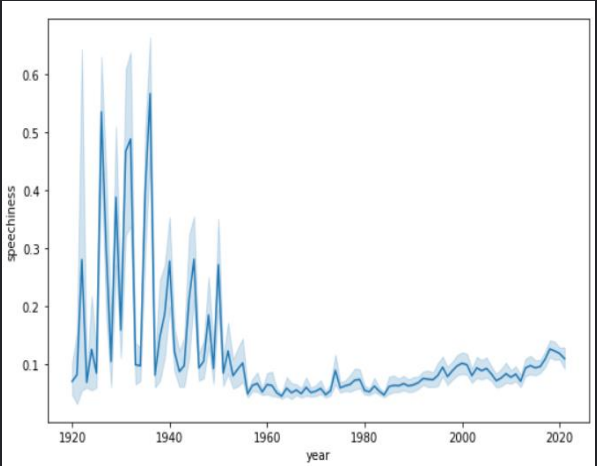
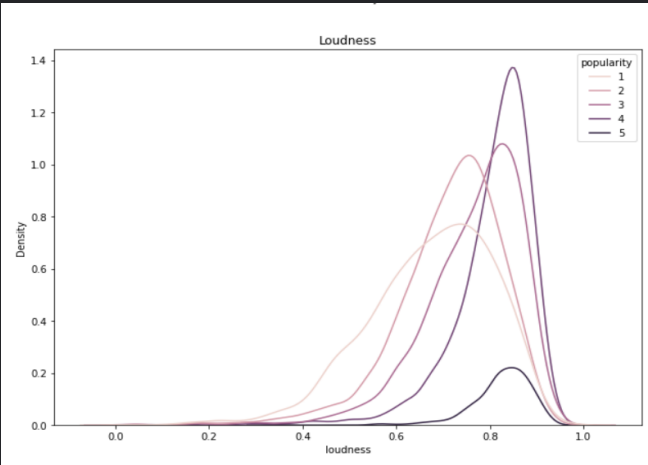
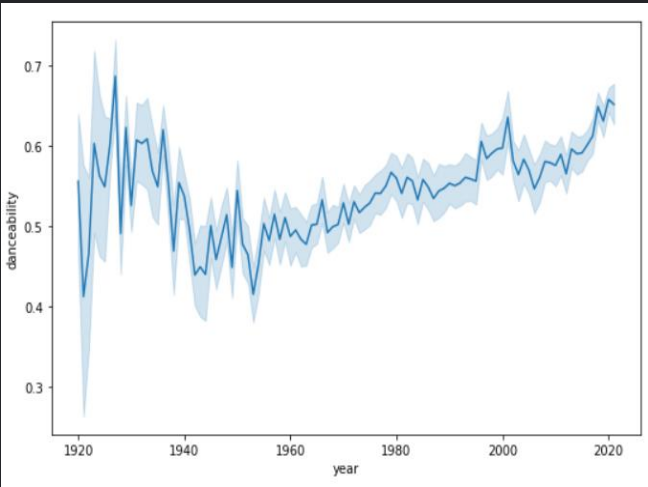
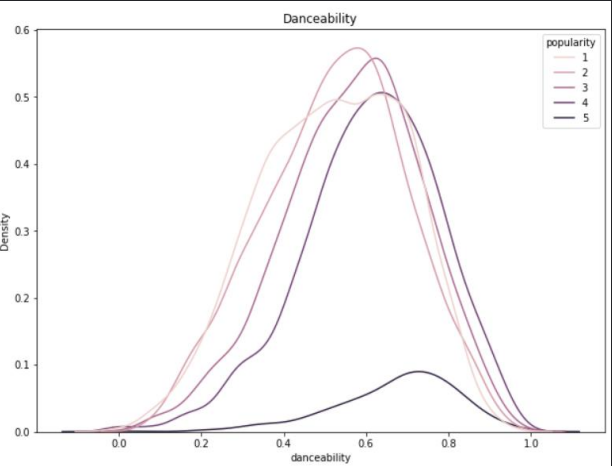
# BASIC WORKFLOW



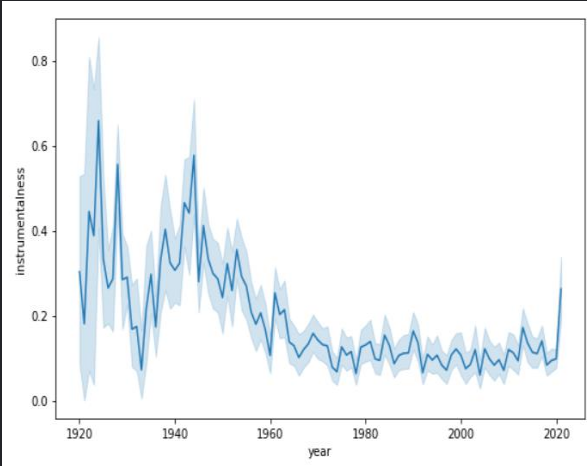
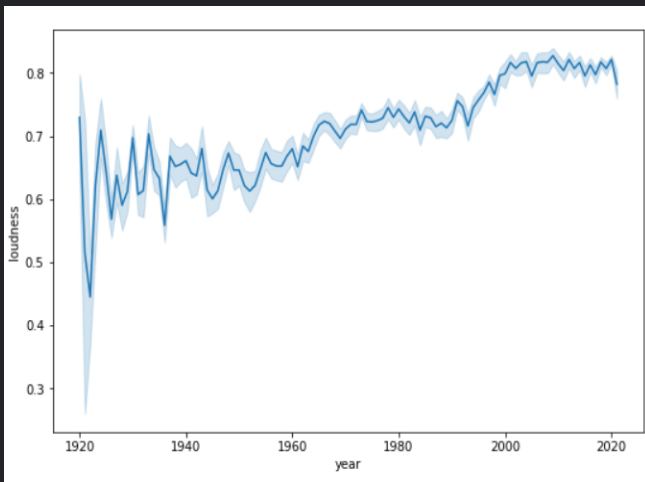
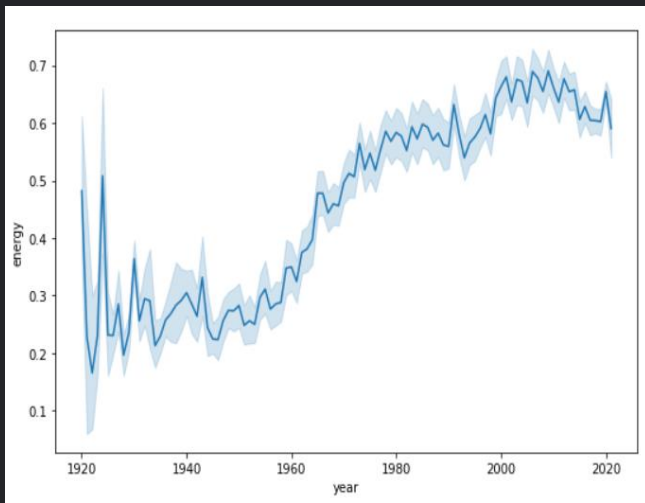
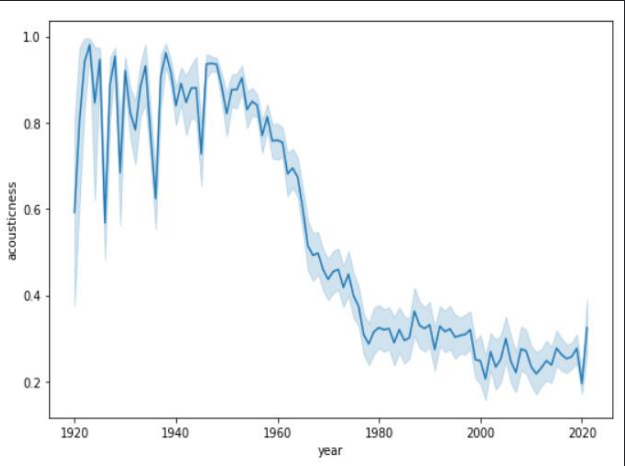


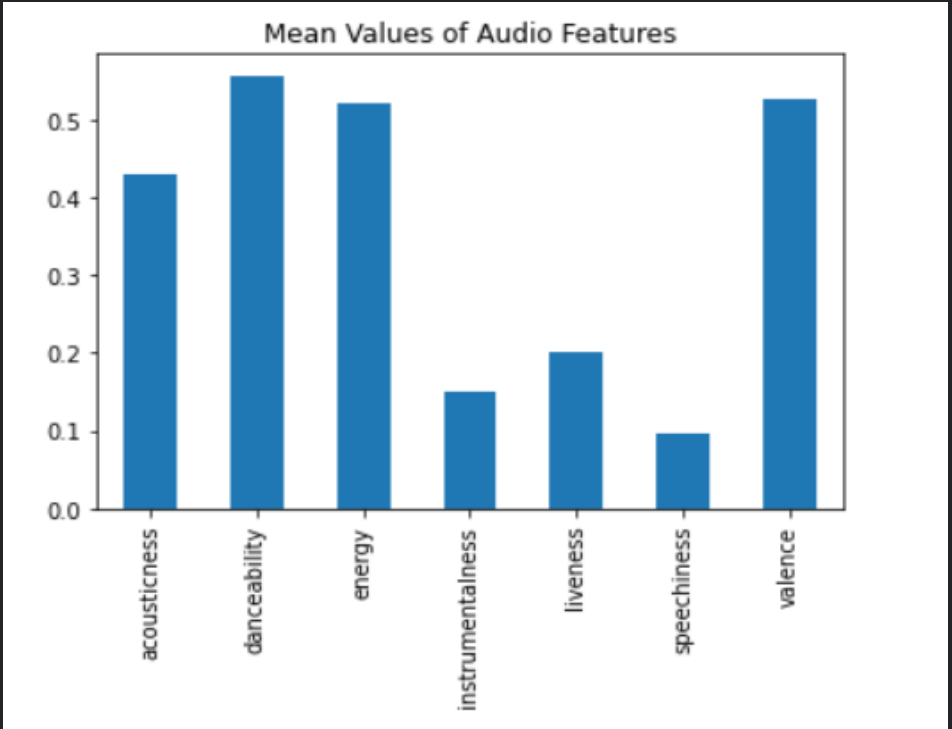
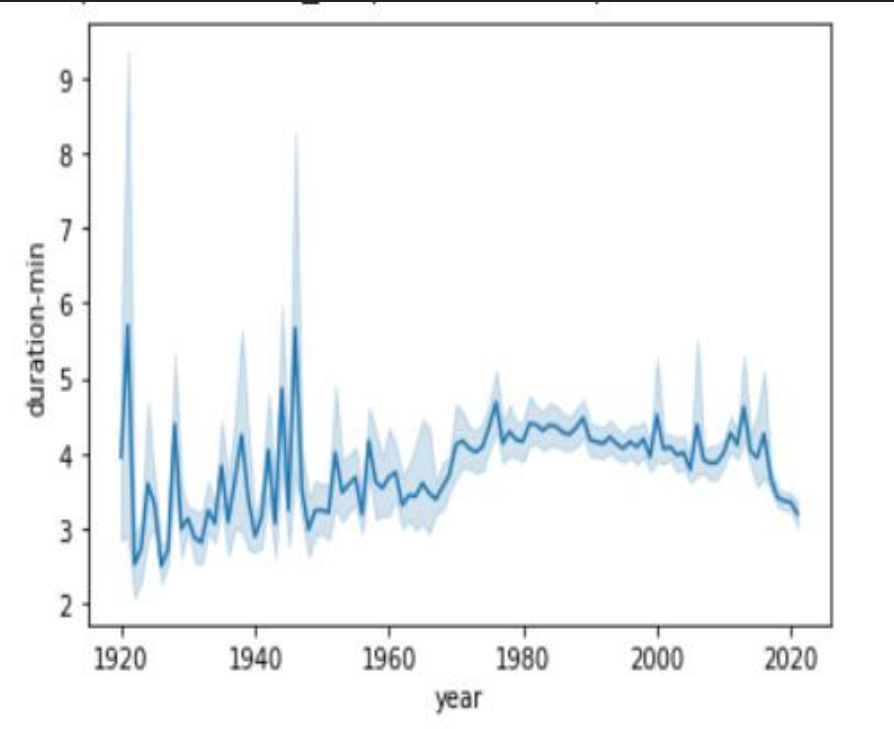
Danceability v/s Popularity

# EXPLORATORY DATA ANALYSIS

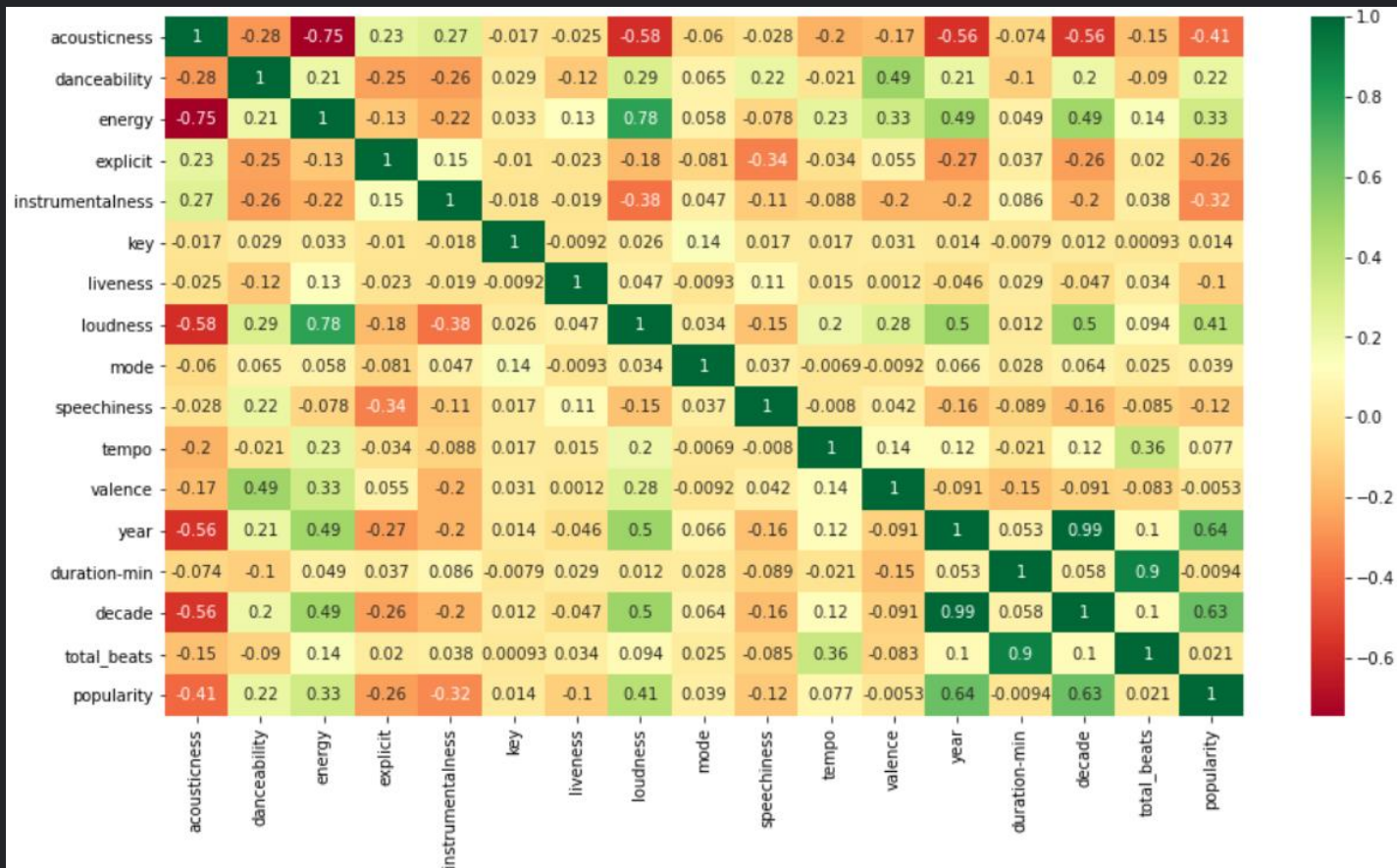


# EXPLORATORY DATA ANALYSIS





# HEATMAP





# FEATURE ENGINEERING



Duration

Tempo



Total Beats

Year



Decade

## WORKFLOW OF MODELS



**Basic  
Model**

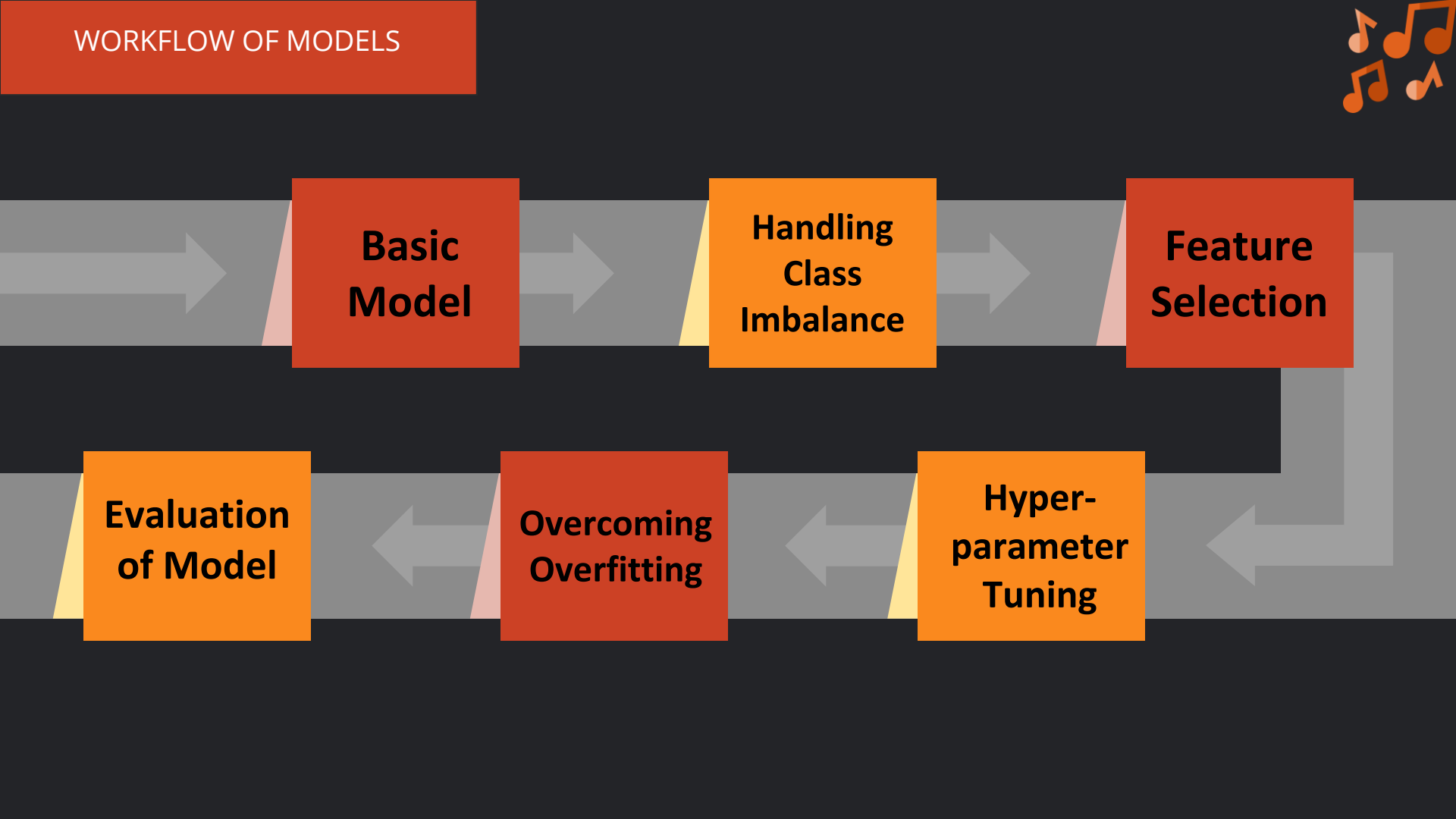
**Handling  
Class  
Imbalance**

**Feature  
Selection**

**Evaluation  
of Model**

**Overcoming  
Overfitting**

**Hyper-  
parameter  
Tuning**



# SUPPORT VECTOR CLASSIFICATION



## Baseline Accuracy:

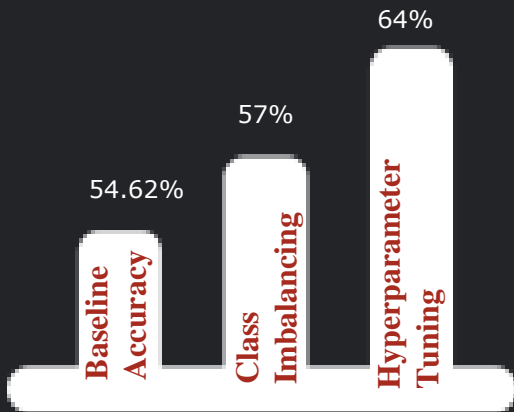
- Removed ID and Release Date
- Trained with default parameter
- Kernel used – Linear

Cross Validation Accuracy : 54.62 %

## Class Imbalancing

- Applied SMOTE
- Trained with default parameter.
- Kernel used-Polynomial

Cross Validation Accuracy : 57 %



## FEATURE SELECTION

Features Dropped	1.id 2.release_date 3.key 4.model
Features Added	total_beats (tempo*duration)

## HYPERPARAMETER TUNING

Kernel	Polynomial
Degree of Polynomial	8
C	1
Cross Val Accuracy	64%

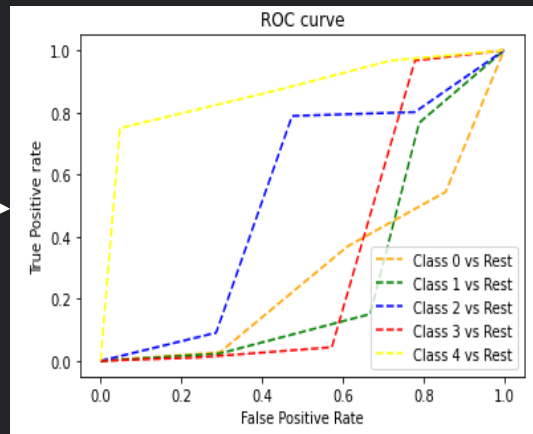
# SUPPORT VECTOR CLASSIFICATION



## CONFUSION MATRIX

278	192	182	0	39
118	453	70	0	15
173	19	444	0	169
2	85	3	0	1
13	183	41	0	577

## ROC-AUC CURVE



	precision	recall	f1-score
0	0.48	0.40	0.44
1	0.49	0.69	0.57
2	0.60	0.55	0.57
3	0.00	0.00	0.00
4	0.72	0.71	0.71
Accuracy			0.57

# XGBOOST



## Baseline Accuracy:

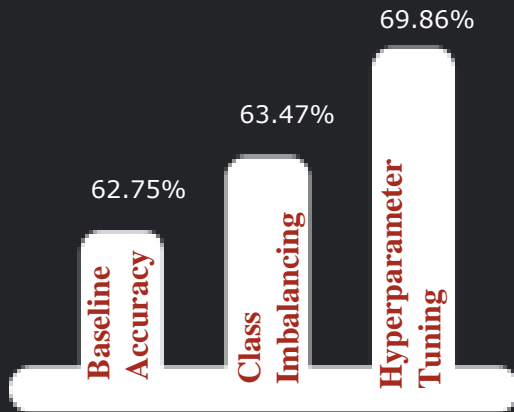
- Removed ID and Release Date
- No oversampling

Cross Val Accuracy : 62.75 %

## Class Imbalancing

- Applied SMOTE
- Removed ID, release date

Cross Val Accuracy : 63.47 %



## FEATURE SELECTION

Features Dropped	1.id 2.release_date 3.key 4.model
Features Added	-total_beats (tempo*duration) -decade

## HYPERPARAMETER TUNING

colsample_bytree	0.7
gamma	0.3
learning_rate	0.15
max_depth	15
n_estimators	1000
min_child_weight	1
Cross Val Accuracy	69.86%

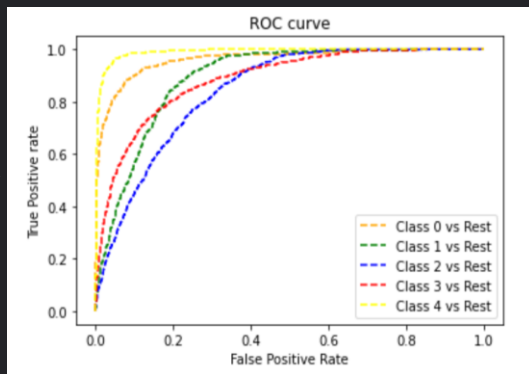
# XGBOOST



## CONFUSION MATRIX

766	94	39	46	10
103	605	194	17	0
35	309	490	137	7
30	106	203	529	123
4	3	11	45	97

## ROC-AUC CURVE



	precision	recall	f1-score
0	0.817	0.802	0.809
1	0.542	0.658	0.594
2	0.523	0.501	0.512
3	0.683	0.534	0.599
4	0.869	0.936	0.901
Accuracy			0.686



## Baseline Accuracy:

- Removed ID and Release Date
- boosting type : gdbt

Cross Val Accuracy : 62.00 %

## Class Imbalancing

- Applied SMOTE
- Removed ID, Release Date
- boosting type : dart

Cross Val Accuracy : 62.09 %

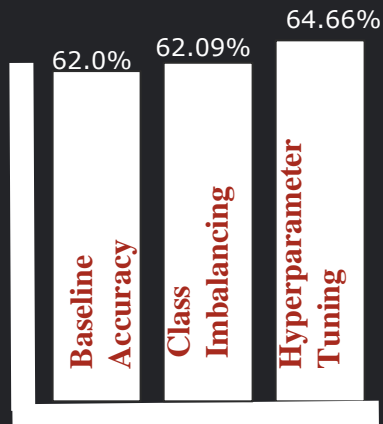
## HYPERPARAMETER TUNING

learning_rate	0.01
early_stopping	30
lambda_l1	1.0
lambda_l2	1.0
num_boost_rounds	2000
Cross Val Accuracy=64.66%	

## FEATURE SELECTION

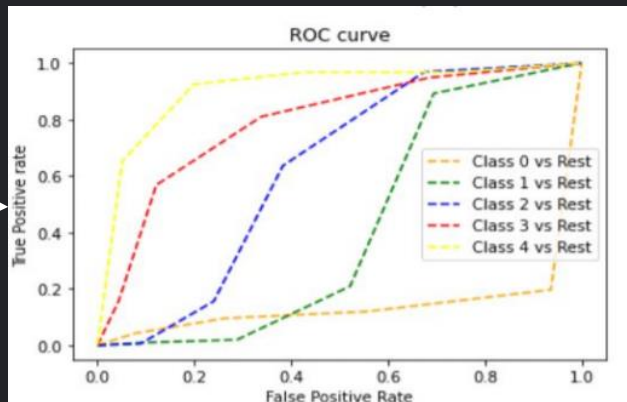
### Features Dropped

1.id  
2.release\_date  
3.key  
4.model  
5.explicit



**CONFUSION  
MATRIX**

638	61	19	42	34
83	526	146	8	7
23	240	347	109	4
35	94	164	278	107
3	0	4	25	60

**ROC-AUC  
CURVE**

Class	Precision	Recall	f1-score
Very Low	0.816	0.800	0.810
Low	0.571	0.683	0.622
Average	0.510	0.480	0.495
High	0.602	0.410	0.488
Very High	0.283	0.652	0.604



# ADABOOST



## Baseline Accuracy:

- Removed ID and Release Date
- No oversampling

Cross Val Accuracy : 59 %

## Class Imbalancing

- Applied SMOTE
- Removed ID, Release Date

Cross Val Accuracy : 57.4 %

## HYPERPARAMETER TUNING

n_estimators	learning_rate	Accuracy
500	0.1	60.88
100	0.2	60.57
100	0.3	61.01
100	0.4	68.60
Cross Val Accuracy=64%		

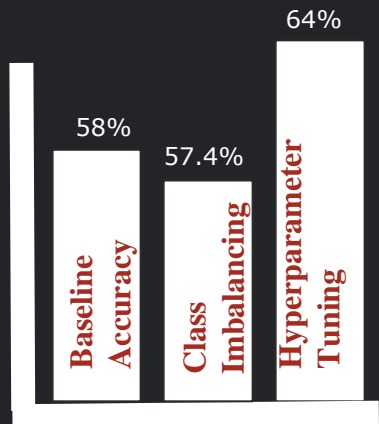
## FEATURE SELECTION

### Features Dropped

1.id  
2.release\_date  
3.key  
4.model  
5.explicit

### Features Added

total\_beats  
(tempo\*duration)



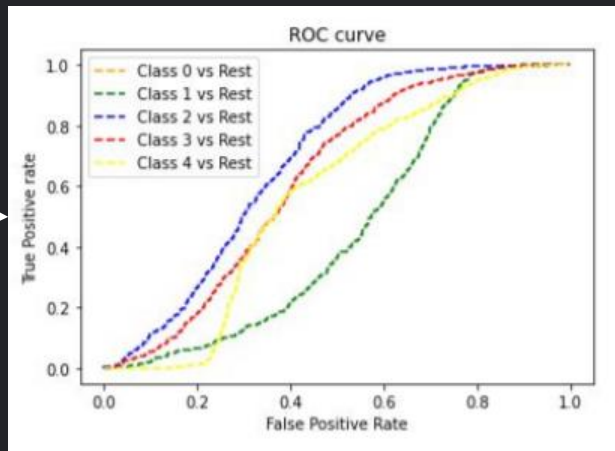
# ADBOOST



## CONFUSION MATRIX

530	67	22	51	9
68	408	137	12	0
17	190	309	94	0
11	64	117	401	55
1	3	7	24	624

## ROC-AUC CURVE



Class	Precision	Recall	f1-score
Very Low	0.845	0.781	0.812
Low	0.557	0.653	0.601
Average	0.522	0.507	0.514
High	0.688	0.619	0.652
Very High	0.907	0.945	0.926

## RANDOM FOREST CLASSIFIER



01

Features removed: key, mode, explicit, ID, release\_date

02

Oversampling resulted in overfitting

03

Applied repeated stratified k-fold cross validation

04

Hyperparameters tuned

**Accuracy of baseline model = ~50%**



## HYPERPARAMETER TUNING



NO. OF ESTIMATORS	MEAN ACCURACY
100	~68%
500	~70%
500-900	~70%(run time was increasing)



## HANDLING CLASS IMBALANCE

TOOLS USED	MEAN ACCURACY
Cost-sensitive learning	~50%
Default SMOTE oversampling	~68%



## OVERSAMPLING STRATEGY

NO. OF EXAMPLES IN EACH CLASS	MEAN ACCURACY
3700	~72%
4000	~74%
5000	~77%

# RANDOM FOREST CLASSIFIER

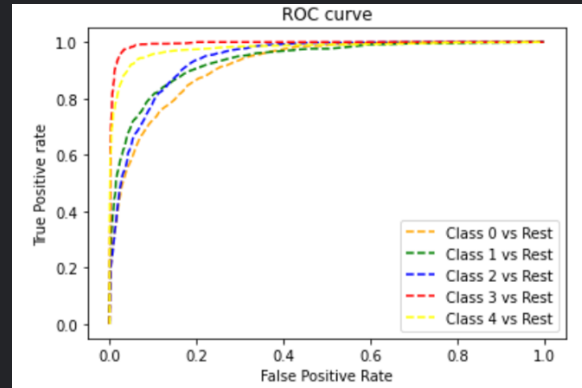


## CONFUSION MATRIX

639	103	203	0	23
147	718	85	66	10
151	10	727	0	69
8	32	8	991	0
28	41	68	8	865

	precision	recall	f1-score
0	0.657	0.660	0.658
1	0.794	0.700	0.744
2	0.666	0.760	0.710
3	0.931	0.954	0.942
4	0.895	0.856	0.875

## ROC-AUC CURVE



## FINAL HYPERPARAMETERS

HYPERPARAMETERS	TUNED VALUE
Estimators	100
Splits	10
Repeats	3
Cross Val Accuracy	~79%

**MODEL 1**

Features removed	Explicit, release_date, tempo, key, year
New Features	Total Beats, Duration-min
Scaling	Standard Scaler
Oversampling	SMOTE

Number of hidden layers	3 (5->32->16->8->5)
Activation in hidden layers	relu
Activation in output layer	softmax
Loss function	Categorical cross entropy
optimizer	Adam
batch size	5
epochs	2000
Cross validation accuracy	53%

**MODEL 2**

Features removed	Explicit, release_date, tempo, key, year
New Features	Total Beats, Duration-min
Scaling	Standard Scaler
Oversampling	SMOTE

Number of hidden layers	5 (5->128->64->32->16->8->5)
Activation in hidden layers	relu
Activation in output layer	softmax
Loss function	Categorical cross entropy
optimizer	Adam
batch size	5
epochs	2000
Cross validation accuracy	65%

## COMPARISON



64%



**SUPPORT  
VECTOR  
MACHINE**

69%



**XGBOOST**

64%



**LightGBM**

68%



**ADABOOST**

79%



**RANDOM  
FOREST**

65%



**NEURAL  
NETWORK**

## WHY

- **Kernel Trick**
- Don't need to hand-perform complex transformation
- Can capture non linear relationships due to this

- **Boosting**
- Has been observed to give excellent accuracy
- Minimal data preprocessing required
- Faster runtime

- Very fast execution compared to XGB
- Leaf wise growth, can lead to better accuracy

- **Boosting** has been observed to give superior performance

- The basis is Bagging
- Bagging reduces variance and solves overfitting

- Neural Networks are great at capturing patterns and generating decision boundaries
- Might eliminate the need to hand engineer features





WHY DON'T THESE WORK WELL?	
SVM	Impacted greatly by correct choice of hyperparameters Hyperparameter tuning is tricky and time consuming Overfitting due to Oversampling Possibly due to presence of outliers
Boosting Algorithms	Boosting algorithms are prone to overfitting Overfitting due to Oversampling
Neural Networks	Overfits due to oversampling and relatively smaller dataset

WHY DOES RANDOM FOREST WORK WELL?
<ul style="list-style-type: none"><li>- Random forest - Bagging algorithm</li><li>- Biggest problem - overfitting due to oversampling</li><li>- Bagging reduces variance - model less prone to overfitting</li><li>- Robust to outliers, requires min preprocessing</li><li>- Faster runtime allows extensive hyperparameter tuning</li></ul>

## FUTURE WORK



1	Alternate methods for dealing with class imbalance
2	Data about song artists and producers
3	Careful removal of univariate and multivariate outliers
4	Meticulous hyperparameter tuning, especially for SVM



THANK  
YOU!