

# Model dubokog učenja za klasifikaciju MIDI datoteka

---

**Rod, Bartol**

**Master's thesis / Diplomski rad**

**2024**

*Degree Grantor / Ustanova koja je dodijelila akademski / stručni stupanj:* **University of Zagreb, Faculty of Electrical Engineering and Computing / Sveučilište u Zagrebu, Fakultet elektrotehnike i računarstva**

*Permanent link / Trajna poveznica:* <https://um.nsk.hr/um:nbn:hr:168:251682>

*Rights / Prava:* [In copyright](#)/[Zaštićeno autorskim pravom.](#)

*Download date / Datum preuzimanja:* **2025-03-20**



*Repository / Repozitorij:*

[FER Repository - University of Zagreb Faculty of Electrical Engineering and Computing repository](#)



SVEUČILIŠTE U ZAGREBU  
FAKULTET ELEKTROTEHNIKE I RAČUNARSTVA

DIPLOMSKI RAD br. 369

**MODEL DUBOKOG UČENJA ZA KLASIFIKACIJU MIDI  
DATOTEKA**

Bartol Rod

Zagreb, lipanj 2024.

SVEUČILIŠTE U ZAGREBU  
FAKULTET ELEKTROTEHNIKE I RAČUNARSTVA

DIPLOMSKI RAD br. 369

**MODEL DUBOKOG UČENJA ZA KLASIFIKACIJU MIDI  
DATOTEKA**

Bartol Rod

Zagreb, lipanj 2024.

## DIPLOMSKI ZADATAK br. 369

Pristupnik: **Bartol Rod (0036508261)**

Studij: Računarstvo

Profil: Računarska znanost

Mentor: doc. dr. sc. Stjepan Šebek

Zadatak: **Model dubokog učenja za klasifikaciju MIDI datoteka**

### Opis zadatka:

U ovom diplomskom radu zadatak je iskoristiti modele dubokog učenja u svrhu klasifikacije MIDI datoteka, koje su uobičajeni formati za glazbene zapise. Cilj je razviti model koji može efikasno razvrstati MIDI datoteke prema klasičnim glazbenicima. Koristi se duboka neuronska mreža koja se trenira na velikom skupu raznolikih MIDI zapisa iz skupa podataka MAESTRO. Evaluacija će se provesti usporedbom performansi različitih arhitektura dubokih modela i analizom njihove sposobnosti generalizacije na novim glazbenim podacima. Krajnji cilj je razvoj efikasnog sustava za automatsku klasifikaciju glazbenih datoteka koji bi imao određenu primjenu u glazbenoj industriji, istraživačkim projektima ili umjetničkoj domeni.

Rok za predaju rada: 28. lipnja 2024.

*Hvala na čekanju.*

# Sadržaj

<b>1. Uvod</b>	<b>3</b>
<b>2. MIDI protokol</b>	<b>5</b>
2.1. Tehničke karakteristike	5
2.1.1. MIDI portovi i kablovi	5
2.1.2. MIDI poruke	6
2.2. MIDI datoteka	7
<b>3. Skup podataka GiantMIDI-Piano</b>	<b>8</b>
3.1. Zastupljenost podataka	8
3.2. Odabir skladatelja	9
3.3. Podjela podataka	11
<b>4. Stroj potpornih vektora i klasifikator nasumičnih šuma</b>	<b>12</b>
4.1. Ekstrakcija značajki	12
4.1.1. <i>Pretty_midi</i>	12
4.1.2. Analiza glavnih komponenata	13
4.2. Učenje klasifikatora	15
4.2.1. Stroj potpornih vektora	16
4.2.2. Klasifikator slučajnih šuma	19
4.3. Rezultati	21
4.3.1. Matrice zabune	23
<b>5. Klasifikacija metodama obrade prirodnog teksta</b>	<b>27</b>
5.1. Tokenizacija	28
5.1.1. REMI	28

5.1.2. MidiTok . . . . .	29
5.2. Povratne neuronske mreže . . . . .	32
5.2.1. LSTM i GRU . . . . .	33
5.2.2. Pozornost . . . . .	35
5.3. Učenje modela . . . . .	36
5.4. Rezultati . . . . .	40
5.4.1. Matrice zabune . . . . .	42
<b>6. Zaključak . . . . .</b>	<b>44</b>
<b>Literatura . . . . .</b>	<b>45</b>
<b>Sažetak . . . . .</b>	<b>46</b>
<b>Abstract . . . . .</b>	<b>47</b>

# 1. Uvod

Povijest prijenosa i zapisa glazbenih djela odražava tehnološki i kulturološki napredak ljudskog društva. Od usmene predaje, do zapisa u obliku pisma, tiskarskog stroja i digitalnim medijima glazba se prenosi i čuva iz generacije u generaciju prateći napredak tehnologije i kulture. Svrha zapisivanja glazbe je raznolika; radi li se o edukaciji, zabavi ili umjetnosti, svaki oblik prijenosa glazbe pronade svoj medij. Danas možemo reći da se prijenos glazbe dijeli u dvije vrste medija: oni namijenjeni slušatelju i oni namijenjeni glazbeniku. Prvi tip medija prenosi zvučne valove s ciljem što vjernijeg dostavljanja glazbe slušatelju. To uključuje različite formate kao što su CD, digitalni audio formati poput Wav-a ili MP3-a, te streaming servisi. Drugi tip medija nastoji glazbu zapisati kako bi se interpretirala i izvodila, te je namijenjen izvođačima. To uključuje notne zapise, tabulature ili digitalne formate poput MIDI, engl. *Musical Instrument Digital Interface*, protokola.

MIDI protokol nije samo medij za prijenos glazbenog zapisa, već je primarno namijenjen računalima i definira skup instrukcija kojima se određuje ponašanje glazbenog hardvera. MIDI omogućava komunikaciju između različitih elektroničkih instrumenata i računala, omogućujući kontrolu nad različitim aspektima zvuka, uključujući visinu tona, jačinu, trajanje i dinamiku. Također, omogućava i programiranje složenih glazbenih aranžmana te sekvenciranje, čineći ga ključnim alatom u modernoj glazbenoj produkciji. Melodije i ritam, tj. sve ono što obuhvaća zapis glazbe u njenom klasičnom smislu samo su dio MIDI protokola.

Ljudska potreba za strukturom i uređenjem jedan su od osnovnih problema koji čovjek pokušava riješiti. Podjela glazbe u žanrove ili neku drugu kategoriju često je predmet rasprave i nije jasno koje djelo svrstati u koju kategoriju. Klasifikacija MIDI datoteka stoga se može svesti i na klasifikaciju informacija o glazbi zapisanoj u njima. To uklju-



čuje analizu parametara kao što su tempo, skala, ključ i aranžman koji se mogu dobiti iz informacije sadržane u MIDI datoteci.

U ovom radu prvo se opisuje MIDI protokol i njegova aplikacija te MIDI datoteke i na koji način je u njima pohranjena glazbena informacija. Zatim je predstavljen skup podataka *GiantMIDI-Piano* [1] i kako će se koristiti te kriterij klasifikacije klasičnih skladatelja. Također, pokazano je i kako u programskom jeziku *Python* iz MIDI datoteke izlučiti bitne značajke glazbenog zapisa. Istraživanje je podijeljeno u dva dijela: u prvom se provodi klasifikacija korištenjem algoritama strojnog učenja, a u drugom metodama dubokog učenja. Klasifikacija strojnim učenjem istražena je korištenjem SVM-ova, engl. *Support Vector Machines* i RF klasifikatorom, engl. *Random Forest Classifier*, a zatim se problem svodi na NLP, engl. *Natural Language Processing*, kao što je i predloženo u radu [2] i uče duboke neuronske mreže, RNN-ovi, eng. *Recurrent Neural Network*. Nakon provedenih metodi rezultati su pokazani u obliku usporedbe preciznosti (engl. *accuracy*) i *F1-score*-a.

## 2. MIDI protokol

MIDI, engl. *Musical Instrument Digital Interface* tehnički je standard koji opisuje komunikacijski protokol, digitalno sučelje i konektore koji povezuju razne elektroničke glazbene instrumente, računala i druge povezane uređaje u svrhu stvaranja, uređivanja i reprodukcije glazbe. Dakle, MIDI sam po sebi ne prenosi zvuk, već informacije o glazbenim notama, tempu, dinamici i drugim parametrima koji su potrebni za reprodukciju glazbenog djela te definira specifične parametre za prijenos informacija. Razvijen je početkom 1980-ih kao odgovor na potrebu za standardizacijom komunikacije između različitih glazbenih instrumenata i uređaja. Prije MIDI-a, različiti proizvođači imali su svoje vlastite sustave koji nisu bili međusobno kompatibilni. Prvi put javno je predstavljen 1983. godine na godišnjem Sajmu glazbenih instrumenata, NAMM, u Anaheimu, Kalifornija i rezultat je suradnje glavnih proizvođača elektronskih instrumenata, uključujući Roland, Yamaha, Korg i druge.

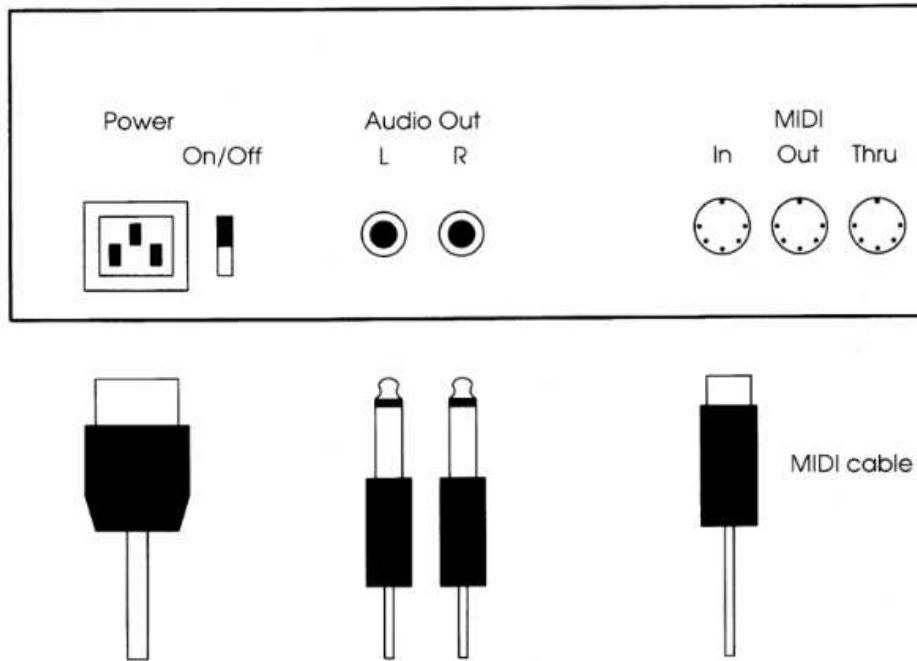
### 2.1. Tehničke karakteristike

Tehničke karakteristike MIDI protokola mogu se svrstati u dvije grupe zahtjeva koje proizvođač mora ispuniti [3, *The MIDI Specification*]. Prvi zahtjevi odnose se na fizičku vezu i tu spadaju MIDI Portovi i MIDI kablovi [3, *The MIDI Specification*]. Drugi zahtjevi se odnose na podatkovne strukture i njihovu uređenost u protokolu te ih možemo nazvati MIDI porukama [3, *The MIDI Specification*].

#### 2.1.1. MIDI portovi i kablovi

MIDI port ili sučelje je mjesto preko kojeg dva inače nekompatibilna uređaja ostvaruju komunikaciju dok su kablovi fizički prijenosnici voltaže u vrijednostima od 0 ili 1 [3, *MIDI Ports*]. Specifikacija opisuje tri tipa MIDI portova, tipove kablova njihovo kori-

štenje kao što je ilustrirano na slici [3, *MIDI Ports*]. MIDI uređaji tradicionalno koriste 5-pinski DIN konektore, iako novije verzije također koriste USB ili bežične veze. Na slici 2.1. možemo vidjeti tri vrste 5-pinskih DIN konektora na pozadini nekog MIDI hardvera.



**Slika 2.1.** Ilustracija MIDI uređaja sa sučeljima kablovima [3, *MIDI Ports*]

U ovom radu ne dotičemo se implementacije i korištenja hardver uređaja te njihovih sučelja, stoga ćemo sada prijeći na opis MIDI poruka.

### 2.1.2. MIDI poruke

MIDI poruke mehanizam su serijalizacije i tumačenja informacije prenesene s jednog uređaja na drugi preko komunikacijskih kanala protokola. Neke od MIDI poruke prikazane su i opisane u tablici 2.1. O njima će biti više riječi u kasnijim poglavljima i bitan su dio obrade MIDI datoteka u svrhu klasifikacije.

MIDI poruke	
Naziv poruke	Opis
<i>Note Off</i>	Nota je zaustavljena
<i>Note On</i>	Nota je počela
<i>Control Change</i>	Promijeni kontrolu parametra
<i>Program Change</i>	Promijeni instrument
<i>Pitch Bend</i>	Promijeni visinu tona
...	...

Tablica 2.1. Neke MIDI poruke i njihova značenja

Kanal za prijenos MIDI poruka podržava 16 neovisnih kanala po vezi, što omogućava istovremenu kontrolu nad različitim instrumentima te se podaci mogu slati brzinom prijenosa od 31.25 kbps. MIDI poruke kakve ćemo koristiti u ovom istraživanju sadržane su u MIDI datotekama.

## 2.2. MIDI datoteka

MIDI datoteke su digitalni formati koji pohranjuju informacije o glazbenim izvedbama u obliku MIDI poruka. One ne sadrže stvarni audio zapis, već skup uputa koje govore elektronskim instrumentima kako reproducirati glazbu. Standardni MIDI datotečni format koristi ekstenziju .mid ili .midi. MIDI datoteke temeljni su alat u glazbenoj industriji i primjena im je široka. Od komponiranja glazbe u prilagođenim softverima gdje se note i drugi parametri mogu jednostavno mijenjati, premještati i kopirati gdje se svaki aspekt izvedbe uključujući visinu tona, trajanje note i dinamiku može precizno prilagoditi, do orkestracije i notacije glazbe gdje se MIDI datoteke koriste kao izvor u koji se prepisuje glazbeni zvuk u MIDI poruke. U tu svrhu nastao je niz softvera za obradu MIDI datoteka i upravo bi korisnici tih softvera imali značajne koristi od klasifikacije MIDI datoteka.

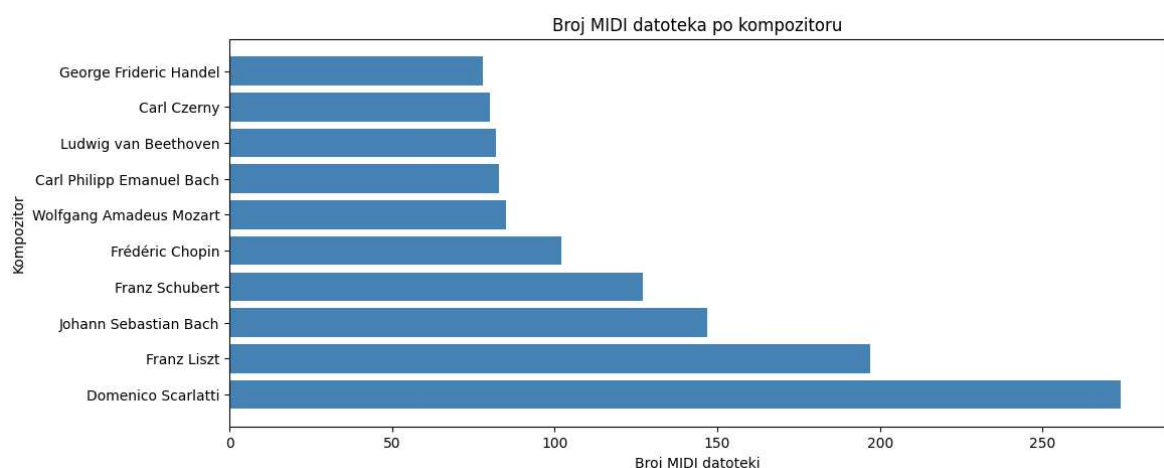
Standardna MIDI datoteka sastoji se od tri glavna dijela: *Header Chunk*-a koji sadrži osnovne informacije o datoteci, poput formata datoteke, broja traka i razlučivosti *tick*-a (najmanja jedinica mjerenja vremena u MIDI datoteci), *Track Chunk*-a gdje svaka traka (*track*) sadrži niz događaja (*events*) koji definiraju glazbenu izvedbu te *Events*-a koji specificiraju glazbene upute. Ukratko, možemo reći da svaka MIDI datoteka sadrži skup MIDI poruka i upravo će se iz informacija o glazbi iz njih dobiti korisne informacije za njihovu klasifikaciju.

### 3. Skup podataka GiantMIDI-Piano

Skup podataka *GiantMIDI-Piano* sastoji se od približno 39 milijuna noti, prepisanih iz audio u MIDI datoteke [1, *Abstract*]. Sadrži oko 11 tisuća djela napisanih za klavir svako pridruženo nekom od 2786 skladatelja [1, *Abstract*]. 90% MIDI datoteka u skupu stvoreno je prijepisom audio zapisa u MIDI, a ostalih 10% je upisano kao MIDI sekvenca [1, *Abstract*]. Ovaj skup podataka ćemo koristiti za učenje klasifikatora.

#### 3.1. Zastupljenost podataka

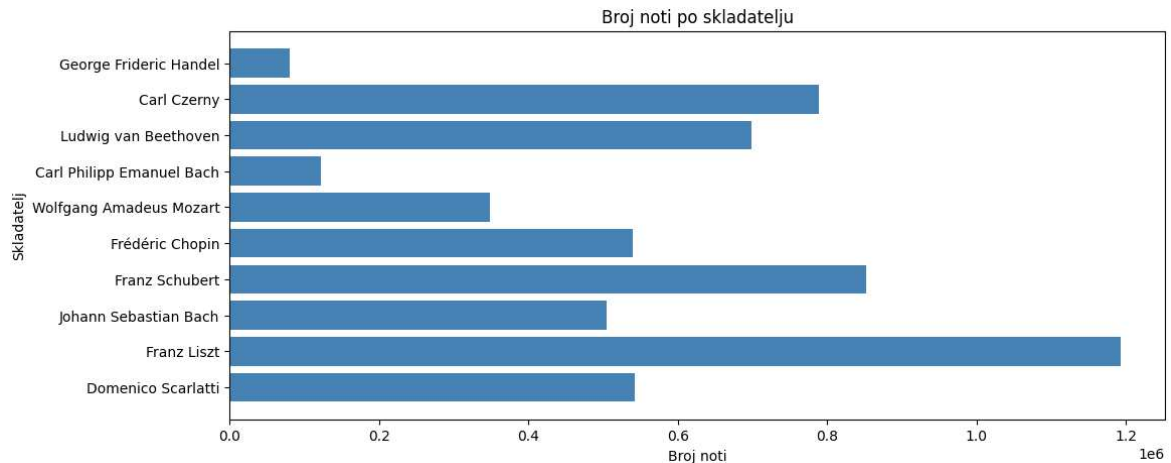
Uvjet klasifikacije bit će 10 skladatelja klasičnih djela. Za kratki pokaz klasi odabrano je prvih 10 najzastupljenijih skladatelja u skupu podataka. Na slici, 3.1., prikazan je broj djela prvih 10 skladatelja. Ovdje se primjećuje nebalansiranost primjera klasi ovakvim odabirom. Pri radu s ovim tipom podataka savršenu balansiranost ne možemo očekivati te se nećemo previše koncentrirati na pretjerano uređivanje zasebnih klasi.



**Slika 3.1.** Broj MIDI datoteki po 10 najzastupljenijih skladatelja u skupu podataka

Pošto skladbe nisu jednake duljine bitno je vidjeti koliko je informacije, broj noti,

sadržano u svim MIDI datotekama pojedinog skladatelja što je prikazano na slici 3.2. Promatranjem ova dva grafa (3.1. i 3.2.) možemo zaključiti da količina informacije u klasama ne ovisi samo o količini MIDI datoteka već i broju noti. Prilikom odabira skladatelja bitna su nam dakle dva saznanja: broj zastupljenih MIDI datoteka i broj noti.



Slika 3.2. Broj noti zastupljenih u skupu podataka svakog skladatelja

## 3.2. Odabir skladatelja

Nakon navedenih spoznaji o podacima možemo odabrati skladatelje koje ćemo klasificirati. Dakle, u obzir uzimamo broj MIDI datoteka i broj noti zastupljenih u podacima svakog skladatelja. Odabir je učinjen sljedećim kodom.

```

1 # Obrisati pri najzastupljeniji zapis
2 del(composers_by_notes[0])
3 del(composer_by_midi[0])
4
5 # Napravi presjek
6 targets = list(set(composers_by_notes) & set(composer_by_midi))
7
8 # Odaberi prvih 10 skladatelja
9 targets = targets[:10]

```

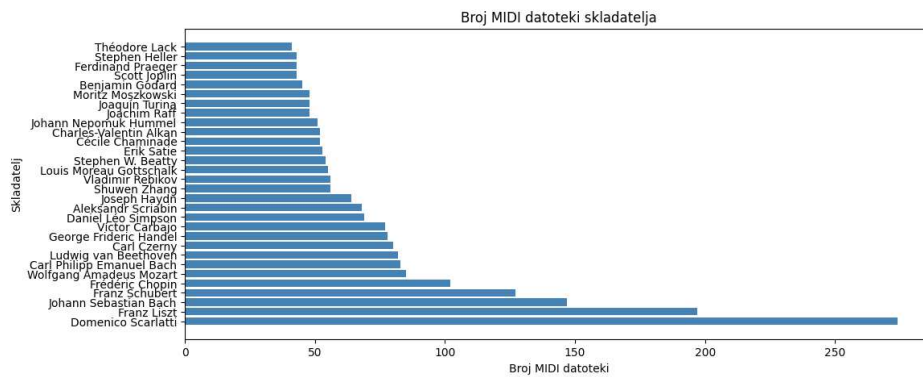
Liste *composers\_by\_midi* i *composers\_by\_notes* sadrže 30 prvih skladatelja sortirano po broju MIDI datoteka, tj. noti. Iz njih je uklonjen prvi skladatelj kako bismo postigli blagu balansiranost. Na slikama 3.3. i 3.4. vidljivo je kako ti zapisi strše. Zatim, u listu *targets*

učinjen je presjek i odabrano je prvih 10 skladatelja koji čine odabrane skladatelje za klasifikaciju.

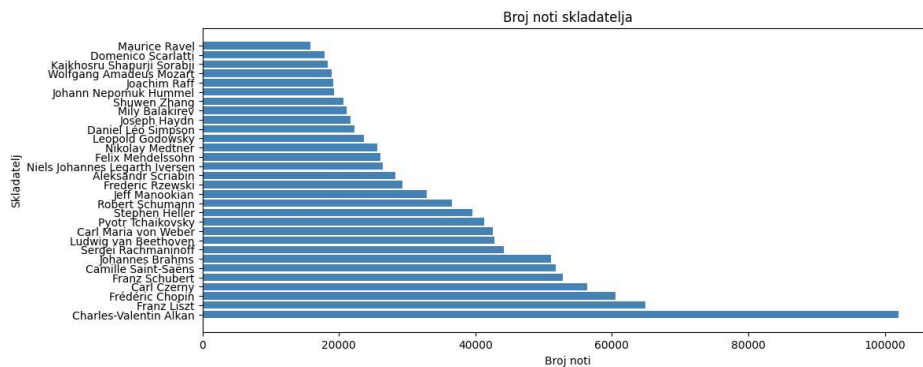
```

1 targets = {
2     'Franz Liszt': '0',
3     'Wolfgang Amadeus Mozart': '1',
4     'Stephen Heller': '2',
5     'Aleksandr Scriabin': '3',
6     'Daniel Leo Simpson': '4',
7     'Ludwig van Beethoven': '5',
8     'Johann Nepomuk Hummel': '6',
9     'Carl Czerny': '7',
10    'Shuwen Zhang': '8',
11    'Joseph Haydn': '9'
12 }

```



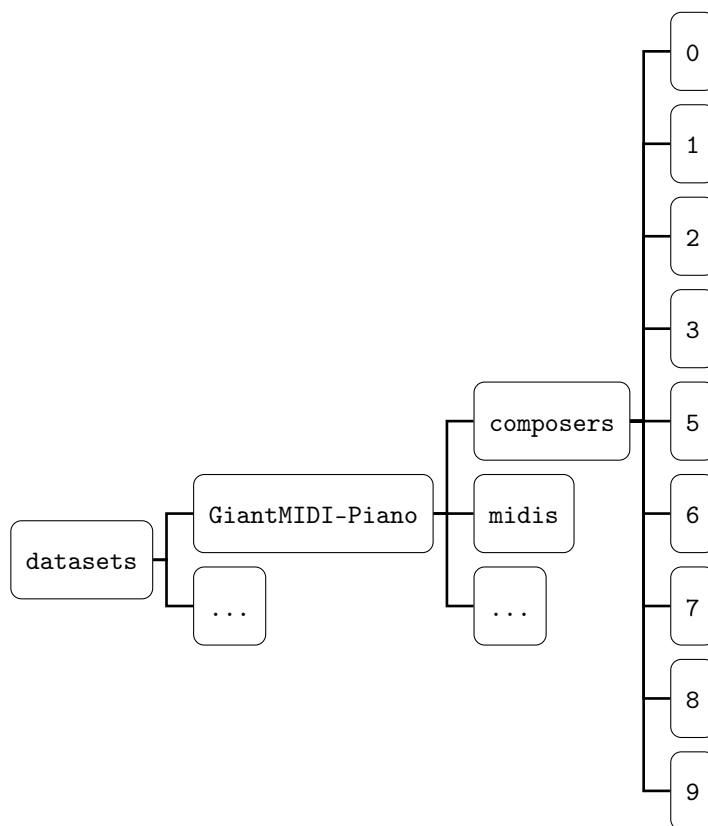
Slika 3.3. 30 najzastupljenijih skladatelja po broju MIDI datoteki



Slika 3.4. 30 najzastupljenijih skladatelja po broju noti

### 3.3. Podjela podataka

Prije no što počnemo s klasifikacijom preoblikovat ćemo skup podataka u novu datotečnu strukturu kako bismo iz imena lakše izlučili klasu. Unutar mape skupa podataka *GiantMIDI-Piano* stvorena je nova mapa *composers*, a u njoj 10 novih mapi nazvane rednim brojem klase skladatelja. U svaku od 10 tako kreiranih mapi iskopirane su MIDI datoteke pojedinog skladatelja. Ovakva podjela skladatelja u datoteke bit će nam od koristi u poglavlju *Klasifikacija metodama obrade prirodnog teksta 5*. Datotečna struktura prikazana je u nastavku.





## 4. Stroj potpornih vektora i klasifikator nasumičnih šuma

Klasifikacija podataka algoritmima stroja potpornih vektora, SVM (engl. *Support vector machine*) i klasifikatorom nasumičnih šuma, RF (engl. *Random forest classifier*) zahtijeva podatke s jednakim brojem značajki što znači da zapisi u skupu podataka moraju imati jednako dugu vektorsku reprezentaciju pridruženu pojedinoj klasi. U slučaju MIDI podataka to nije slučaj. Čak i kada bi se datoteke srezale na manje dijelove opet nismo sigurni što su naše značajke. Možemo pretpostaviti da je značajka redni broj note, a što ako je u pitanju akord, tj. više noti je pritisnuto u isto vrijeme, u tom slučaju također griješimo. Isto možemo reći ako je datoteke rascjepkamo na više njih, tada redni broj značajke nikako ne odgovara rednom broju note ili glazbene jedinice te se gubi trag melodiji. Stoga prije no što pokušamo klasificirati naše podatke moramo smisliti mehanizam ekstrakcije značajki.

### 4.1. Ekstrakcija značajki

Problem različite duljine skladbi pokušat će se riješiti korištenjem algoritma analize glavnih komponenata, PCA (engl. *Principal component analysis*) koji će biti opisan u ovom poglavlju. No, prije nego što iskoristimo PCA trebamo nekako iz podataka izvući informacije o notama. U tome će nam pomoći *Python* biblioteka *pretty\_midi*.

#### 4.1.1. *Pretty\_midi*

*Pretty\_midi* [4] biblioteka je razvijena u programskom jeziku Python za olakšano baratanje MIDI protokolom. Također, omogućava korisniku da stvori tok MIDI poruka iz MIDI datoteke i filtrira samo potrebne poruke.

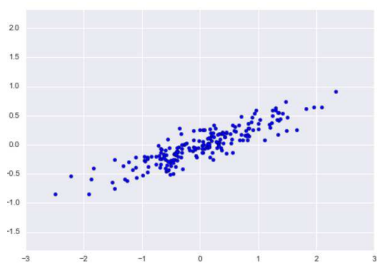
Svaka MIDI datoteka u sebi sadrži informaciju o notaciji skladbe u obliku noti. Dakle,

svaka nota ima svoju visinu tona (*pitch*), jačinu tona (*velocity*), početak (*note\_start*) i kraj (*note\_end*) trajanja note. Pošto su *note\_start* i *note\_end* značajke vremenske i rastuće ujedinit ćemo ih u novu značajku trajanja note *duration*. Dobivenom značajkom *duration* ćemo zamijeniti značajku *note\_end*.

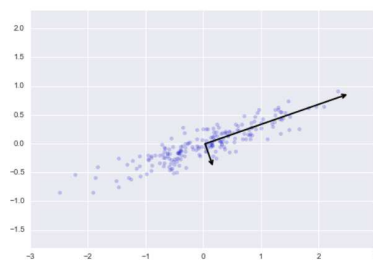
Nakon što znamo što tražimo, iz svake MIDI datoteke uz pomoć biblioteke *pretty\_midi* izvlačimo četiri glavne značajke iz svake note: *pitch*, *velocity*, *note\_start* i *duration*. Tako dobivene vrijednosti za svaku notu spremaju se u listu. Dakle, za skup podataka MIDI datoteka veličine ( $N$ ) sada smo dobili skup podataka dimenzija ( $N, n_i, 4$ ), gdje je  $n_i$  proizvoljan za svaku MIDI datoteku. Ovdje uočavamo problem varirajuće dimenzije  $n_i$  svakog MIDI zapisa. Upravo taj problem ćemo riješiti algoritmom PCA.

### 4.1.2. Analiza glavnih komponenta

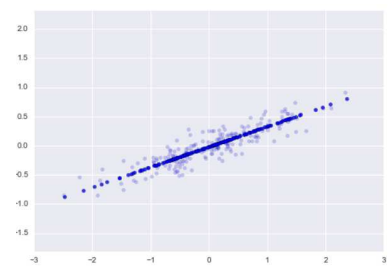
Analiza glavnih komponenta ili PCA, engl. *Principal component analysis* jednostavna je i brza nenadzirana metoda strojnog učenja redukcije dimenzionalnosti podataka [5, *Introducing Principal Component Analysis*]. Ovdje je bitno napomenuti kako se radi o metodi koja smanjuje dimenzionalnost podataka stoga imamo ograničenje minimalnom duljinom skladbe u skupu podataka. Algoritam ovdje neće biti pokazan već će ga se koristiti kao *crnu kutiju*. Na slikama 4.1., 4.2. i 4.3. slikovito je prikazano kako algoritam u suštini djeluje. U podacima 4.1. prepoznaju se glavne komponente 4.2. te se ona manje važnije svodi na nulu 4.3. i time se reduciraju podatci.



**Slika 4.1.** Skup podataka prije primjene PCA metode [5, *Introducing Principal Component Analysis*]



**Slika 4.2.** Glavne komponente [5, *Introducing Principal Component Analysis*]



**Slika 4.3.** Skup podataka nakon obrade PCA metodom [5, *Introducing Principal Component Analysis*]

Dakle, kako bismo uspješno saželi podatke dobivene iz MIDI datoteka potrebno je učiniti više koraka PCA algoritma u više dimenzija. PCA je izvršen po uzoru na odlomak

iz knjige [5, *Introducing Principal Component Analysis*] i pseudokod za sažimanje MIDI informacije pokazan je u nastavku:

```
1 reduce_features(notes: list, num_features: int) -> list:
2     pca_1 = PCA(n_components=1)
3     features = pca_1.fit_transform(notes)
4
5     chunks = floor(len(features) / num_features)
6     features = features[:int(chunks * num_features)]
7     features = features.reshape(num_features, chunks)
8
9     pca_2 = PCA(n_components=1)
10    features = pca_2.fit_transform(features)
11
12    return features
```

Početno polje *notes* dimenzije je  $dim(n, 4)$ , sadrži  $n$  note, a svaka nota sadrži 4 vrijednosti. Prvo korištenje PCA algoritma reducira 4 vrijednosti svake note u jednu komponentu  $n\_components$  i dobivamo novo polje *features* dimenziji  $dim(n, 1)$ . Nakon tog koraka grupiramo vrijednosti polje **features** u skupine, *chunks* te tako dobivamo broj skupini ekvivalentan broju željenih značajki *num\_features*. U posljednjem koraku još jednim korakom PCA algoritma reduciramo skupine na jednu vrijednost te dobivamo polje *features* veličine *num\_features*.

Ovakvim korištenjem PCA algoritma reducirali smo dimenzije MIDI podatka na sljedeći način:

$$dim(n_i, 4) \rightarrow dim(n_i, 1) \rightarrow dim(n, 1),$$

gdje je  $n_i$  količina noti u  $i$ -tom MIDI zapisu, a  $n$  količina značajki nakon provođenja algoritma.

Bitno je napomenuti da ovakvim pristupom ne dobivamo smislene informacije o okolini neke note već pokušavamo naprosto pogoditi vrijednosti koje će biti svojstvene okolini određene note. Također, pošto je glazba slijed noti, ne možemo nikako znati dali je veličina *chunk*-a smisljena za određenu notu, tj. jesmo li veličinom *chunk*-a pogodili

neku melodiju ili akord. S ovim činjenicama možemo pokušati naučiti naše klasifikatore stroja potpornih vektora i slučajnih šuma.

## 4.2. Učenje klasifikatora

Nakon podjele skupa podataka u klase u datotečnu strukturu kao što je prikazao na grafu u potpoglavlju *Podjela podataka* 3.3. jednostavno možemo učitati puteve do datoteka u jednu listu. Kada smo tako podijelili podatke sada ih možemo razvrstati u skup za učenje i skup za treniranje u omjeru 80/20. Također, oznaku MIDI datoteke sada lako možemo izlučiti iz puta do datoteke uz pomoć neke arbitrarne funkcije. To se može učiniti na sljedeći način:

```
1 midi_paths = list(Path(dataset_path, "by_composer"))
2 midi_paths = midi_path.glob("**/*.mid")
3 random.shuffle(midi_paths)
4
5 train_files = midi_paths[:int(len(midi_paths) * 0.8)]
6 test_files = midi_paths[
7     int(len(midi_paths) * 0.8):int(len(midi_paths) * 1.0)
8 ]
```

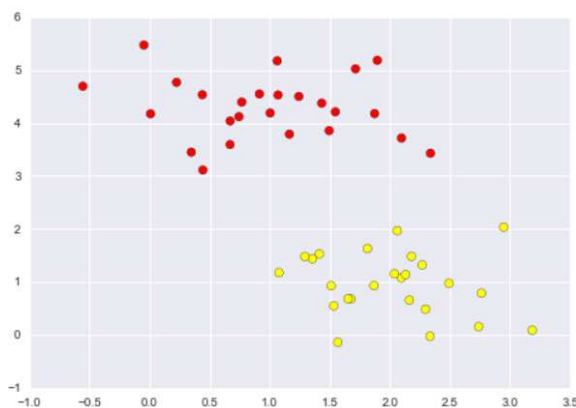
Prije no što na učitanim podacima primijenimo PCA potreban nam je još jedan parametar, a to je minimalan broj noti u MIDI datoteci. Tako ćemo, pošto činimo redukciju podataka, neke datoteke jednostavno zanemariti.

```
1 number_of_features = n
2 min_notes = number_of_features
3
4 train_midis, train_labels = load_files(train_files, min_notes)
5 test_midis, test_labels = load_files(test_files, min_notes)
6
7 X_train = extract_features_batch(train_midis, number_of_features
8     )
9 X_test = extract_features_batch(test_midis, number_of_features)
```

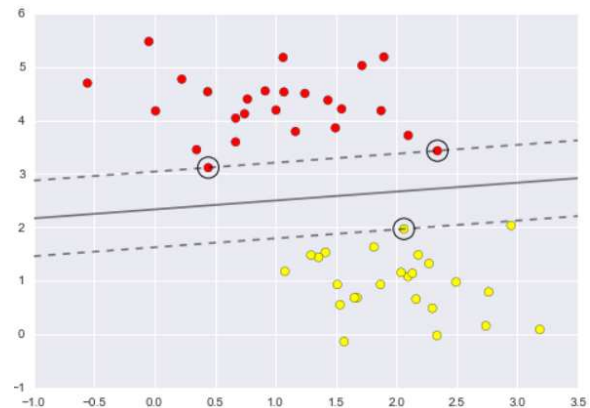
U odsječku koda funkcija *load\_files* uz pomoć biblioteke *pretty\_midi* učitava podatke kao što je predstavljeno u poglavlju 4.1.1. te iz putanje do datoteke izluči oznaku klase dok funkcija *extract\_features* vrši PCA algoritam predstavljen u poglavlju 4.1.2. U klasifikaciji će se probati više vrijednosti varijable *num\_features*: 50, 75, 150, 250 i 500 te za svaku vrijednost naučit oba klasifikatora, SVM i RF.

### 4.2.1. Stroj potpornih vektora

Stroj potpornih vektora, skraćeno SVM, engl. *Support Vector Machine*, iznimno je moćan i fleksibilan algoritam nadziranog učenja korišten za klasifikaciju ili regresiju [5, *In-Depth: Support Vector Machines*]. Algoritam u suštini funkcionira tako da između dvije odvojene klase podataka nastoji povući pravac koji maksimizira minimalnu marginu između najbližih podataka u pojedinim klasama.



**Slika 4.4.** Skup podatak koji sadrži dvije klase [5, *In-Depth: Support Vector Machines*]



**Slika 4.5.** Podatci u skupu odvojeni marginom [5, *In-Depth: Support Vector Machines*]

Na slikama 4.4. i 4.5. pokazan je skup podataka odvojen pravcem koji maksimizira minimalnu marginu između najbližih podataka dvije klase. Vektori koji se nalaze točno na margini, zaokruženi na slici, nazivaju se potporni vektori. U slučaju naše klasifikacije MIDI datoteka postoji 10 klasi, no taj problem je rješiv algoritmima OvO (*One-vs-One*) ili OvA (*One-vs-All*) koji su već implementirani u sklopu SVM algoritma iz biblioteke *skit-learn* [6].

SVM algoritam moguće je konfigurirati parametrima. Parametre kojima će SVM klasifikator biti konfiguriran su: *kernel*, *C* i *gamma*. Za svaku konfiguraciju parametara u parametarskom *grid*-u naučit će se model za pojedinu vrijednost varijable *num\_features*

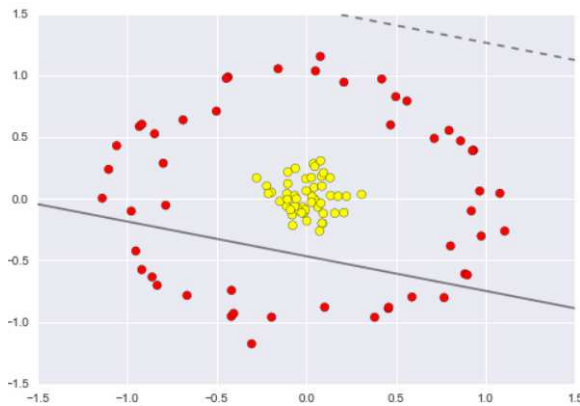
predstavljene u poglavlju 4.2. Značenje pojedinog parametar i na što njihova promjena utječe bit će predstavljeno u nastavku. Također, matematika SVM-ova neće biti predstavljena, već samo kako ih koristiti u sklopu problema klasifikacije korištenjem već implementiranog algoritma iz biblioteke *skit-learn* [6].

```

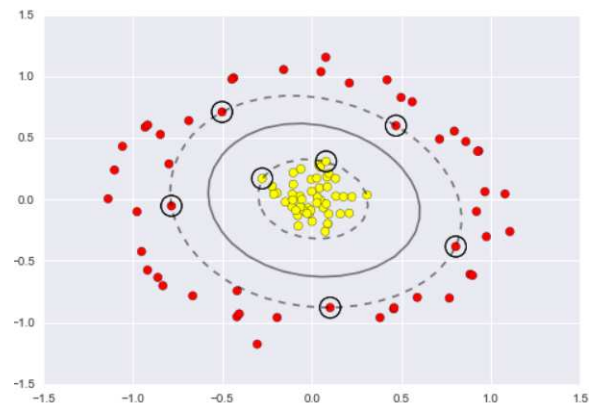
1 param_grid = {
2     'kernel': [rbf]
3     'C': [1, 5, 10, 50],
4     'gamma': [0.0001, 0.0005, 0.001, 0.005]
5 }

```

*Kernel* (ili jezgra) parameter omogućava nam projekciju podataka u višedimenzionalni prostor definiran polinomima ili Gausovim baznim funkcijama i time definira nelinearnu granicu između podataka [5, *In-Depth: Support Vector Machines, Beyond linear boundaries: Kernel SVM*]. Na slikama 4.6. i 4.7. slikovito je prikazano što *rbf* jezgra čini. Upravo tu jezgrenu funkciju ćemo koristiti u konfiguraciji klasifikatora.



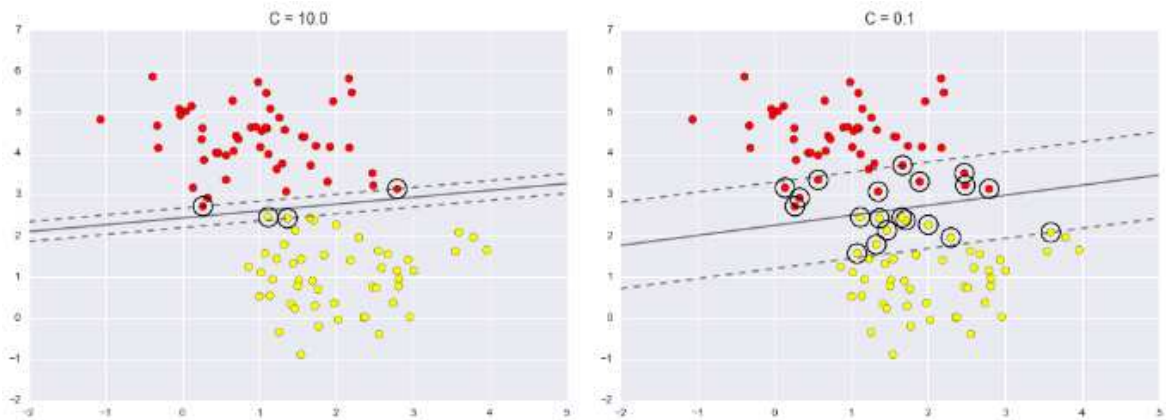
**Slika 4.6.** Skup podatak koji ima nelinearnu granicu između klasi [5, *In-Depth: Support Vector Machines*]



**Slika 4.7.** Nelinearnost modelirana *rbf* jezgrom [5, *In-Depth: Support Vector Machines*]

C parametar kontrolira takozvanu mekoću margine. U stvarnom svijetu, našem primjeru klasifikacije, granica između klasi nije jasna linija, tj. podaci jedne klase ulaze u prostor druge pošto su često marginalni primjeri sličnih vrijednosti. Iz navedenog razloga nekim podacima moramo pustiti da uđu unutar margine, tj. u prostor između potpornih vektora naših klasa. Taj pristup naziva se omekšavanje margine i kontroliran je parametrom C. Parametar C poprima vrijednosti između 0 i  $+\infty$  i što je parametar manji

marginu je mekša, tj. više primjera puštamo u marginu, a što je parametar veći marginu je tvrđa i manje primjera se nalazi u njoj. Može se reći da ovaj parametar regularizira naš model, tj. smanjenjem parametra  $C$  model postaje manje sklon prenaučeni skupu podataka za učenje. Na slici 4.8. prikazan je utjecaj veličine parametra  $C$  na marginu.



**Slika 4.8.** Utjecaj parametra  $C$  na mekoću margine [5, *In-Depth: Support Vector Machines*]

Parametar  $\gamma$  koristan je samo ako je za jezgru SVM modela odabrana rbf jezgra. Taj parametar kontrolira zakrivljenost linije naše jezgre [5, *In-Depth: Support Vector Machines*]. Što je  $\gamma$  veći krivulja je zakrivljenija i prilagođenija na podacima za učenje. Stoga smanjenjem ovog parametra također utječemo na regularizaciju modela.

Predstavljenim parametrima SVM klasifikatora sada možemo naučiti model na sljedeći način:

```

1 from sklearn.svm import SVC
2 from sklearn.model_selection import GridSearchCV
3
4 estimator = SVC(kernel='rbf', class_weight='balanced')
5 grid = {
6     'C': [10e-3, 10e-2, 10e-1, 0.1, 0.5, 1, 5, 10, 50, 100, 500,
7         1000],
8     'gamma': [10e-3, 10e-2, 10e-1, 0.1, 0.5, 1, 5, 10, 50, 100,
9         500, 1000]
10 }
11 grid_search = GridSearchCV(estimator, grid)

```

```
11 grid_search.fit(X_train, y_train)
12 y_pred = grid_search.predict(X_test)
```

Varijabla *estimator* stvara objekt klase SVC koji u sebi sadrži algoritme za računanje potpunih vektora. Klasa *GridSearchCV* namjesti estimator i *grid* koje smo zadali, a metodom *fit* *grid*-om i skupovima za treniranje *X\_train* i *y\_train* optimira se *estimator*. Metoda *predict* predviđa oznake za neoznačeni skup *X\_test* i dobivamo oznake *y\_pred*.

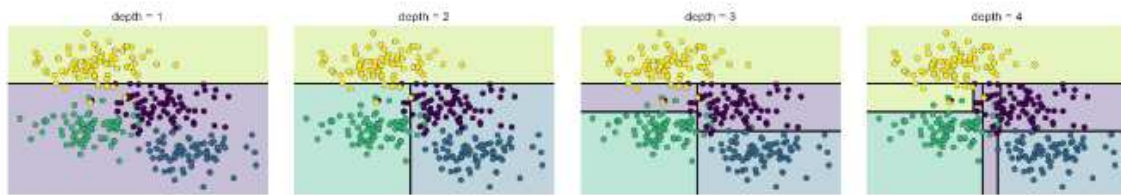
Rezultati u obliku točnosti, engl. *accuracy* i F1-vrijednosti prikazani su u potpoglavlju 4.3.

## 4.2.2. Klasifikator slučajnih šuma

Klasifikator slučajnih šumi ili skraćeno RF, engl. Random forest, za razliku od SVM-a neparametarski je klasifikator [5, *In-Depth: Decision Trees and Random Forests*]. Također, RF klasifikator je ansambl metoda što znači da se oslanja na skupljene podatke skupa jednostavnijih klasifikatora [5, *In-Depth: Decision Trees and Random Forests*]. Logika iza ansambla je ta da je *zbroj veći od pojedinačnih dijelova*, dakle glasanje među više klasifikatora donosi bolju odluku od svakog klasifikatora pojedinačno [5, *In-Depth: Decision Trees and Random Forests*]. RF primjer je ansambla skupa klasifikatora zvanih stablima odluke.

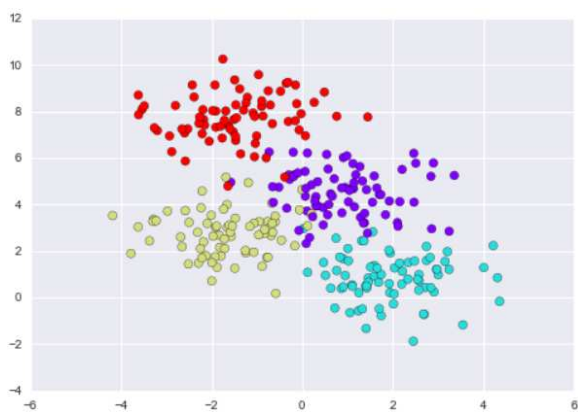
Stabla odluke veoma su intuitivan algoritam klasifikacije podataka [5, *In-Depth: Decision Trees and Random Forests*]. Možemo ih zamisliti kao skup pitanja kojima dolazimo do zaključku o objektu klasifikacije. Takva binarna podjela odluke ako su pitanja dobro koncipirana čini jako efikasnu klasifikaciju. U strojnom učenju pitanja generalno poprimaju oblik podjele podataka po nekoj osi u prostoru parametara. Dakle za svaku dubinu stabla podaci su podijeljeni nekom hiperravninom u prostoru podataka i tako su stvorene zone u kojima pripada pojedinačna klasa kao što je prikazano na slici 4.9.



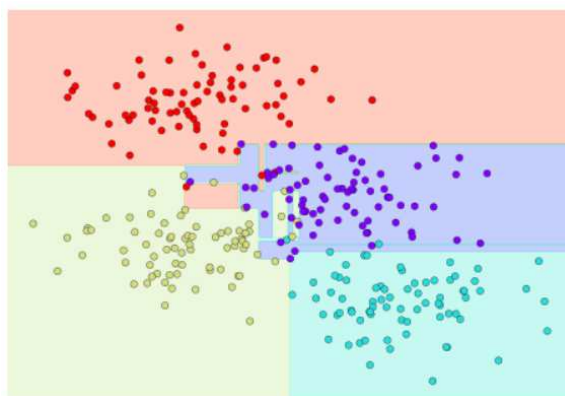


**Slika 4.9.** Promjena granica klasifikacije stablom odluke dubinom stabla *depth* [5, *In-Depth: Decision Trees and Random Forests*]

Na slikama 4.10. i 4.11. kao što je to pokazano u knjizi [5, *In-Depth: Decision Trees and Random Forests*] vidi se kako stablo odluke razvrstava klase u svoja područja unutar prostora značajki. Također, vidljivo je kako ta podjela nije savršena pošto su se neki primjeri zavukli i uzrokovali svoja područja unutar drugih klasi i tu možemo reći da se klasifikator prenaučio. Takva prenaučenosť jedna je od značajki stabla odluke pošto je jako lako ući preduboko u stablo i izazvati rupe kao na slici 4.11. [5, *In-Depth: Decision Trees and Random Forests*]. Taj problem prenaučenosťi rješava se korištenjem algoritma nasumičnih šumi [5, *In-Depth: Decision Trees and Random Forests*].



**Slika 4.10.** Nepodijeljen skup podataka [5, *In-Depth: Decision Trees and Random Forests*]



**Slika 4.11.** Skup podijeljen stablom odluke [5, *In-Depth: Decision Trees and Random Forests*]

Ideja da više prenaučeni klasifikatora može biti iskorišteno kao jedan kako bi se smanjila prenaučenosť i povećala generalizacija naziva se *bagging*. U praksi stabla odluke najbolje funkcioniraju kada su nasumična, unosenjem stohastičnosti u podjelu podataka hiperravninom. [5, *In-Depth: Decision Trees and Random Forests*] U biblioteci *skit-learn* [6] stohastičnost je implementirana klasom *RandomForestClassifier* koji takvu nasumičnost automatizira; sve što je potrebno učiniti jest dati broj klasifikatora varijablom *n\_estimators* [5, *In-Depth: Decision Trees and Random Forests*]. U nastavku je is-

ječkom koda prikazano kako iskoristiti algoritam nasumičnih šuma koristeći biblioteku *skit-learn*. Također, kao i u prijašnjem poglavlju 4.2.1. namješten je grid i prikazan je cijeli postupak učenja klasifikatora.

```
1 from sklearn.ensemble import RandomForestClassifier
2 from sklearn.model_selection import GridSearchCV
3
4 estimator = RandomForestClassifier()
5 grid = {
6     'n_estimators': [50, 100, 250, 500, 750, 1000]
7 }
8
9 grid_search = GridSearchCV(estimator, grid)
10 grid_search.fit(X_train, y_train)
11 y_pred = grid_search.predict(X_test)
```

Rezultati ovog algoritma u obliku točnosti, engl. *accuracy* i F1-vrijednosti također su prikazani sljedećem potpoglavlju 4.3.

### 4.3. Rezultati

Rezultati su prikazani u tablicama *Accuracy* 4.1. i *F1-macro* vrijednostima 4.2. te matricama zabune u potpoglavlju *Matrice zabune* 4.3.1. Prije nego što pokažemo rezultate objasniti ćemo ove pojmove i što oni znače.

*Accuracy* (točnost) je udio ispravno klasificiranih primjera među svim uzorcima. Daje nam uvid u to koliko je puta model pogodio točnu klasu. *Accuracy* nas može lako zavarati jer model može imati veliku točnost predviđajući samo najzastupljeniju klasu, a to je često tako pri treniranju modela nebalansiranim skupom podataka. Računa se na idući način:

$$Accuracy = \frac{n_{ispravno\_klasificirani\_uzorci}}{n_{upupan\_broj\_uzoraka}}$$

*F1-macro* vrijednost je srednja vrijednost *F1-score* vrijednosti za svaku klasu. *F1-score*

za pojedinu klasu računa se kao harmonijska sredina između preciznosti (*Precision*) i osjetljivosti (*Recall*) te klase. *F1-macro* je bolji u situacijama gdje postoji neuravnoteženost među klasama, jer uzima u obzir preciznost i osjetljivost svake pojedine klase. Računa se na sljedeći način:

$$F1\_score_i = 2 \cdot \frac{Precision_i \cdot Recall_i}{Precision_i + Recall_i}$$

$$F1\_macro = \frac{1}{n} \cdot \sum_i^n F1\_score_i.$$

*Precision* (preciznost) klase *i* odnosi se na udio točno pozitivnih predikcija klase *i* (*TP*) u odnosu na sve predikcije koje je model klasificirao točno (*TP + FP*), gdje su *FP* svi negativni primjeri koje je model klasificirao točno. Preciznost se računa na sljedeći način:

$$Precision = \frac{TP}{TP + FP}.$$

*Recall* (osjetljivost) odnosi se na udio ispravno klasificiranih pozitivnih uzoraka u odnosu na sve ispravno klasificirane uzorke (*TP + TN*), gdje su *TN* svi negativni uzorci koje je model klasificirao negativno. Računa se ovako:

$$Recall = \frac{TP}{TP + TN}.$$

Iz tablica *Accuracy* 4.1. i *F1-macro* možemo očitati najviši *accuracy* za modele s 250 parametara i iznose 0.42 i 0.40. Također najviše *F1-macro* vrijednosti očitavaju se za iste modele te iznose 0.31 i 0.33. Ovakva razlika u ovim vrijednostima nalaže nam da skup podataka nije dobro balansiran i vjerojatno se model prilagodio jednoj klasi. Kako bismo vjernije prikazali rezultate treniranih modela, također smo rezultate pokazali i u obliku matrica zabune gdje se nebalansiranost podataka može lakše vidjeti.

<i>Accuracy</i>					
Param.	50	75	150	250	500
<i>SVM</i>	0.35	0.31	0.38	0.42	0.37
<i>RF</i>	0.35	0.33	0.37	0.40	0.31

Tablica 4.1. *Accuracy* ovisno o broju parametara i klasifikatoru

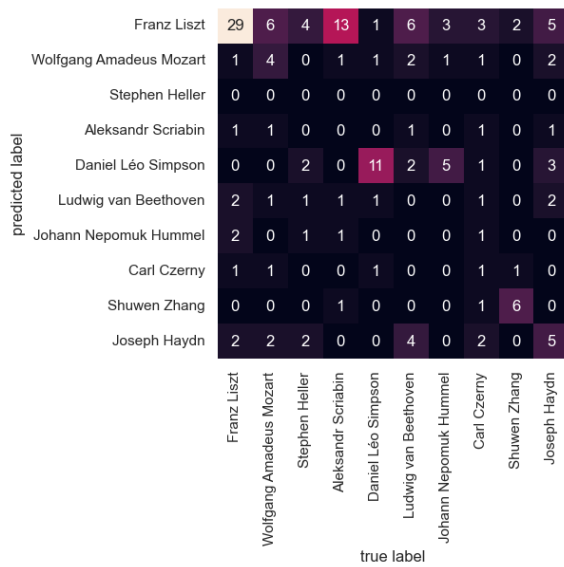
<i>F1-macro</i>					
Param.	50	75	150	250	500
<i>SVM</i>	0.25	0.22	0.30	0.31	0.25
<i>RF</i>	0.22	0.24	0.31	0.33	0.23

Tablica 4.2. *F1-macro* ovisno o broju parametara i klasifikatoru

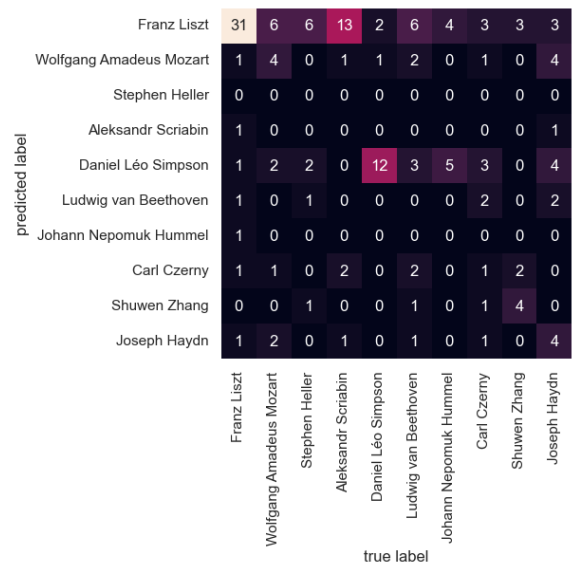
### 4.3.1. Matrice zabune

Matrica zabune (*confusion matrix*) je matrica iz koje možemo ocijeniti performanse modela. Ona prikazuje stvarne i predviđene klase i time omogućava uvid u to koliko je primjera završilo u točnoj, tj. krivoj klasi. U našim matricama zabune koje slijede točne klase nalaze se na stupcima, a njihove predikcije u redcima. Broj koji se nalazi na dijagonali matrice označava broj točno klasificiranih primjera pojedine klase. Što je više primjera zastupljeno na dijagonali matrice zabune to model bolje generalizira.

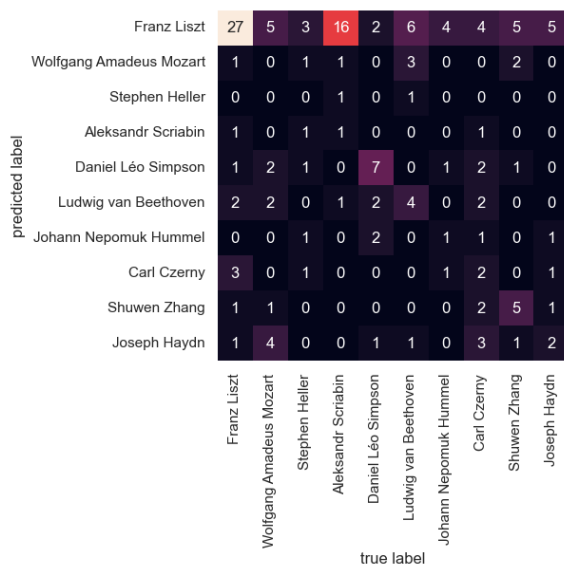
U sljedećim matricama zabune uočavamo kako je najviše primjera krivo klasificirano u klase *Franz Liszt* i *Daniel Leo Simpson*. To se može uočiti u skoro svim matricama zabune. U tablicama *Accuracy* i *F1-macro* modela SVM i RF s 250 parametara ne uočavamo znatno različite vrijednosti, no u njihovim matricama zabune možemo vidjeti kako ne klasificiraju primjere jednoliko. Dok je SVM većinu primjera smjestio u klasu *Franz Liszt* i jedino dobro prepoznao klasu *Daniel Leo Simpson* model RF donekle je popunio dijagonalu i tu možemo reći da RF s 250 parametara bolje generalizira od SVM-a s 250 parametara.



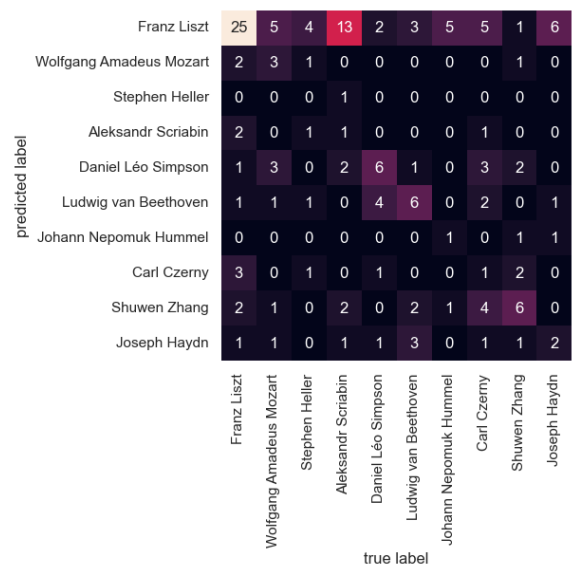
**Slika 4.12.** Matrica zabune SVM modela s 50 parametara



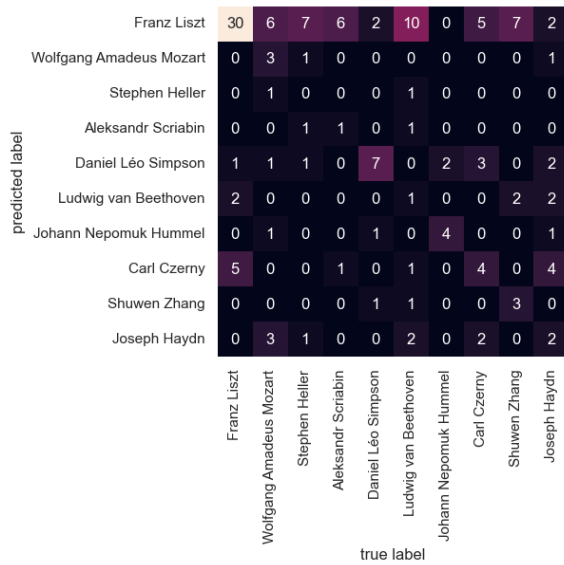
**Slika 4.13.** Matrica zabune RF modela s 50 parametara



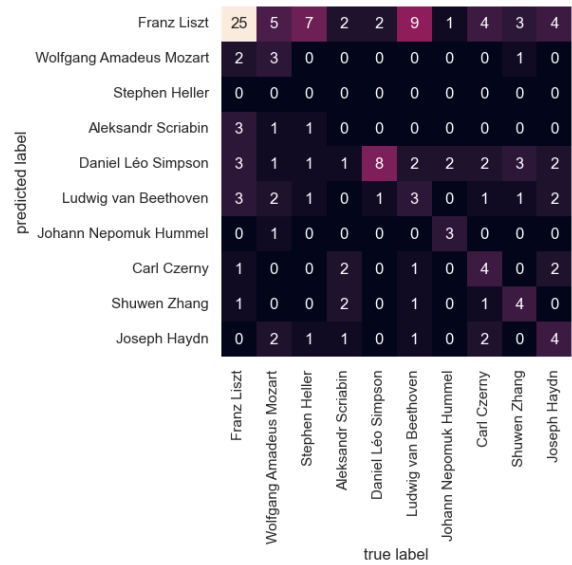
**Slika 4.14.** Matrica zabune SVM modela s 75 parametara



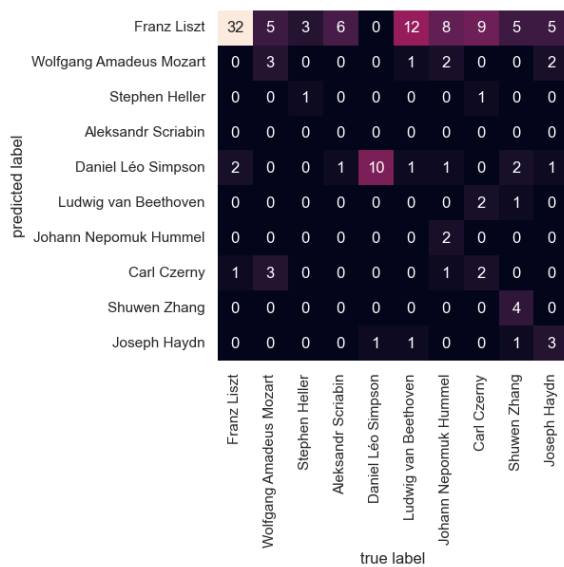
**Slika 4.15.** Matrica zabune RF modela s 75 parametara



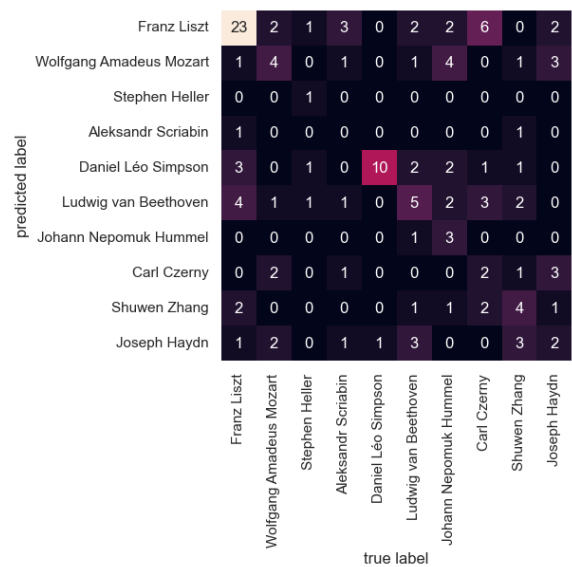
**Slika 4.16.** Matrica zabune SVM modela s 150 parametara



**Slika 4.17.** Matrica zabune RF modela s 150 parametara



**Slika 4.18.** Matrica zabune SVM modela s 250 parametara



**Slika 4.19.** Matrica zabune RF modela s 250 parametara



## 5. Klasifikacija metodama obrade prirodnog teksta

MIDI datoteke skup su instrukcija koje se kao što smo objasnili u poglavlju 2. mogu svesti na skup poruka. Među tim porukama lako se filtriraju one svojstvene za glazbeni zapis samog djela, tablica 2.1. Također, ništa nas ne sprječava da te poruke zapišemo u tekstualnom zapisu u neku datoteku i tako dobijemo tekst nastao iz MIDI poruka. Sada taj tekst možemo dovesti na ulaz nekog dubokog modela i naučiti ga. Dakle, ovim pristupom uklonili smo potrebu za dobivanjem i svođenjem MIDI datoteka na skup parametara, kao u poglavlju 4. te ih pretvorili u sekvence znakova i time sačuvali bitnu vremensku komponentu glazbe.

Pristup obradi MIDI datoteka prirodnim tekstom podijeljen je na dva glavna dijela: tokenizacija MIDI datoteka i model klasifikatora tokena. Tokenizacijom skupa podataka dobivamo *rječnik* tokena, engl. *vocab*, kojima razlučujemo smislene cjeline u našem tekstu MIDI poruka. Za MIDI datoteke algoritmi koji tokeniziraju MIDI poruke već postoje i jedan od njih je REMI. REMI tokenizator [7] poslužit će nam za tokenizaciju melodijskih, ritmičnih i harmonijskih dijelova MIDI poruka. Zatim će se u sklopu tokenizacije predstaviti i biblioteka *MidiTok* [8] napisan u programskom jeziku *Python*. *MidiTok* pruža nam pripremu MIDI datoteka za tokenizaciju te već enkapsulira upotrebu i učenje *REMI* klasifikatora algoritmom po izboru. Također, *MidiTok*-om podatci se mogu lako pripremiti i pretvoriti u sekvence za obradu.

Nakon učenja tokena iz skupa za učenje i pretvorbe podataka u sekvence, njima možemo naučiti duboke mreže. U ovom radu istražiti će se povratne neuronske mreže pri tom koristeći GRU, engl. *Gated recurrent unit* i LSTM, engl. *Long short-term memory* ćelije. Također, u njih će se i ugraditi mehanizam pozornosti, *attention*. Naučeni modeli bit će tablično uspoređeni i opisani kao i u poglavlju 4.3.



## 5.1. Tokenizacija

Tokenizacija skupa podataka je pronalazak smislenih leksičkih jedinica, tokena, svojstvenih za obradu podatka u kojima se pojavljuju i neophodan je dio obrade prirodnog teksta. Pošto MIDI datoteke možemo pretvoriti u tekstualni zapis te sadrže glazbene uzorke MIDI poruka, u njima također možemo prepoznati tokene. No, kako bismo uopće stvorili rječnik tokena, potrebno je prvo odrediti početne tokene na temelju kojih ćemo izgraditi cijeli rječnik. Odabir početnih tokena i učenje novih odradit će se u dva koraka korištenjem tokenizatora REMI i biblioteke koja enkapsulira REMI i daje nam na izbor više tokenizacijskih algoritama *MidiTok* kao što su BPE ili Unigram.

### 5.1.1. REMI

Predstavljen u radu *Pop Music Transformer (Huang and Yang)* [7] REMI, *REvamped-MIDI*, model je transformera za prepoznavanje melodijskih i ritmičkih strukturi u MIDI porukama. REMI nastoji iz previše zrnate uređenosti MIDI poruka prepoznati kompleksnije strukture i njima ih zamijeniti ne narušavajući time ritam i harmoniju zapisane glazbe. Neki učestali REMI tokeni prikazani su u tablici 5.1. Broj REMI tokena jedne MIDI datoteke ne odgovara njenom broju poruka.

REMI token	
Naziv tokena	Opis
<i>Bar</i>	Početak nove glazbene jednine
<i>Position_&lt;n&gt;</i>	Pozicija note
<i>Pitch_&lt;n&gt;</i>	Visina tona
<i>Velocity_&lt;n&gt;</i>	Dinamika tona
<i>Duration_&lt;n&gt;</i>	Trajanje note
<i>Rest_&lt;n&gt;</i>	Glazbena pauza
<i>Chord_&lt;type&gt;</i>	Tip akorda
...	...

Tablica 5.1. česti REMI tokeni i njihova značenja

Na primjer, neki niz REMI tokena može izgledati ovako:

```
[..., 'Rest_7.0.2', 'Rest_0.3.8', 'Position_27', 'Pitch_62', 'Velocity_79', 'Duration_0.2.8', 'Position_29', 'Pitch_64', 'Velocity_79', 'Duration_0.3.8', 'Position_31', 'Pitch_65', 'Velocity_79', 'Duration_0.3.8', ...]
```

## 5.1.2. MidiTok

REMI tokenizator neće se koristiti izravno, već u sklopu *Python* paketa *MidiTok* [8]. *MidiTok* tokenizira MIDI datoteke enkapsulirajući pritom najpoznatije MIDI tokenizatore među kojima je i REMI. Pruža i algoritme poput BPE, *Unigram* i *WordPiece* za učenje tokenizatora. U nastavku ćemo pokazati kako se paket koristi, odabrat ćemo REMI tokenizator i BPE algoritam učenja tokenizacije. Također, pokazat ćemo i kako se podatci mogu pripremiti za učenje modela.

Prije nego što pokažemo kako korištenjem *MidiTok* paketa možemo naučiti tokenizator objasniti ćemo BPE algoritam.

BPE, engl. *Byte-Pair Encoding* originalno je osmišljen za kompresiju teksta, a nakon toga korišten u mnogim današnjim modelima transformera kao što su *GPT*, *GPT-2*, *BART* itd. BPE tokenizacija podijeljena je u sljedećih 5 koraka:

1. Određivanje inicijalnih tokena: U ovom koraku tekst se razbije na najosnovnije jedinice, tokene, najčešće znakove tj. slova ili skupove slova te je svakoj jedinici dodijeljen je jedan bajt. U slučaju obrade MIDI datoteka one su podijeljene na REMI tokene, 5.1. te je svakom tokenu dodijeljen identifikacijski broj, *id*, to jest preko *id*-a jedan bajt.
2. Brojanje frekvenciji tokena: Nakon definiranja osnovnih tokena algoritam uparuje sve tokene u tekstu (u slučaju MIDI datoteka to je skup datoteka za učenje) te ih broji.
3. Spajanje najčešćeg para: Nakon brojanja parova najučestaliji par se spaja u novi token. U slučaju REMI tokenizatora to se učini na način da se bajtovi dva tokena spoje u novu riječ, a njoj se dodijeli novi *id*.
4. Ponavljanje koraka 2. i 3.: Sada s novom riječi u rječniku koraci 2. i 3. se ponavljaju dok nije postignuta željena veličina rječnika, tj. *vocab*-a.
5. Spremanje rječnika: Nakon što smo dobili zadovoljavajuću veličinu rječnika ili je algoritam prošao korak 4. bez promjene generiran je rječnik (*vocab*) kojim možemo tokenizirati novi tekst.

Učenje rječnika, skupa tokena, *vocab* korištenjem *MidiTok* paketa prikazano je u odsječku koda koji slijedi:

```
1 from miditok import REMI, TokenizerConfig
2
3 tokenizer_type = "BPE"
4 vocab_size = 10000
5
6 config_params = {
7     "pitch_range": (0, 127),
8     "num_velocities": 16,
9     "special_tokens": ["PAD", "BOS", "EOS"],
10    "use_chords": True,
11    "use_rests": True,
12    "use_tempos": False,
13    "use_time_signatures": False,
14    "use_programs": False
15 }
16
17 config = TokenizerConfig(**config_params)
18 tokenizer = REMI(config)
19 tokenizer.train(vocab_size=vocab_size, model=tokenizer_type,
    files_paths=train_files)
```

Iz paketa *MidiTok* prvo učitamo *REMI* i *TokenizerConfig*. U varijable *tokenizer\_type* i *vocab\_size* pohranimo ime tokenizatora, u našem slučaju BPE, i broj riječi u rječniku. *Config\_params* definira parametre kojima ćemo inicijalizirati objekt *TokenizerConfig* tj. njime ćemo zadati željene početne *REMI* tokene. *Pitch\_range* će definirati *Pitch* tokene, *num\_velocities* *velocity* tokene, itd. iz tablice 5.1. Također, ovdje vidimo i specijalne tokene *PAD*, *BOS* i *EOS*, oni nam pomažu pri učenju modela kako bi dobili zadovoljavajuće veličine ulaznog niza, no o njima ćemo pričati u poglavlju *Povratne neuronske mreže*, 5.2. Konfiguracija se zatim proslijedi u objekt *REMI* koji stvara tokenizator *tokenizer*. Tokenizator se metodom *train* trenira na skupu podataka za učenje *train\_files* korištenjem BPE algoritma te dobivamo rječnik (*vocab*) veličine *vocab\_size*.

U nastavku ćemo još pokazati kako paketom *MidiTok* korištenjem naučenog tokenizatora možemo pripremiti MIDI datoteke za učenje. Dakle, cilj je svaku MIDI datoteku rascjepkati u manje jednake dijelove kako bismo ih lakše doveli na ulaz modela koji učimo. Isječak koda koji slijedi pokazuje upravo to:

```
1 from miditok.pytorch_data import
2     DatasetMIDI,
3     split_files_for_training
4 from src.util.dataset import get_label_from_file_path
5
6 max_seq_len = 512
7 chunks_path = Path(dataset_path, "midi_chunks")
8
9 train_chunks = split_files_for_training(
10     train_files,
11     tokenizer,
12     Path(chunks_path, "train"),
13     max_seq_len)
14
15 train_dataset = DatasetMIDI(
16     files_paths=train_chunks,
17     tokenizer=tokenizer,
18     max_seq_len=max_seq_len,
19     bos_token_id=tokenizer["BOS_None"],
20     eos_token_id=tokenizer["EOS_None"],
21     func_to_get_labels=get_label_from_file_path
22 )
```

Iz paketa *MidiTok* učitane su klasa *DatasetMIDI* i metoda *split\_file\_for\_training*. U varijable *max\_seq\_len* pohranjen je veličina sekvenci u koje MIDI datoteke dijelimo, a u varijablu *chunks\_path* put do direktorija u koji ćemo spremiti novonastale sekvence. Metodom *split\_file\_for\_training* podijelili smo datoteke, te klasom *DatasetMIDI* stvorili skup podataka za treniranje. Bitno je primjetiti kako sad koristimo tokene *BOS*, engl. *Beginning of Sequence* i *EOS*, engl. *End of Sequence* koji označuju kada nova sekvenca započinje, odnosno završava. Također, *DatasetMIDI* prima i funkciju za izluči-

vanje oznake klase MIDI datoteke `func_to_get_labels`, ona je implementirana u funkciji `get_label_from_file_path` koja iz putanje do sekvence izluči oznaku klase. Oznaka klase u programskom jeziku *Python* jednostavno izluči iz strukture direktorija pokazane u poglavlju Podjela podataka, 3.3.

Nakon predstavljanja bitnih alata za stvaranje rječnika tokena iz MIDI datoteka možemo krenuti učiti duboke modele.

## 5.2. Povratne neuronske mreže

Povratne neuronske mreže ili RNN, engl. *Recurrent Neural Networks* tip su neuronskih mreža za rad sa sekvencijalnim podacima. Pogodne su za zadatke gdje je redosljed informacije bitan. Obrada prirodnog jezika, prepoznavanje govora, generiranje teksta ili obrada vremenskih serija neki su od zadataka koje RNN može riješiti.

RNN dakle pretpostavlja međusobnu zavisnost o podacima koji mu dolaze na ulaz, u sebi sadrži mehanizam pamćenja prijašnjih vrijednosti te prilikom računanja izlaza ga uzima u obzir. Taj mehanizam RNN-a možemo nazvati skrivenim stanjem. Matematički, skriveno stanje računa se na idući način:

$$h_t = \sigma(W_h \cdot h_{t-1} + W_x \cdot x_t + b),$$

gdje su:

- $h_t$ : skriveno stanje u trenutku  $t$
- $h_{t-1}$ : skriveno stanje u trenutku  $t-1$
- $x_t$ : ulazne vrijednosti u trenutku  $t$
- $W_h$  i  $W_x$ : težinske matrice
- $b$ : slobodni parametar (*bias*)
- $\sigma$ : aktivacijska funkcija, npr. tanh ili ReLU

Unatoč svojim korisnim svojstvima RNN-ovi nailaze nekoliko problema. Jedan od

tih problema su eksplodirajući ili nestajući gradijenti. Tijekom treniranja modela, dovođenjem sve više primjeraka sekvenci učenje gradijenata metodom propagacije unatrag zbog učestalog matričnog množenje neke gradijente svede na veoma mali broj (nestajanje gradijenata) dok neke na vrlo velik broj (rast gradijenata). Time je treniranje otežano. Također, zbog ovih problema mreža ne može *pamtiti* preduge sekvence. RNN-ovi su stoga zamijenjeni ili nadopunjeni naprednijim arhitekturama kao što su LSTM, GRU ili u današnje vrijeme transformeri. U nastavku ćemo malo jasnije objasniti LSTM i GRU.

### 5.2.1. LSTM i GRU

LSTM, engl. *Long Short-Term Memory* mreža razvijena je kako bi riješila problem eksplodirajućih ili nestajućih gradijenata. Kao što naziv govori LSTM ćelija ima sposobnost pamtiti dugoročne ovisnosti u podacima. Sastoji se od četiri ključnih mehanizama koji kontroliraju informaciju:

1. Stanje ćelije  $C_t$ : Vrijednost koja pamti dugoročne informacije o ulaznim sekvencama. Može se zamisliti kao neka vrsta memorije koja se prenosi iz koraka u korak. Ovaj mehanizam opisan je sljedećom jednačbom:

$$C_t = f_t \cdot C_{t-1} + i_t \cdot C'_t,$$

gdje su  $f_t$  zaboravna vrata,  $i_t$  ulazna vrata,  $C_{t-1}$  prethodno stanje ćelije, a  $C'_t$  doprinos stanju ćelije.

2. Ulazna vrata  $i_t$  i doprinos stanju ćelije  $C'_t$ : Ulazna vrata odlučuju koje informacije će biti dodane u ćelijsko stanje. Ove vrijednosti računaju se sljedećim jednačbama:

$$i_t = \sigma(W_{ihh} \cdot h_{t-1} + W_{ixh} \cdot x_t + b_f)$$

$$C'_t = \tanh(W_{chh} \cdot h_{t-1} + W_{cxh} \cdot x_t + b_{ch}),$$

gdje su  $W_*$  težinske matrice,  $b_*$  slobodni parametri,  $h_{t-1}$  skriveno stanje iz prethodnog koraka, a  $x_t$  trenutni ulaz.

3. Zaboravna vrata  $f_t$ : Ova vrata odlučuju koje informacije treba zaboraviti iz ćelij-

skog stanja. Računaju se sljedećom jednađžbom:

$$f_t = \sigma(W_{fhh} \cdot h_{t-1} + W_{fxh} \cdot x_t + b_{fn}).$$

4. Izlazna vrata  $o_t$ : Izlazna vrata odlučuju koje informacije će biti generirane na izlaz u trenutnom vremenskom koraku sekvence. Računa se sljedećim jednađžbama:

$$o_t = \sigma(W_{ohh} \cdot h_{t-1} + W_{oxh} \cdot x_t + b_{oh})$$

Sada kada smo izlistali mehanizme LSTM ćelija kojima se rješava problem eksplodirajućih ili nestajućih gradijenata vidimo koliko je LSTM ćelija kompleksnija od klasične RNN ćelije. Također, ove jednađžbe nećemo objašnjavati kako i zašto rješavaju navedeni problem, one su navedene samo da nam predoče koliko LSTM ćelija ima više parametara  $W_*$  za naučiti i koliko je složena. LSTM ćelija može biti zamijenjena svojom jednostavnijom inačicom s dvojim vratima - GRU.

GRU, engl. *Gated Recurrent Unit* kao što i LSTM, varijanta je RNN ćelije te također rješava problem dugoročne ovisnosti. Jednostavnijeg je dizajna od LSTM-a zadržavajući pritom sličnu učinkovitost. Sastoji se od sljedećih vrata:

1. Vrata ažuriranja  $u_t$ : Vrata ažuriranja određuju koliko prethodnog skrivenog stanja  $h_{t-1}$  treba prenijeti u novo skriveno stanje  $h_t$ . Ova vrata omogućavaju GRU da dugoročno zadrži informacije:

$$u_t = \sigma(W_{uhh} \cdot h_{t-1} + W_{uxh} \cdot x_t + b_{uh}).$$

2. Vrata resetiranja  $r_t$ : Odlučuju koliko prethodnog skrivenog sloja ćelija treba zaboraviti u računanju nove vrijednosti skrivenog sloja:

$$r_t = \sigma(W_{rhh} \cdot h_{t-1} + W_{rxh} \cdot x_t + b_{rh}).$$

3. Povratno stanje  $h_t$  sada se računa na sljedeći način:

$$h_t = u_t \cdot h_{t-1} + (1 - u_t)h'_t$$

$$h'_t = \sigma(W_{hh} \cdot (r_t \cdot h_{t-1}) + W_{xh} \cdot x_t + b_h).$$

### 5.2.2. Pozornost

Mehanizam pozornosti, engl. *attention*, može se koristiti u zadacima klasifikacije povratnim neuronskim mrežama (LSTM, GRU) kako bi se poboljšala sposobnost modela da prepozna ključne dijelove ulazne sekvence bitne za klasifikaciju. Pozornost dakle omogućava modelu da se usredotoči na određene dijelove sekvence. U kontekstu klasifikacije podataka korištenjem RNN-ova mehanizam pozornosti funkcionira na sljedeći način:

1. Generiranje skrivenih stanja: Kao što znamo RNN (LSTM ili GRU) ćelije iz skupa ulaznih podataka  $X = [x_1, x_2, x_3, \dots, x_t]$  generiraju skup skrivenih stanja  $H = [h_1, h_2, h_3, \dots, h_t]$  do trenutka  $t$ .
2. Računanje težina pozornosti: Za svaki  $t$  iz  $H$  izračuna se vektor pozornosti  $e$  ili *score* koji mjeri važnost pojedinog skrivenog stanja  $h_t$ :

$$e_t = v^T \cdot \tanh(W_h h_t + b_h),$$

gdje su  $W_h$  i  $b_h$  parametri koji se uče,  $v$  vektor za ocjenu važnosti.

3. Normalizacija *score*-ova: Normalizacijom *score*-ova  $e_t$  korištenjem funkcije *softmax* dobivamo težine pozornosti  $\alpha_t$ :

$$\alpha_t = \text{softmax}(e_t).$$

4. Računanje kontekstnog vektora  $c$ :

$$c = \sum_t (\alpha_t \cdot h_t).$$



5. Klasifikacija: Tako izračunat kontekstni vektor  $c$  sada se koristi kao ulaz u završni klasifikacijski sloj modela. Na primjer izlaz modela  $y$  mogao bi računati ovako:

$$y = \text{softmax}(W_c \cdot c + b_c).$$

### 5.3. Učenje modela

Nakon što smo u poglavlju *Tokenizacija* 5.1. pokazali kako se MIDI datoteke mogu korištenjem tehnike obrade prirodnog teksta pretvoriti u smislene tokene te u poglavlju *Povratne neuronske mreže* 5.2. objasnili kako funkcioniraju LSTM i GRU modeli s ugrađenom pozornošću možemo pripremiti podatke i naučiti duboke modele.

U nastavku je pokazan isječak koda u kojem se određuju hiperparametri treniranja modela, učitavanje sekvenci, odabir modela i učenje na skupu za učenje. Također, nakon svake epohe model će se evaluirati na neviđenom skupu za evaluaciju. Nakon treniranja odabran je model koji najbolje generalizira (ima najmanji gubitak na skupu za evaluaciju) te je njime testiran skup za testiranje i pokazane su performanse tog modela.

Hiperparametri modela:

```
1 epochs = 20
2 batch_size = 64
3
4 learning_rate = 0.0005
5 max_norm = 1.5
6
7 vocab_size = vocab_size
8 embedding_dim = 128
9 model_name = "lstm_custom"
10 layer_type = "lstm"
11 layer_dims = [512, 256, 128]
12 dropout = 0.3
13 num_classes = 10
14
15 model = CustomRNN(
```

```

16     vocab_size ,
17     embedding_dim ,
18     layer_types ,
19     layer_dims ,
20     dropout ,
21     num_classes ,
22     attention=True
23 )
24 criterion = nn.CrossEntropyLoss()
25 optimizer = optim.Adam(model.parameters(), lr=learning_rate)

```

U odsječku koda iznad pokazano je kako odabrati hiperparametre, inicijalizirati model i optimizator modela. Prvo je određen broj epohi u varijabli *epoches* i taj broj će biti 20 u svim modelima koje ćemo učiti, također, varijable *batch\_size* uvijek će biti 64. *Batch\_size*-om određujemo koliko ćemo sekvenci u jednom potezu provući kroz model bez evaluiranja, dakle u svakom koraku optimizacije modela kroz njega je provučeno *batch\_size* sekvenci veličine *max\_seq\_len* (potpoglavlje *MidiTok*) 5.1.2. *Learning\_rate* i *max\_norm* parametri su optimizatora i njima je varijablom *learning\_rate* određena brzina učenja tj. varijablom *max\_norm* najveća vrijednost norme vektora gradijenta čime dodatno sprječavamo eksplodirajuće gradijente. Nakon toga definirani su parametri svojstveni za model koji treniramo, a to su:

1. *Vocab\_size*: Veličina rječnika tokena, ova veličina je objašnjena u poglavlju *MidiTok* 5.1.2.
2. *Embedding\_dim*: dimenzije vektorske reprezentacije svakog tokena. Dakle, prilikom učenja svaki token dobiva svoj vektor nasumičnih vrijednosti koji se kroz korake učenja također uči. Nakon nekog vremena učenja tokeni koji se nalaze u međusobnoj blizini imat će bliže vrijednosti u svojoj vektorskoj reprezentaciji. Ove vrijednosti su u našem slučaju nasumične, no one se također mogu pametnije namijesiti ili naučiti, ali to nije tema ovog rada.
3. *Model\_name*: Ime modela koji treniramo, ovaj parametar poslužit će nam za nazivanje i kategorizaciju modela kako bismo ih lakše usporedili.

4. *Layers\_type*: Tip RNN ćelije koje će model koristiti. Ova varijabla poprima vrijednosti *LSTM* ili *GRU*.
5. *Layers\_dims*: Dimenzije pojedine RNN ćelije.
6. *Dropout*: Vrijednost kojom se regularizira model. Određuje koliki postotka težini će se isključiti tijekom učenja.
7. *Num\_classes*: Broj klasi koje klasificiramo. Određen je u poglavlju *Odabir sklada- telja 3.2.* i iznosi 10.

Funkcijama *CustomRNN*, *nn.CrossEntropyLoss* i *optim.Adam* kreirani su mehanizmi objekti kojima omogućavamo duboko učenje. *Model*-om je stvorena sama mreža, *criterion* računa gubitak mreže, a *optimizer* strategiju optimiranja parametara.

Kreiranje objekata za učitavanje podataka:

```
1 collator = DataCollator(  
2     tokenizer.pad_token_id,  
3     copy_inputs_as_labels=True  
4 )  
5 train_loader = DataLoader(  
6     train_dataset,  
7     batch_size=batch_size,  
8     shuffle=True,  
9     collate_fn=collator  
10 )  
11 valid_loader = DataLoader(  
12     valid_dataset,  
13     batch_size=batch_size,  
14     shuffle=True,  
15     collate_fn=collator  
16 )
```

*Collator* nam omogućava da između svake sekvence stavi *pad\_token* kako bi se naj-avila pojava nove sekvence. *Train\_loader* i *valid\_loader* učitavaju skupove podataka u obliku tokena.

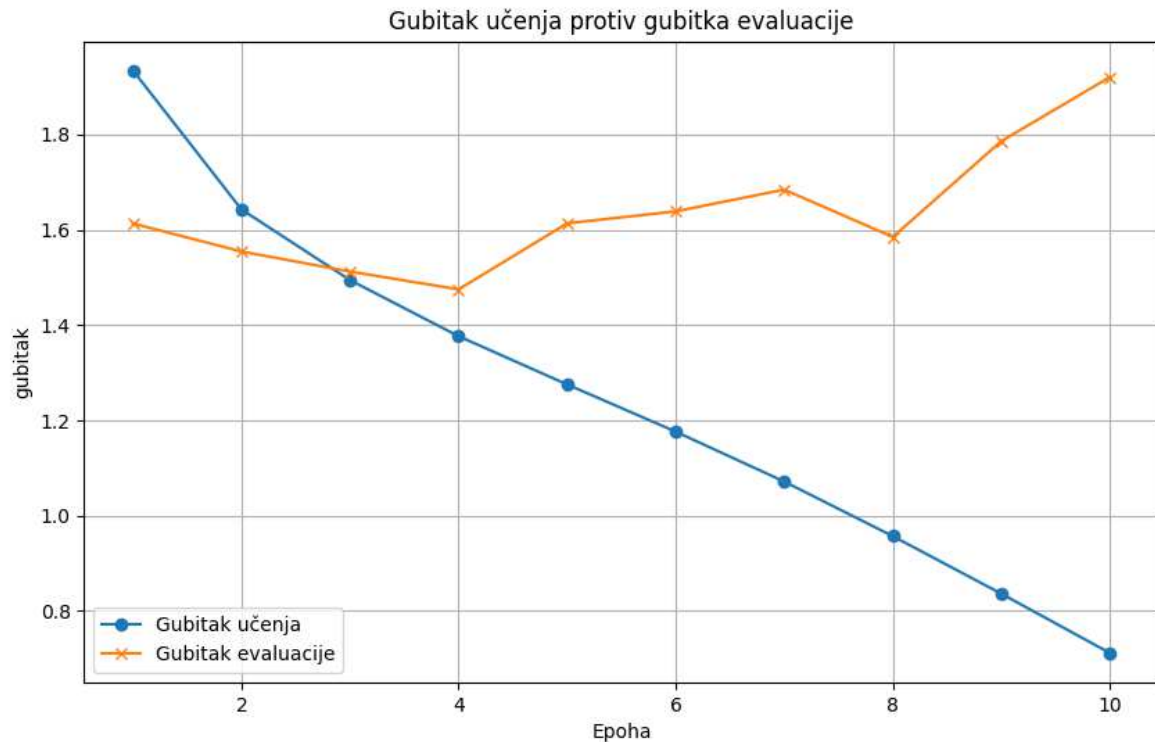
## Učenje modela:

```
1 for epoch in range(epochs):
2     print(f"Epoch: {epoch + 1}")
3     train(
4         model,
5         train_loader,
6         optimizer,
7         criterion,
8         max_norm
9     )
10    results = evaluate(
11        model,
12        valid_loader,
13        criterion,
14        num_classes
15    )
16
17    model_path = Path(
18        f"models/{model_name}/{tokenizer_name}/{epoch + 1}"
19    )
20    os.makedirs(model_path, exist_ok=True)
21    save_checkpoint(
22        epoch + 1,
23        model,
24        optimizer,
25        results,
26        model_path
27    )
```

Učenje modela odvija se u epohama, u svakoj epohi modelu je pokazan cijeli skup podataka. Dakle, za učenje u 10 epohi kroz model je skup podataka prošao 10 puta. Funkcijom *train* uče se parametri modela na skupu za učenje. Implementacija ove funkcije nije pokazana. Funkcijom *evaluate* na skupu za evaluaciju provjeravaju se performanse modela. Računa se točnost i F1-vrijednost te gubitak modela. Nakon svake epohe model

je spremljen na putanju pohranjenu u varijablu *model\_path* kako bi se rezultati mogli usporediti.

Na slici 5.1. prikazana je promjena gubitna na skupu za učenje i skupu za validaciju modela postupcima opisanim u prethodnim koracima. Na slici možemo jasno vidjeti trenutak kada je model na skupu za evaluaciju imao najmanji gubitak i možemo reći da u tom trenu model najbolje generalizira, a taj trenutak je nakon 4. epohe učenja.



**Slika 5.1.** Kretanje gubitka na skupu za učenje i skupu za evaluaciju po epohama

Sada taj model možemo učitati, evaluirati ga na skupu za testiranje i pokazati njegove performanse. Na ovaj način naučit ćemo nekoliko modela i prikazati njihove performanse u poglavlju koje slijedi.

## 5.4. Rezultati

Rezultate treniranja ćemo prikazati u 3 skupine modela po veličini i 2 po tipu ćelije, dakle usporedit ćemo 6 modela. Za svaki tip ćelije odabrat ćemo tri veličine modela: *small*, *medium* i *large*. Kao mjere uspješnosti modela odabrat ćemo *accuracy* i *F1-macro* te matrice zabune. Veličine slojeva pojedinog modela vidimo u nastavku:

```

1 {
2   "model": {
3     "type": "small",
4     "embeddings": 32,
5     "layers": [128]
6   }
7 }
8 {
9   "model": {
10    "type": "medium",
11    "embeddings": 64,
12    "layers": [256, 128]
13  }
14 }
15 {
16   "model": {
17     "type": "large",
18     "embeddings": 128,
19     "layers": [256, 128, 128]
20   }
21 }

```

U tablicama *Accuracy* 5.2. i *F1-score* 5.3. prikazani su rezultati treniranih dubokih modela. Možemo odmah uočiti znatnu razliku u mjerenjima plitkih modela pokazanih u poglavlju *Stroj potpornih vektora* 4. Možemo reći da su *accuracy* i *F1-score* za otprilike 10% viši što nam ukazuje kako je ova metoda efikasnija. Povećanjem slojeva mreže ne možemo reći kako se mjerenja povećavaju, no to je možda zato što smo za pojedini model povećali i veličinu rječnika ovim slijedom: 5000, 7500 i 15000. Kada bismo uzeli jednak rječnik za svaku veličinu možda bismo dobili drugačije vrijednosti. Razlika između tipa ćelije također nije velika i možemo opaziti kako GRU funkcionira malo bolje od LSTM-a.

<i>Accuracy</i>			
Type.	Small	Medium	Large
<i>GRU</i>	0.49	0.52	0.51
<i>LSTM</i>	0.48	0.50	0.50

Tablica 5.2. *Accuracy* ovisno o broju parametara i tipu ćelije

<i>F1-score</i>			
Type.	Small	Medium	Large
<i>GRU</i>	0.37	0.45	0.44
<i>LSTM</i>	0.35	0.40	0.42

Tablica 5.3. *F1-score* ovisno o broju parametara i tipu ćelije

### 5.4.1. Matrice zabune

U matricama zabune modela GRU i LSTM uočavamo kako duboki modeli bolje generaliziraju od plitkih. Matrice potvrđuju i manju razliku u vrijednosti iz tablici *Accuracy* i *F1-macro*. Porastom parametara modela ne uočavamo znatnu razliku, no to se može objasniti kako smo za svaki model povećali i vokabular tokena. Uočavamo kako je model sposobniji uz *Franz Liszta* prepoznati i ostale skladatelje kao što su *Ludwig van Beethoven*, *Carl Czerny*, *Shuwen Zhang* i *Joseph Hayden* što je za ovu vrstu problema impresivno.

	Franz Liszt	Wolfgang Amadeus Mozart	Stephen Heller	Aleksandr Scriabin	Daniel Léo Simpson	Ludwig van Beethoven	Johann Nepomuk Hummel	Carl Czerny	Shuwen Zhang	Joseph Hayden
Franz Liszt	320	1	4	8	17	48	5	37	79	2
Wolfgang Amadeus Mozart	9	28	5	0	0	46	4	5	0	59
Stephen Heller	5	0	2	0	6	4	7	1	13	1
Aleksandr Scriabin	24	0	1	6	0	14	1	3	4	0
Daniel Léo Simpson	0	1	25	0	75	4	15	7	5	5
Ludwig van Beethoven	41	6	5	2	1	256	21	22	5	19
Johann Nepomuk Hummel	9	4	22	0	11	52	24	17	7	5
Carl Czerny	54	1	13	2	1	43	11	150	36	6
Shuwen Zhang	10	0	3	0	1	3	2	3	27	1
Joseph Hayden	0	50	9	0	0	8	5	9	1	73

	Franz Liszt	Wolfgang Amadeus Mozart	Stephen Heller	Aleksandr Scriabin	Daniel Léo Simpson	Ludwig van Beethoven	Johann Nepomuk Hummel	Carl Czerny	Shuwen Zhang	Joseph Hayden
Franz Liszt	347	1	11	12	17	70	10	46	104	1
Wolfgang Amadeus Mozart	4	51	9	0	0	62	5	6	2	60
Stephen Heller	6	0	2	0	3	5	6	1	0	0
Aleksandr Scriabin	10	0	1	3	0	17	0	3	1	0
Daniel Léo Simpson	2	0	13	0	77	4	16	9	6	3
Ludwig van Beethoven	38	4	9	2	2	243	22	38	12	31
Johann Nepomuk Hummel	6	2	16	0	10	32	16	15	18	5
Carl Czerny	42	0	15	0	1	26	10	122	16	4
Shuwen Zhang	17	0	5	1	2	1	6	5	18	0
Joseph Hayden	0	33	8	0	0	18	4	9	0	67

Slika 5.2. Matrica zabune GRU *small* modela

Slika 5.3. Matrica zabune LSTM *small* modela

		Franz Liszt	297	1	12	6	14	75	2	33	31	4
	Wolfgang Amadeus Mozart	2	63	5	0	0	20	0	9	1	29	
	Stephen Heller	13	0	6	0	4	18	9	7	11	1	
	Aleksandr Scriabin	22	1	3	6	0	14	4	3	8	2	
	Daniel Léo Simpson	0	0	13	1	68	7	13	5	0	1	
	Ludwig van Beethoven	37	9	10	3	3	215	21	26	3	28	
	Johann Nepomuk Hummel	18	1	15	1	11	45	23	22	4	13	
	Carl Czerny	38	0	11	0	4	25	9	119	11	4	
	Shuwen Zhang	16	0	2	1	0	2	1	9	96	1	
	Joseph Haydn	2	13	7	0	1	29	6	8	0	80	
predicted label			Franz Liszt	Wolfgang Amadeus Mozart	Stephen Heller	Aleksandr Scriabin	Daniel Léo Simpson	Ludwig van Beethoven	Johann Nepomuk Hummel	Carl Czerny	Shuwen Zhang	Joseph Haydn
			true label									

Slika 5.4. Matrica zabune GRU *medium* modela

		Franz Liszt	291	2	8	5	15	65	3	35	44	2
	Wolfgang Amadeus Mozart	1	48	2	1	0	33	4	5	0	58	
	Stephen Heller	22	1	5	2	6	7	8	6	8	1	
	Aleksandr Scriabin	27	0	4	6	0	16	1	5	8	0	
	Daniel Léo Simpson	0	0	12	0	70	4	12	6	3	0	
	Ludwig van Beethoven	27	9	21	2	0	228	22	23	4	29	
	Johann Nepomuk Hummel	9	3	18	0	7	35	24	17	5	16	
	Carl Czerny	38	0	11	1	1	42	5	129	20	5	
	Shuwen Zhang	30	0	2	1	5	4	4	8	72	0	
	Joseph Haydn	0	25	1	0	1	16	5	7	1	52	
predicted label			Franz Liszt	Wolfgang Amadeus Mozart	Stephen Heller	Aleksandr Scriabin	Daniel Léo Simpson	Ludwig van Beethoven	Johann Nepomuk Hummel	Carl Czerny	Shuwen Zhang	Joseph Haydn
			true label									

Slika 5.5. Matrica zabune LSTM *medium* modela

		Franz Liszt	230	0	7	6	16	31	6	16	18	1
	Wolfgang Amadeus Mozart	4	67	8	1	0	61	3	14	2	52	
	Stephen Heller	27	0	7	0	5	15	10	7	8	4	
	Aleksandr Scriabin	24	0	2	5	0	20	1	4	3	0	
	Daniel Léo Simpson	0	0	9	0	60	0	10	5	0	0	
	Ludwig van Beethoven	45	4	5	3	1	189	15	32	0	11	
	Johann Nepomuk Hummel	15	2	19	0	12	55	25	15	3	13	
	Carl Czerny	53	0	5	0	1	23	4	114	7	1	
	Shuwen Zhang	15	0	8	1	1	2	4	6	112	0	
	Joseph Haydn	0	8	8	0	2	23	5	10	1	68	
predicted label			Franz Liszt	Wolfgang Amadeus Mozart	Stephen Heller	Aleksandr Scriabin	Daniel Léo Simpson	Ludwig van Beethoven	Johann Nepomuk Hummel	Carl Czerny	Shuwen Zhang	Joseph Haydn
			true label									

Slika 5.6. Matrica zabune GRU *large* modela

		Franz Liszt	272	1	8	7	18	76	3	33	18	2
	Wolfgang Amadeus Mozart	3	58	0	0	0	35	3	3	0	35	
	Stephen Heller	12	0	9	1	6	19	14	7	1	7	
	Aleksandr Scriabin	27	2	3	4	1	44	5	9	3	4	
	Daniel Léo Simpson	1	0	8	0	54	1	6	3	0	0	
	Ludwig van Beethoven	16	11	11	2	0	161	14	26	1	29	
	Johann Nepomuk Hummel	5	2	22	0	16	38	23	14	2	27	
	Carl Czerny	49	0	12	1	1	35	5	120	15	7	
	Shuwen Zhang	28	0	3	1	2	7	4	6	114	1	
	Joseph Haydn	0	7	2	0	0	3	6	2	0	38	
predicted label			Franz Liszt	Wolfgang Amadeus Mozart	Stephen Heller	Aleksandr Scriabin	Daniel Léo Simpson	Ludwig van Beethoven	Johann Nepomuk Hummel	Carl Czerny	Shuwen Zhang	Joseph Haydn
			true label									

Slika 5.7. Matrica zabune LSTM *large* modela



## 6. Zaključak

Ovaj rad predstavlja problem klasifikacije MIDI datoteke računalima. Predstavljen je MIDI protokol i objašnjena je njegova namjena: od MIDI portova u fizičkom sloju pa sve do MIDI poruka i MIDI datoteka koje se koriste u računalnim programima. Objasnjene su neke MIDI poruke i što one znače u sklopu MIDI protokola. Za zadatak klasifikacije odabran je skup podataka GiantMIDI-Piano te je pokazano kako su odabrane klase za problem klasifikacije. Nadalje, pokazan je način pristupanja podacima u MIDI datotekama korištenjem programskog jezika *Python* i ekstrakcija bitnih značajki za klasifikaciju melodije. Zatim smo naučili plitke klasifikatore SVM i RF *naivnim* pokušajem pronalaženja cjelini u podacima korištenjem PCA algoritma. Nakon toga iskoristili smo sofisticiranije metode obrade prirodnog jezika i MIDI datoteke smo pretvorili u tokene koje smo zatim pohranili u duboke modele te pokazali kako su uspješniji od plitkih modela.

Obrada MIDI datoteka metodama prirodnog jezika pokazala se obećavajućim putem u istraživanju klasifikacije ili čak generiranja MIDI datoteka. U daljnjim istraživanjima za klasifikaciju mogli bi se i primijeniti modeli transformatora koji su se u proteklih par godina pokazali superiornijim pristupom od LSTM i GRU ćelija.

## Literatura

- [1] Q. Kong, B. Li, J. Chen, i Y. Wang, “Giantmidi-piano: A large-scale midi dataset for classical piano music”, *arXiv preprint arXiv:2010.07061*, 2020.
- [2] S. Deepaisarn, S. Chokphantavee, S. Chokphantavee, P. Prathipasen, S. Buaruk, i V. Sornlertlamvanich, “Nlp-based music processing for composer classification”, *Scientific Reports*, sv. 13, br. 1, str. 13228, 2023.
- [3] J. Rothstein, *MIDI: A comprehensive introduction*. AR Editions, Inc., 1995., sv. 7.
- [4] C. Raffel i D. P. Ellis, “Data with pretty\_midi”.
- [5] J. VanderPlas, *Python data science handbook: Essential tools for working with data*. " O’Reilly Media, Inc.", 2016.
- [6] F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg *et al.*, “Scikit-learn: Machine learning in python”, *the Journal of machine Learning research*, sv. 12, str. 2825–2830, 2011.
- [7] Y.-S. Huang i Y.-H. Yang, “Pop music transformer: Beat-based modeling and generation of expressive pop piano compositions”, u *Proceedings of the 28th ACM international conference on multimedia*, 2020., str. 1180–1188.
- [8] N. Fradet, J.-P. Briot, F. Chhel, A. E. F. Seghrouchni, i N. Gutowski, “Midotok: A python package for midi file tokenization”, *arXiv preprint arXiv:2310.17202*, 2023.

# Sažetak

## Model dubokog učenja za klasifikaciju MIDI datoteka

Bartol Rod

Klasifikacija glazbe, točnije MIDI datoteka, zanimljiv je pristup istraživanja metodi obrade prirodnog jezika. U ovom radu pokazuje se mogućnost obrade i klasifikacije MIDI datoteka u tu svrhu. Pokazan je način pretvorbe MIDI poruka u tokene te njihovo korištenje u dubokim modelima. Time je jasno postavljena osnova za daljnja istraživanja ovog pristupa.

**Ključne riječi:** MIDI datoteke, duboko učenje, klasifikacija

# Abstract

## Deep Learning Model for MIDI Classification

Bartol Rod

Music classification, specifically MIDI file classification, is an interesting approach to exploring natural language processing methods. This paper demonstrates the possibility of processing and classifying MIDI files for this purpose. The method of converting MIDI messages into tokens and their use in deep models is shown. This clearly establishes a foundation for further research in this approach.

**Keywords:** MIDI files, deep learning, classification