

Zadanie 1

Úvod do problematiky

Naším cieľom bolo vytvoriť program schopný kategorizovať piesne z platformy Spotify do kategórií podľa nálady. K realizácii tohto úkolu sme mali k dispozícii pripravený dataset vo formáte CSV. Tento dokument detailne popisuje kroky, ktoré sme podnikli počas procesu vývoja riešenia. Tieto kroky zahŕňujú predspracovanie dát, exploratívnu dátovú analýzu, tvorbu a tréning neurónovej siete a analýzu výsledkov.

1. časť:

Načítať dáta, predspracovať, natrénovať jednoduchý sieť a vyhodnotenie.

a. Outliery

Zo špecifikácie datasetu a z obrázka č. 1 je zrejmé, že stĺpce 'loudness' a 'duration_ms' obsahujú vychýlené hodnoty. Tieto vychýlenia bolo možné rýchlo identifikovať výpisom minimálnych a maximálnych hodnôt pre stĺpce s numerickými hodnotami. Následne sme tieto vychýlené hodnoty odstránili.

	danceability	energy	loudness	speechiness	acousticness	instrumentalness	liveness	valence	tempo	duration_ms	popularity	number_of_artists
min	0.000	0.000197	-47.046	0.000	0.000	0.000	0.00967	0.000	0.000	-4.273460e+05	0.0	1.0
max	8.375	1.000000	1.519	0.965	0.996	0.994	0.99700	0.995	241.423	1.930821e+09	82.0	19.0

Z vizuálnej analýzy obrázka môžeme pozorovať, že stĺpec 'loudness' má maximálnu hodnotu 1.519. Rovnako si všímame, že stĺpec 'duration_ms' obsahuje záporné hodnoty, ktoré predstavujú vychýlené hodnoty. V nasledujúcom kroku plánujeme tieto vychýlené hodnoty odstrániť.

	danceability	energy	loudness	speechiness	acousticness	instrumentalness	liveness	valence	tempo	duration_ms	popularity	number_of_artists
min	0.000	0.000197	-47.046	0.000	0.000	0.000	0.00967	0.000	0.000	1.408000e+04	0.0	1.0
max	8.375	1.000000	-0.116	0.965	0.996	0.994	0.99700	0.995	241.423	1.930821e+09	82.0	19.0

Vychýlené hodnoty boli odstránené.

b. Odstránenie stĺpcov a null hodnoty

Na nasledujúcom obrázku je možné vidieť, koľko nulových hodnôt obsahujú jednotlivé stĺpce po odstránení vychýlených hodnôt.

#	Column	Non-Null Count	Dtype
0	danceability	11944 non-null	float64
1	energy	11960 non-null	float64
2	loudness	11957 non-null	float64
3	speechiness	11960 non-null	float64
4	acousticness	11960 non-null	float64
5	instrumentalness	11960 non-null	float64
6	liveness	11960 non-null	float64
7	valence	11960 non-null	float64
8	tempo	11960 non-null	float64
9	duration_ms	11948 non-null	float64
10	popularity	11843 non-null	float64
11	number_of_artists	11839 non-null	float64
12	explicit	11960 non-null	bool
13	name	11960 non-null	object
14	url	11960 non-null	object
15	genres	11960 non-null	object
16	filtered_genres	11960 non-null	object
17	top_genre	11791 non-null	object
18	emotion	11960 non-null	object

Na nasledujúcej obrázke je zreteľné, že po odstránení všetkých hodnôt NaN prideme o 419 záznamov. Vzhľadom na celkový počet dát je táto strata minimálna, a preto nevidíme dôvod odstraňovať žiadne konkrétne stĺpce.

#	Column	Non-Null Count	Dtype
0	danceability	11541 non-null	float64
1	energy	11541 non-null	float64
2	loudness	11541 non-null	float64
3	speechiness	11541 non-null	float64
4	acousticness	11541 non-null	float64
5	instrumentalness	11541 non-null	float64
6	liveness	11541 non-null	float64
7	valence	11541 non-null	float64
8	tempo	11541 non-null	float64
9	duration_ms	11541 non-null	float64
10	popularity	11541 non-null	float64
11	number_of_artists	11541 non-null	float64
12	explicit	11541 non-null	bool
13	name	11541 non-null	object
14	url	11541 non-null	object
15	genres	11541 non-null	object
16	filtered_genres	11541 non-null	object
17	top_genre	11541 non-null	object
18	emotion	11541 non-null	object

V ďalšom kroku boli odstránené stĺpce 'name', 'url', 'genres' a 'filtered_genres'. Stĺpce 'name' a 'url' pre nás nepredstavovali žiadnu hodnotnú informáciu. Čo sa týka 'genres' a 'filtered_genres', obsahovali hodnoty v tvare polí reťazcov (array of strings), s ktorými sme v rámci nášho cieľa nemohli efektívne pracovať.

c. Zakódovanie nečíselných stĺpcov

V datasete nám ostali nečíselné stĺpce s viacerými hodnotami ako sú 'top_genre' a 'emotion'. Keďže stĺpec 'emotion' je určený na cieľovú predikciu, v kontexte predspracovania dát nás zaujíma predovšetkým stĺpec 'top_genre'. Tento stĺpec sme následne transformovali pomocou metódy one-hot encoding, tiež známej ako dummy kodovanie. Tento prístup k interpretácii nečíselných hodnôt bol zvolený preto, že hodnoty v spomínaných stĺpcoch nebolo možné hierarchicky zoradiť.

Na nasledujúcom obrázku je zreteľné, že v dôsledku použitia one-hot encodingu nám pribudlo 32 nových stĺpcov, a to napriek tomu, že sme transformovali len jeden pôvodný stĺpec.

#	Column	Non-Null Count	Dtype
0	danceability	11541 non-null	float64
1	energy	11541 non-null	float64
2	loudness	11541 non-null	float64
3	speechiness	11541 non-null	float64
4	acousticness	11541 non-null	float64
5	instrumentalness	11541 non-null	float64
6	liveness	11541 non-null	float64
7	valence	11541 non-null	float64
8	tempo	11541 non-null	float64
9	duration_ms	11541 non-null	float64
10	popularity	11541 non-null	float64
11	number_of_artists	11541 non-null	float64
12	explicit	11541 non-null	int64
13	emotion	11541 non-null	object
14	genre_ambient	11541 non-null	int64
15	genre_anime	11541 non-null	int64
16	genre_bluegrass	11541 non-null	int64
17	genre_blues	11541 non-null	int64
18	genre_classical	11541 non-null	int64
19	genre_comedy	11541 non-null	int64
20	genre_country	11541 non-null	int64
21	genre_dancehall	11541 non-null	int64
22	genre_disco	11541 non-null	int64
23	genre_edm	11541 non-null	int64
24	genre_emo	11541 non-null	int64
25	genre_folk	11541 non-null	int64
26	genre_forro	11541 non-null	int64
27	genre_funk	11541 non-null	int64
28	genre_grunge	11541 non-null	int64
29	genre_hardcore	11541 non-null	int64
30	genre_house	11541 non-null	int64
31	genre_industrial	11541 non-null	int64
32	genre_j-pop	11541 non-null	int64
33	genre_j-rock	11541 non-null	int64
34	genre_jazz	11541 non-null	int64
35	genre_metal	11541 non-null	int64
36	genre_metalcore	11541 non-null	int64
37	genre_opera	11541 non-null	int64
38	genre_pop	11541 non-null	int64
39	genre_punk	11541 non-null	int64
40	genre_reggaeton	11541 non-null	int64
41	genre_rock	11541 non-null	int64
42	genre_rockabilly	11541 non-null	int64
43	genre_ska	11541 non-null	int64
44	genre_sleep	11541 non-null	int64
45	genre_soul	11541 non-null	int64

V rámci predspracovania dát sme boli nútení transformovať hodnoty v stĺpcoch 'explicit' a 'emotion' na číselné reprezentácie. Túto úlohu sme efektívne riešili pomocou mapovania, kde sme využili slovník (dict) s definovanými hodnotami pre príslušné stĺpce.

d. Rozdelenie dát na train, val a testovaciu množinu

Naše dáta sme rozdelili na trénovaciu, validačnú a testovaciu množinu v pomere 8:1:1.

e. Normalizovanie a škálovanie dát

Mnoho algoritmov strojového učenia konverguje rýchlejšie a dosahuje vyšší výkon, keď sú hodnoty vlastností štandardizované a na podobnom mierkovom rozpätí. Pre tento účel som využil nástroj `StandardScaler`, ktorý škáluje hodnoty tak, aby mali priemer rovný nule a štandardnú odchýlku rovnú jednej, pričom zároveň zachováva tvar pôvodnej distribúcie hodnôt. Toto škálovanie zabezpečuje konzistentnosť mierky všetkých číselných vlastností, čo napomáha efektívnemu trénovaniu modelu a jeho následnej interpretácii.

Na štandardizáciu trénovacieho datasetu som použil metódu `fit` implementovanú v `Scaler`-i. Pre validačný a testovací dataset som následne aplikoval iba metódu `transform` tohto už nastaveného `Scaler`-a.

f. Natrénovanie neurónovej siete

Na trénovanie neurónovej siete sme využili knižnicu Sklearn. Konkrétne sme sa rozhodli pre klasifikátor založený na viacvrstvovej perceptrónovej sieti (MLP). Aby sme dosiahli optimálne výsledky trénovania, bolo potrebné nastaviť konfiguráciu siete. V našom prípade sme definovali štruktúru so 4 skrytými vrstvami, ktoré obsahovali neuróny v nasledujúcom rozložení: 128, 128, 64 a 32.

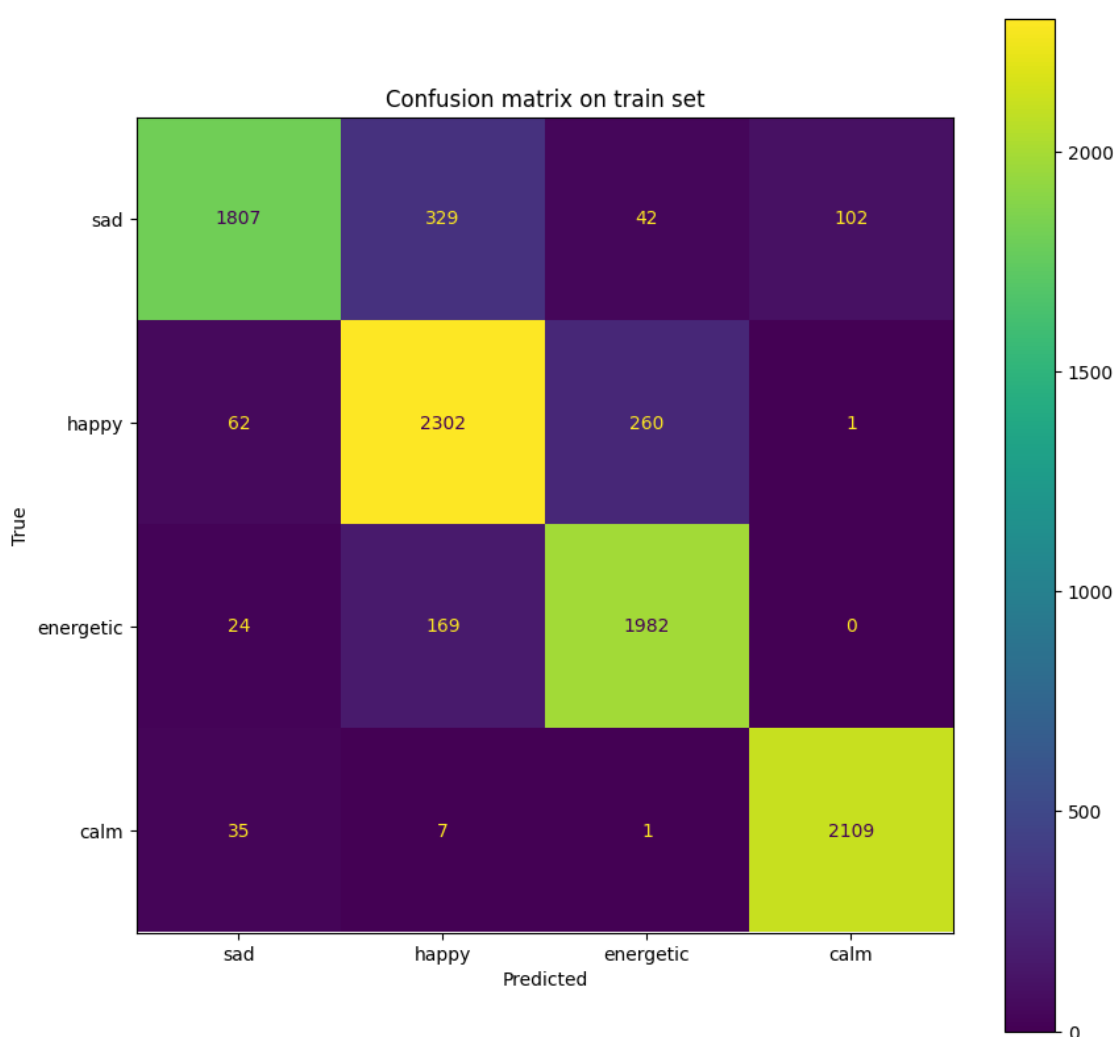
Za účelom zabezpečenia konzistentnosti výsledkov sme atribútu `random_state` prideliť preddefinovanú hodnotu `RANDOM_STATE`.

Maximálny počet iterácií tréovania bol tiež nastavený a pre prevenciu možného pretrénovania siete sme aktivovali funkciu skorého zastavenia tréovania (`early_stopping`) s hodnotou `true`. Táto funkcia efektívne zabraňuje nadmernému zvyšovaniu chyby pri validácii alebo prípadnej stagnácii v zlepšení.

Definovali sme tiež, že 20% tréovacích dát bude využitých ako validačná množina, čo sme špecifikovali v atribúte ``validation_fraction``. Rýchlosť učenia, reprezentovaná hodnotou atribútu ``learning_rate``, bola nastavená na 0.001. Táto hodnota určuje mieru, o akú sú váhy upravované v každej iterácii tréovania.

g. Výsledky tréovania

Z analýzy konfúzných matíc vyplýva, že tréovanie neurónovej siete bolo úspešné. Toto tvrdenie je podložené tým, že najvyššie frekvencie predikovaných hodnôt sa nachádzajú na hlavnej diagonále týchto matíc.

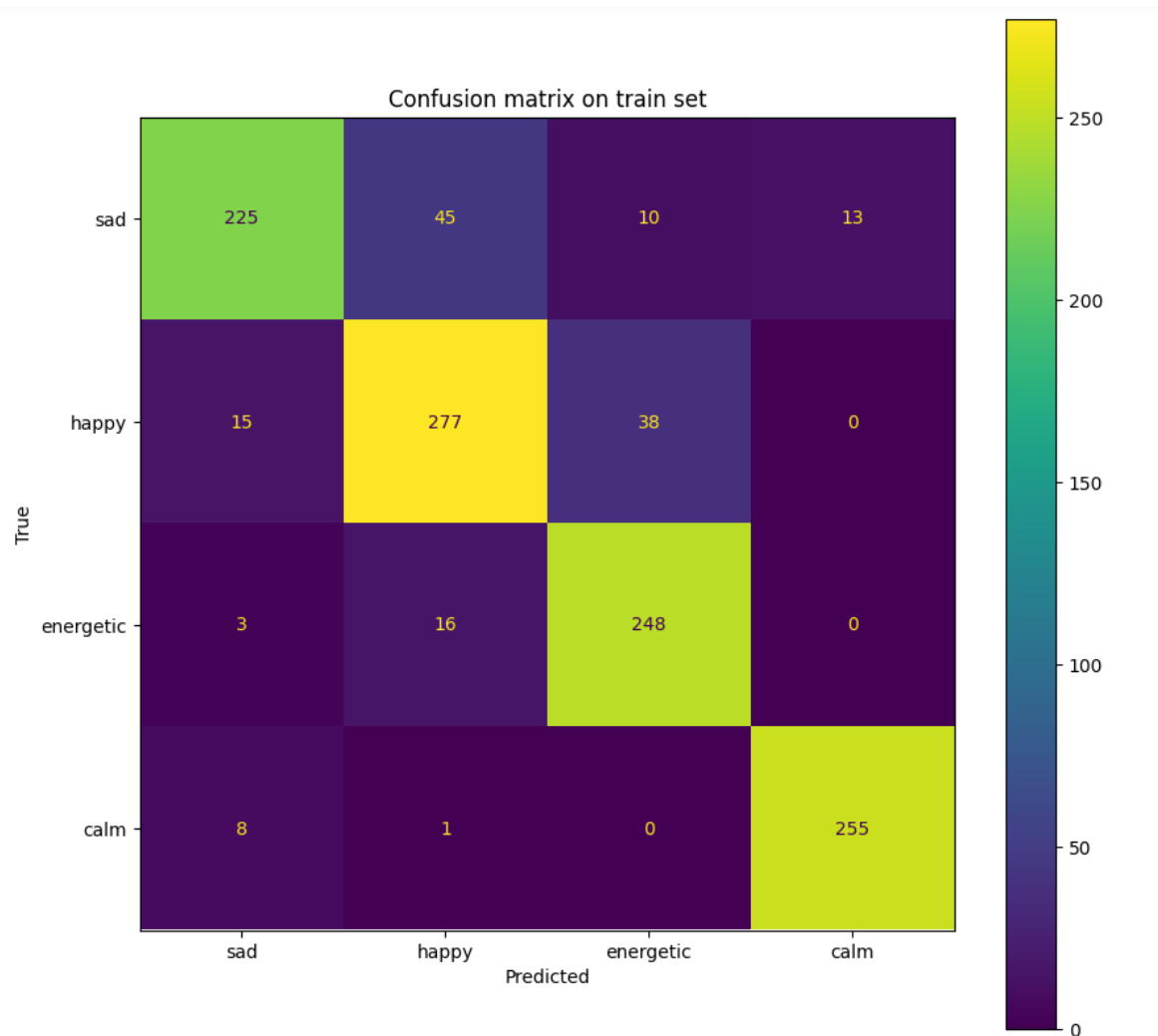


Analýza konfúznej matice pre trénovaciu množinu poskytuje nasledujúce zistenia:

- Emócia 'calm' bola správne predikovaná v 2109 prípadoch. Ostatné predikcie pre 'calm' boli nesprávne, konkrétne 35-krát bola predikovaná ako 'sad', 7-krát ako 'happy' a 1-krát ako 'energetic'.
- Emócia 'energetic' bola správne predikovaná v 1982 prípadoch. Nesprávne predikcie pre 'energetic' boli rozdelené nasledovne: 24-krát bola zamenená za 'sad', 169-krát za 'happy' a žiadny prípad nebol predikovaný ako 'calm'.
- Pre emóciu 'happy' bol správny počet predikcií 2302. Nesprávne predikcie sa rozdeľujú takto: 62-krát bola zmiešaná s 'sad', 260-krát s 'energetic' a 1-krát s 'calm'.
- 'Sad' bola správne identifikovaná v 1807 prípadoch. Chyby pri predikcii 'sad' boli: 329-krát bola zamieňaná za 'happy', 42-krát za 'energetic' a 102-krát za 'calm'.

Na základe výsledkov môžeme konštatovať, že najpresnejšia predikcia emócie v piesni bola v prípade emócie 'happy'. Nasledovala emócia 'calm' s druhou najvyššou presnosťou. Naopak, najnižšiu presnosť sme zaznamenali pri predikovaní emócie 'sad'.

Neurónová sieť dosiahla celkovú presnosť približne 88%, konkrétne hodnotu 88.32%.



Výsledky z testovacej množiny boli v súlade s výsledkami z trénovacej množiny. Najpresnejšia klasifikácia bola dosiahnutá v kategórii 'radosti'. Celková úspešnosť predikcie na testovacej množine bola približne 86%, konkrétne 86.03%.

2. časť:

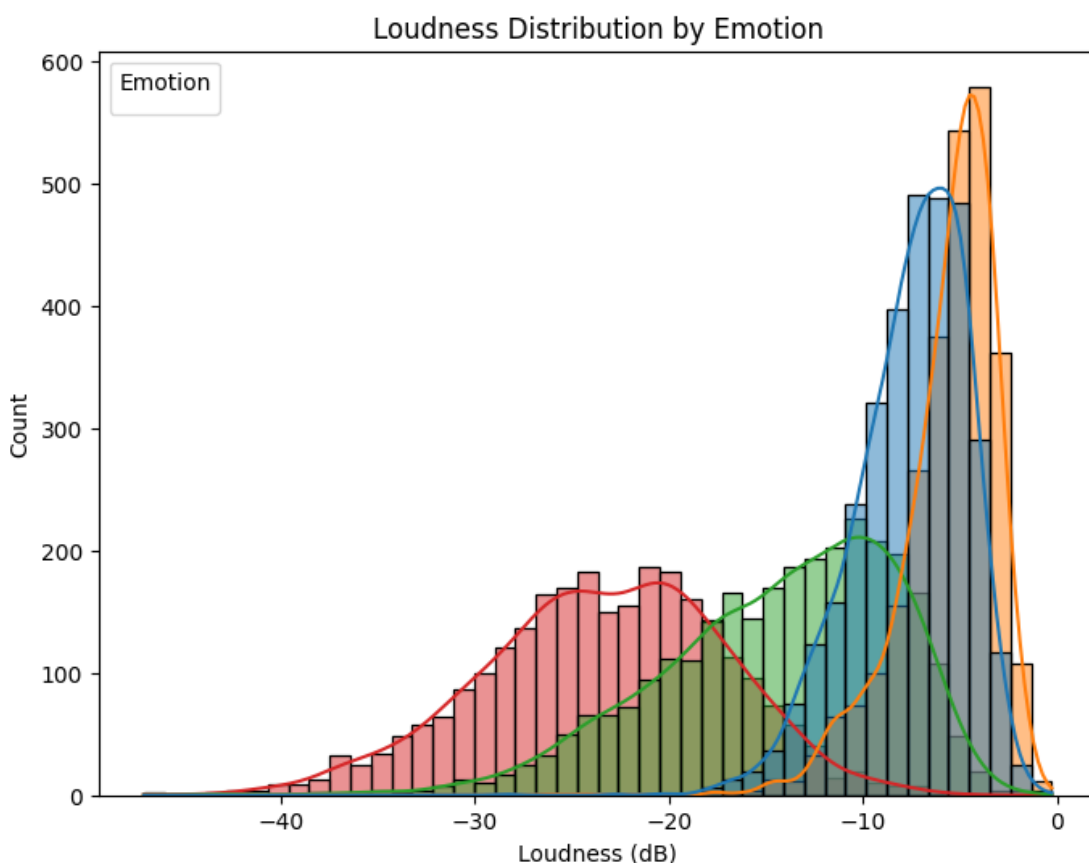
Analýza dát cez EDA.

a. Rozloženie vlastnosti 'loudness' podľa rôznych emócií v stĺpci 'emotion'

Histogram zobrazuje rozloženie vlastnosti "loudness", ktoré je kategorizované podľa rôznych emócií v stĺpci "emotion" dátového súboru.

Stĺpce histogramu sú farebne kódované na základe "emotion" skladieb. To znamená, že uvidíte stĺpce rôznych farieb postavené na seba alebo vedľa seba v závislosti od rozloženia. Každá farba reprezentuje inú emóciu, ako je "sad", "happy", energetic" alebo "calm".

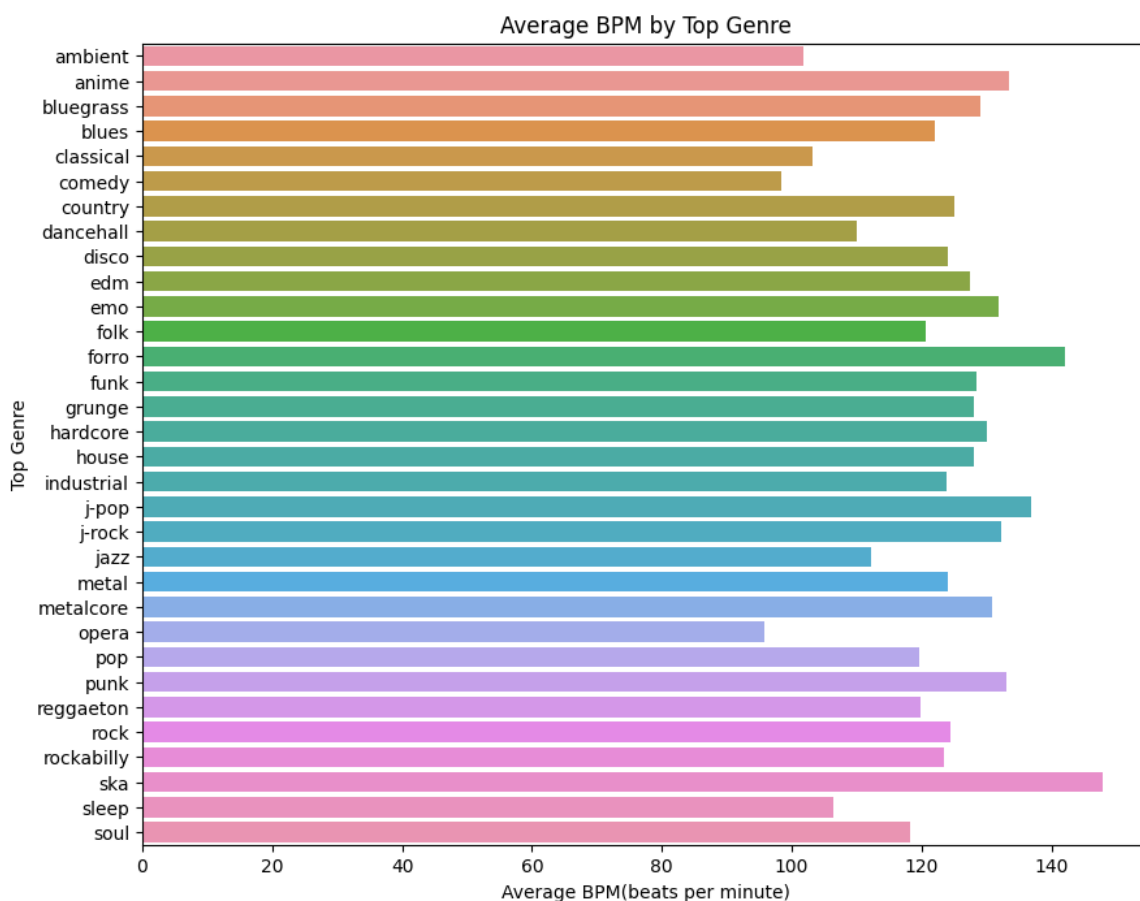
Keďže je kde=True, uvidíte aj krivky hustoty (odhad hustoty jadra) pre každú emóciu nad stĺpcami. Tieto krivky poskytujú vyhladené zobrazenie rozloženia dát.



b. Priemerné tempo (v BPM) skladieb, zoskupených podľa 'top_genre'

Graf je horizontálny stĺpcový diagram, ktorý zobrazuje priemerné tempo (v úderoch za minútu, alebo BPM) skladieb, zoskupených podľa ich "hlavného žánru".

Každý horizontálny stĺpec na diagrame reprezentuje priemerné tempo pre konkrétny žáner. Dĺžka stĺpca zodpovedá priemernej hodnote BPM. Napríklad dlhší stĺpec pre konkrétny žáner naznačuje vyššie priemerné BPM pre skladby v tomto žánre, čo znamená, že skladby v tomto žánre sú vo všeobecnosti rýchlejšie.



rýchlejšie skladby, zatiaľ čo žánre s kratšími stĺpcami naznačujú pomalšie skladby.

Ak je v dĺžkach stĺpcov výrazná variabilita, naznačuje to, že rôzne žánre majú výrazné charakteristiky tempa. Naopak, ak sú mnohé stĺpce

podobnej dĺžky, naznačuje to, že priemerné tempo je pomerne konzistentné naprieč týmito žánrami.

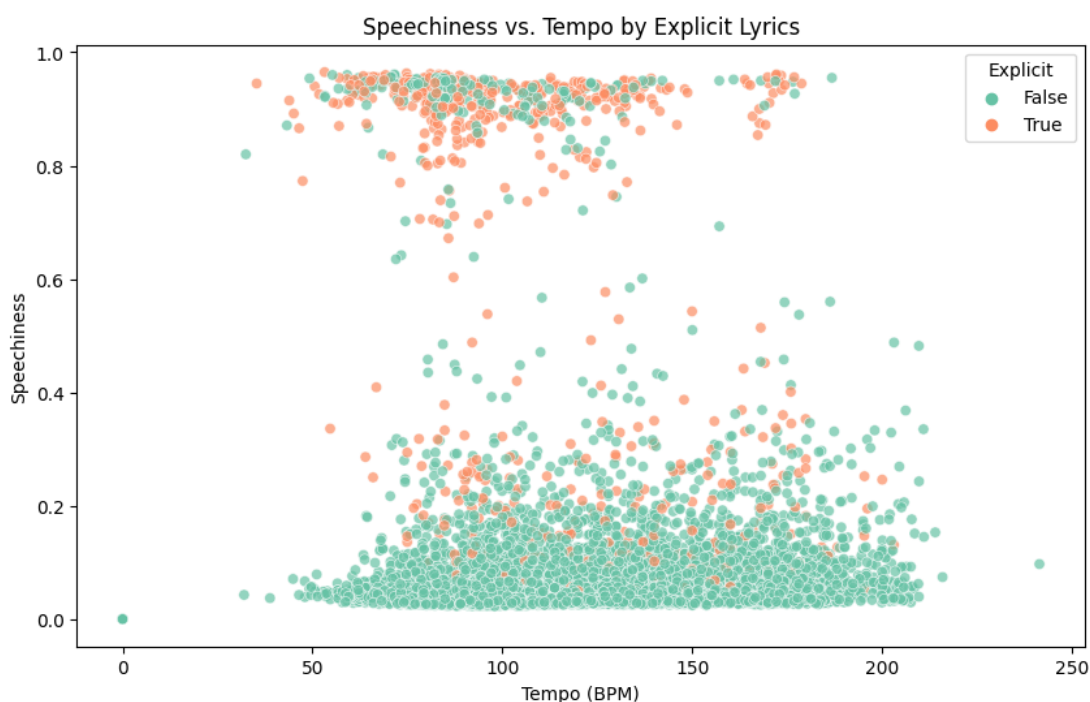
c. Vzťah medzi tempom a speechiness skladieb plus explicitné texty

Graf poskytuje bodový diagram, ktorý ukazuje vzťah medzi tempom a speechiness skladieb, pričom zároveň zvyrazňuje, ktoré skladby majú explicitné texty.

Každý bod na bodovom diagrame reprezentuje jednotlivú skladbu z databázy. Jeho pozícia je určená jej tempom (x-ová súradnica) a 'speechiness' (y-ová súradnica).

Farba každého bodu označuje, či skladba má explicitné texty. Paleta 'Set2' poskytuje odlišné farby pre rôzne hodnoty atribútu 'explicit'.

Napríklad, skladby s explicitnými textami môžu byť ofarbené inak než skladby bez explicitných textov.



Zhrnutie nám umožňuje vidieť rozloženie skladieb z hľadiska ich tempa a speechiness.

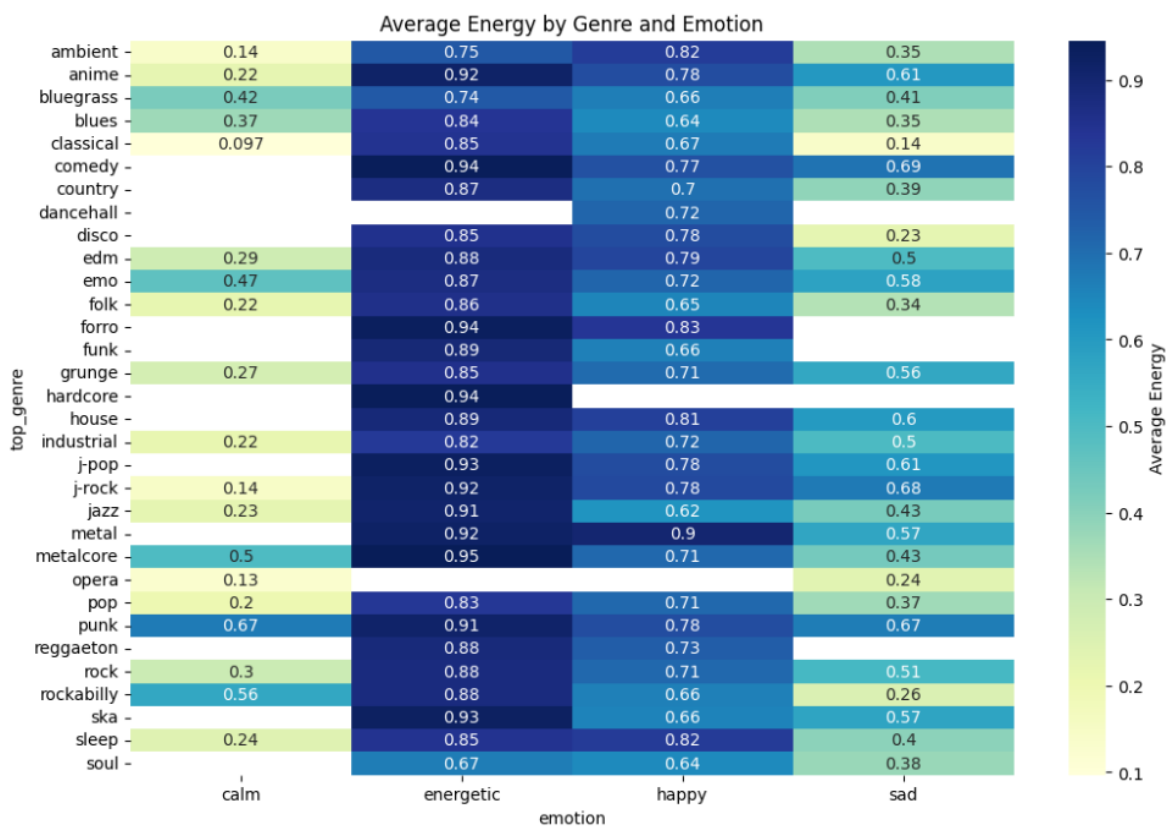
Diferenciácia farieb (na základe explicitných textov) poskytuje vhľady do toho, či skladby s explicitným obsahom majú nejakú koreláciu s tempom alebo 'speechiness'.

Napríklad, ak sú väčšina explicitných skladieb zoskupené v špecifickom regióne (napríklad vysoké tempo a vysoká speechiness), naznačuje to, že skladby s explicitnými textami sú zvyčajne rýchle a podobné reči. Naopak, ak sú explicitné skladby rovnomerne rozložené po celom bodovom diagrame, indikuje to, že neexistuje silný vzťah medzi explicitným obsahom a týmito dvoma vlastnosťami.

d. Priemerná energia skladieb na základe kombinácie top_genre a emotion

Kontingenčná tabuľka (pivot_data) preusporiada údaje v df_eda tak, aby sumarizovala priemernú energiu skladieb na základe kombinácie 'top_genre' a 'emotion'.

Farby v tepelnej mape, určené parametrom cmap='YlGnBu', reprezentujú priemerné energetické úrovne. Tmavšie odtiene značia vyššie energetické úrovne, zatiaľ čo svetlejšie odtiene indikujú nižšie energetické úrovne. Presné mapovanie farieb na energetické úrovne je možné odvodzovať z farebnej lišty na boku tepelnej mapy.



Táto tepelná mapa umožňuje divákovi rýchlo pochopiť, ako sa priemerná energia skladieb líši v rôznych top_genre a emotion. Na prvý pohľad je možné identifikovať, ktoré kombinácie 'top_genre'-'emotion' majú najvyššiu alebo najnižšiu priemernú energiu, čo umožňuje lepšie pochopenie charakteristík skladieb v databáze.

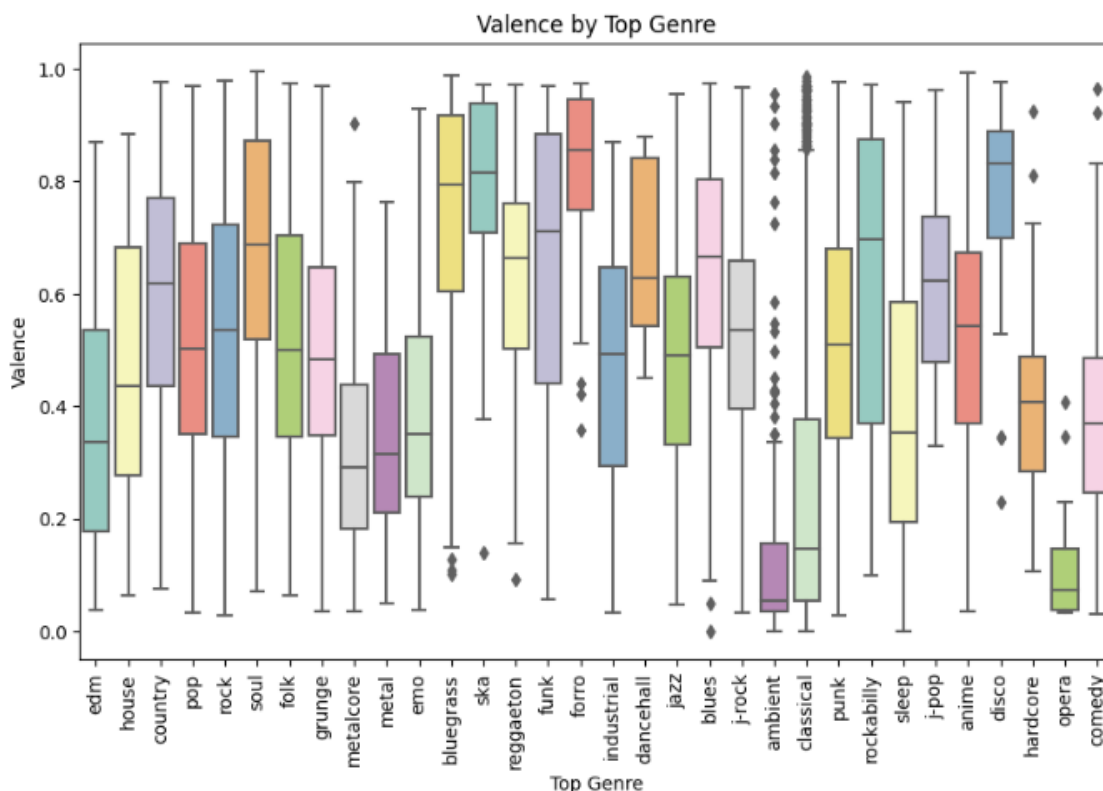
e. Distribúcia hodnôt 'valence' naprieč rôznymi 'top_genre'

Môžeme vidieť boxplot, ktorý je štandardizovaným spôsobom zobrazenia distribúcie dát na základe päťčíselného súhrnu: minimum, prvý kvartil (Q1), medián, tretí kvartil (Q3) a maximum.

Pre každý žáner na x-ovej osi je na y-ovej osi príslušný box, ktorý zobrazuje distribúciu hodnôt valence pre skladby v tomto žánri.

Komponenty boxu sú:

- Spodná čiara boxu: Zastupuje prvý kvartil (Q1) dát, alebo 25. percentil.
- Horná čiara boxu: Zastupuje tretí kvartil (Q3) dát, alebo 75. percentil.
- Čiara vnútri boxu: Zastupuje medián (alebo 50. percentil) dát.
- Fúzy (čiary rozprestierajúce sa nad a pod boxom): Zvyčajne reprezentujú rozsah, v ktorom padajú väčšina dátových bodov, zvyčajne do 1,5-násobku medzikvartilového rozsahu (Q3-Q1) nad Q3 a pod Q1.
- Body mimo fúzov: Sú často považované za odľahlé hodnoty, alebo dátové body, ktoré padajú mimo typického rozsahu dát.



Tento graf vizualizuje distribúciu hodnôt valence naprieč rôznymi top_genre. Skúmaním boxov môžete vidieť, ako sa valence líši podľa žánru, ktoré žánre majú skladby, ktoré sú pozitívnejšie, a ktoré majú skôr negatívne skladby. Rozpätie boxov a odľahlé hodnoty navyše môžu poskytnúť vhľad do variability valence v rámci každého žánru.

3. časť: Výsledky trénovania neurónovej siete (Pytorch)

a. Pretrénovanie siete

Pre pretrénovanie siete sme zvolili tieto parametre:

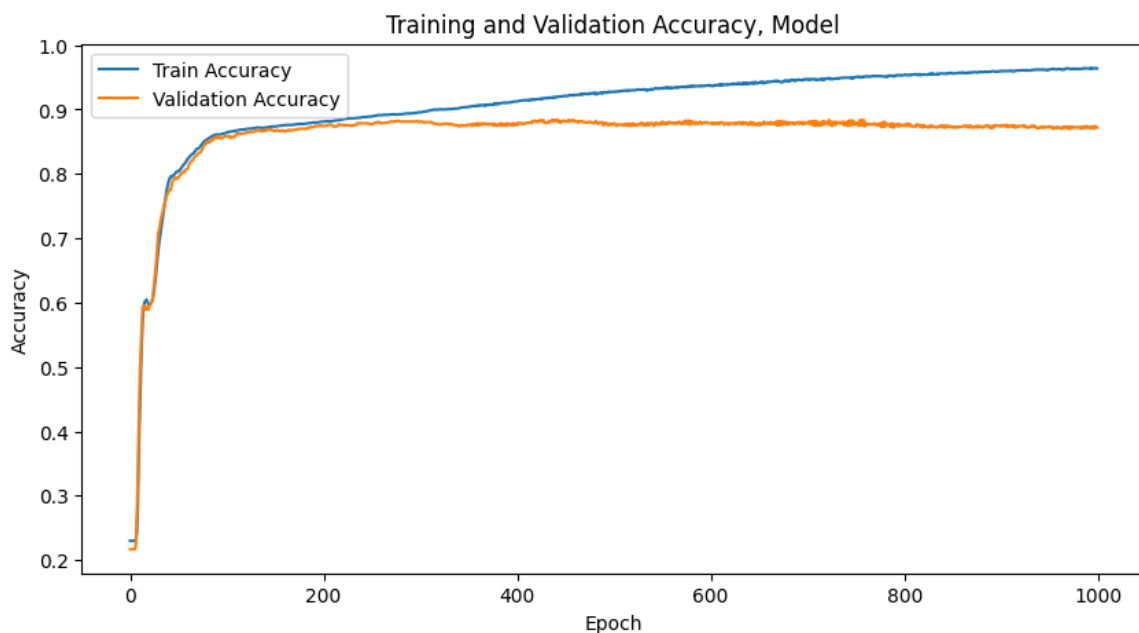
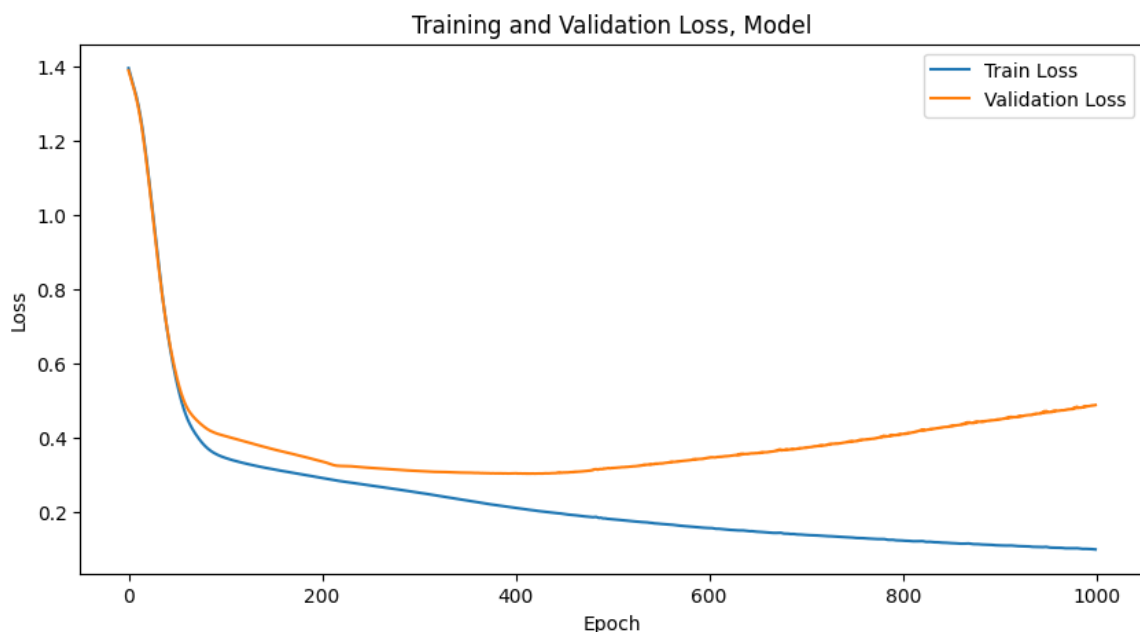
Architektúra siete

- vstupná vrstva, 3 skryté vrstvy, výstupná vrstva
 - vstupná vrstva tvorí rovnaký počet neurónov, aký je počet stĺpcov po predspracovaní datasetu
 - pre skryté vrstvy sme zvolili aktivačnú funkciu ReLU
 - počty neurónov v skrytej vrstve je 128, 64, 32
 - výstupnú vrstvu tvorí rovnaký počet neurónov, ako je pozorovaných hodnôt (4 emotion)
 - pre výstupnú vrstvu sme zvolili aktivačnú funkciu softmax
- kritériálna funkcia - 'categorical_crossentropy'
- použitý solver - Adam, learning rate 0.001
- použité dáta - trénovacie 80%, validačné 10%, testovacie 10%
- počet epoch - 1000

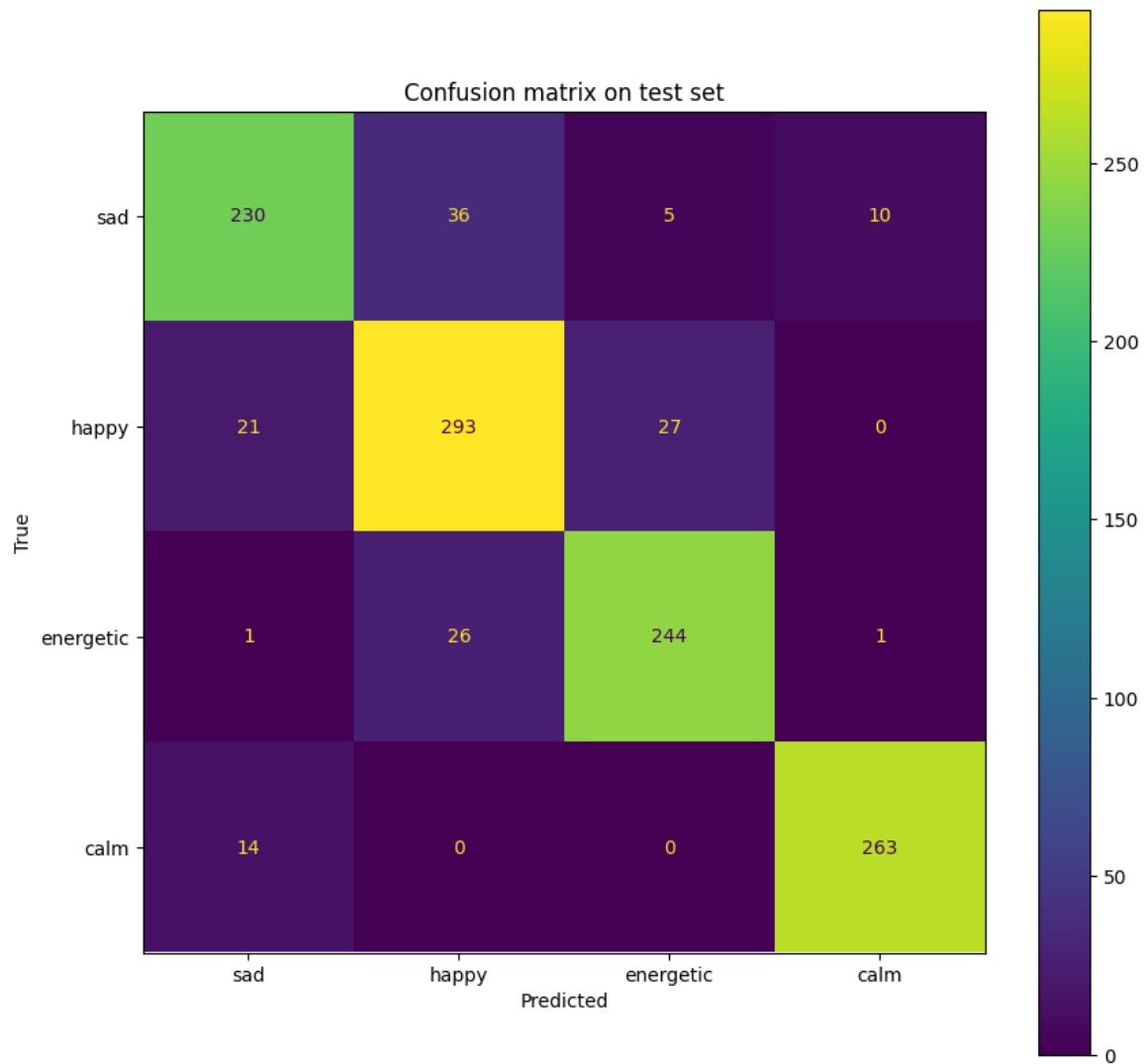
Z grafov vyhodnotenia úspešnosti a priebehu tréovania môžeme vidieť, že na začiatku (v epoche 0) sú obidve straty - trénovacia aj validačná - vysoké, čo značí, že model sa ešte veľa nenaučil.

Modrá čiara predstavuje rýchly pokles trénovacej straty, čo ukazuje rýchle učenie modelu na začiatku. Oranžová čiara, reprezentujúca validačnú stratu, najprv klesá, ale neskôr sa stabilizuje alebo mierne zvyšuje, čo môže naznačovať prispôbenie sa trénovacím dátam.

Rozdiel medzi trénovacou a validačnou stratou sa s časom zväčšuje, čo je znakom pretrénovania. Model sa stáva príliš špecializovaným na trénovacie dáta, čo môže ovplyvniť jeho výkon na nových dátach.



Úspešnosť na tréningových dát je 96.50%, validačných dát 87.62%, testovacích dát 85.82%.



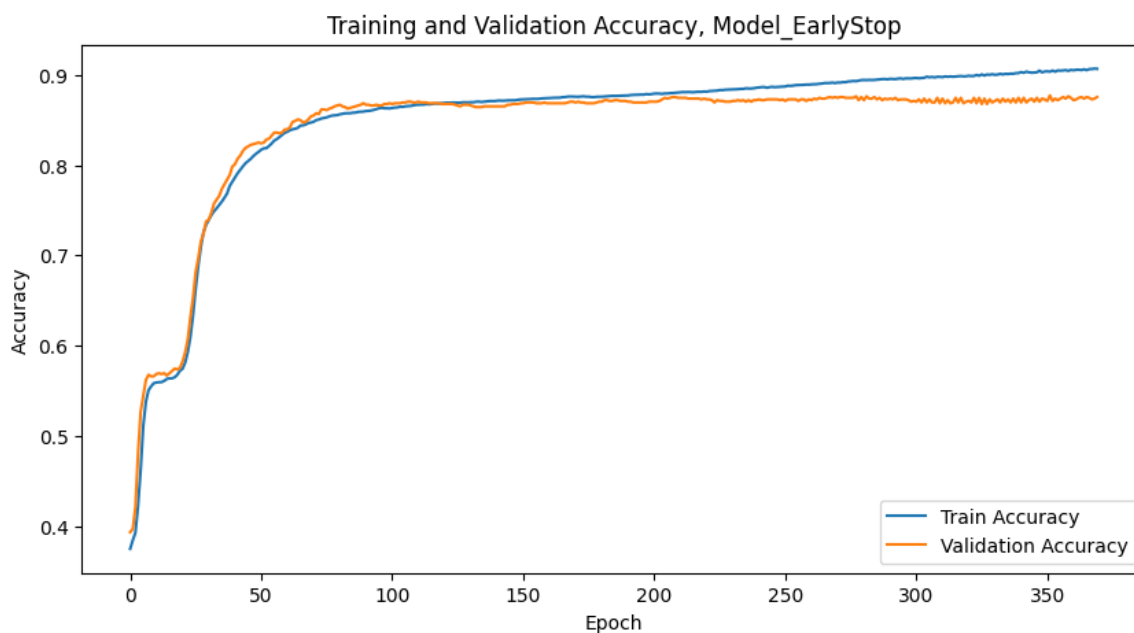
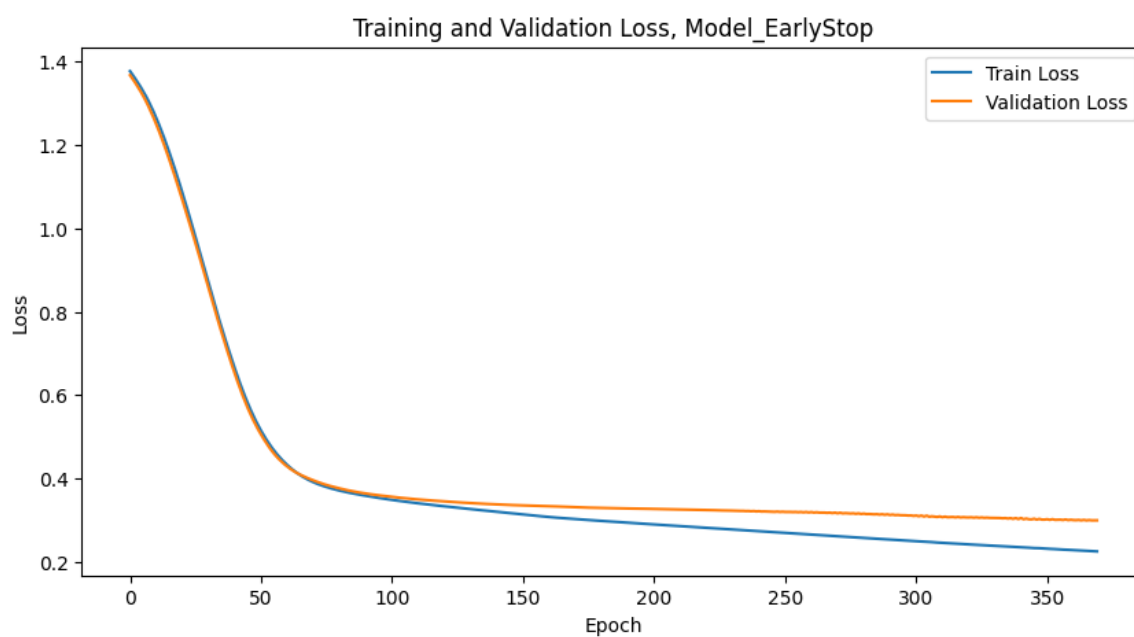
b. Používanie EarlyStoppingu

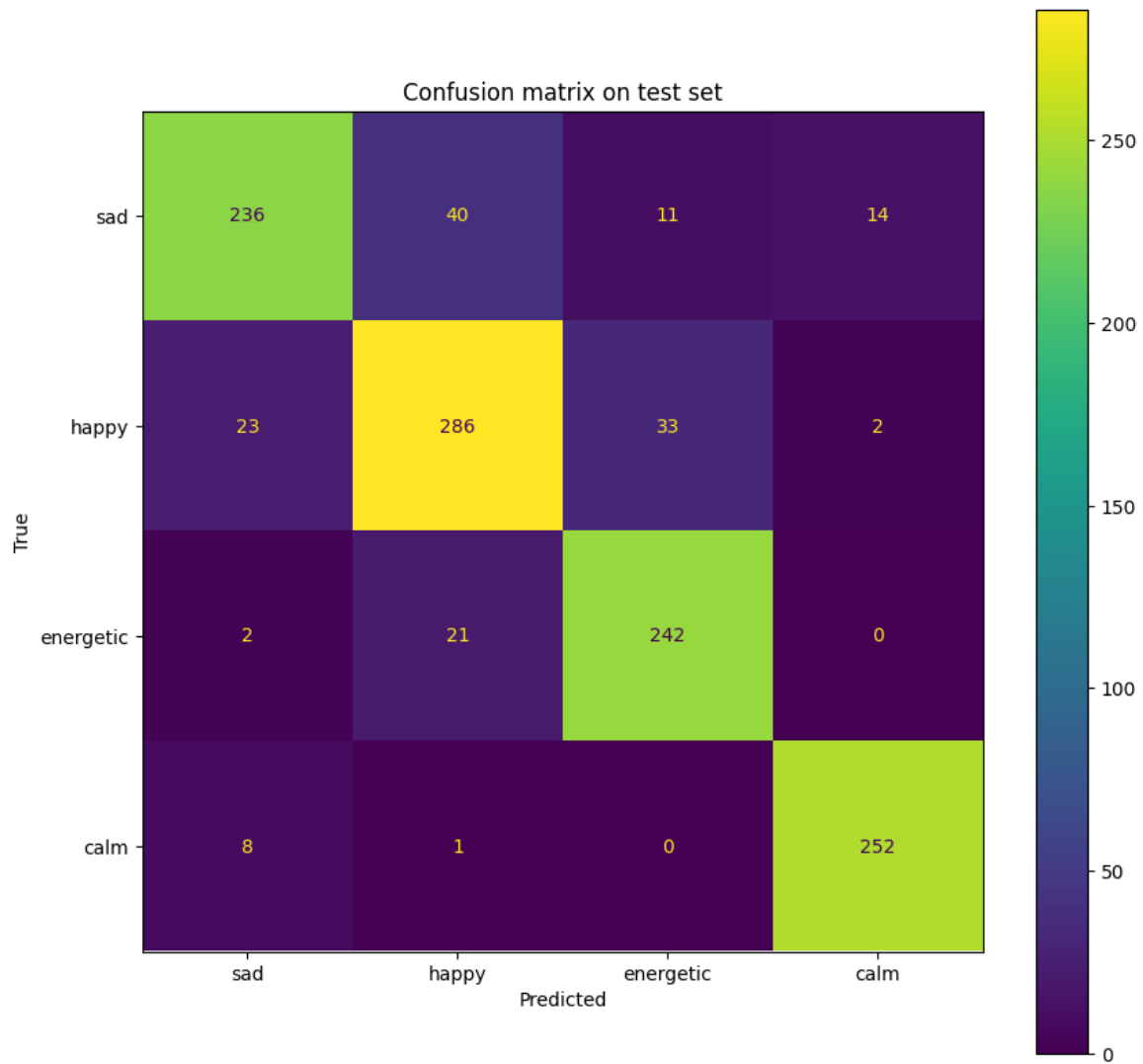
Aby sme predišli pretrénovaniu siete bolo potrebné používať early stopping s parametrami:

- monitorovaná hodnota - val_loss
- počet nezmenených výsledkov pre zastavenie - 12

- obnovenie najlepších váh - true

Úspešnosť vyhodnotenia po zavedení early stoppingu bola na tréningových dát je 90.64%, validačných dát 87.53%, testovacích dát 86.76%.

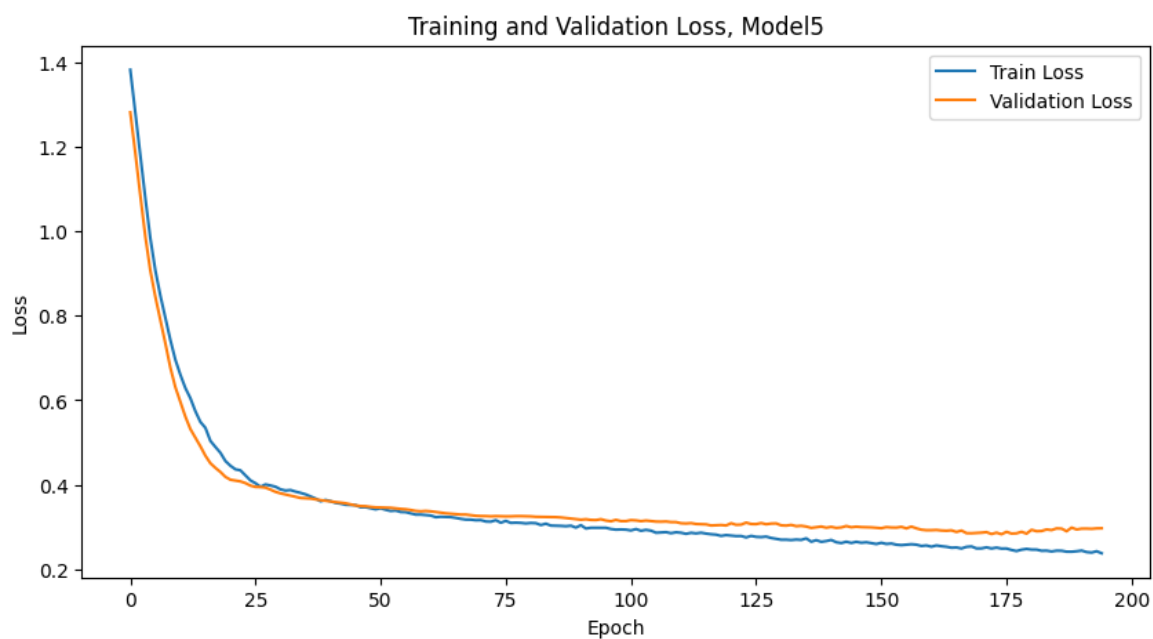


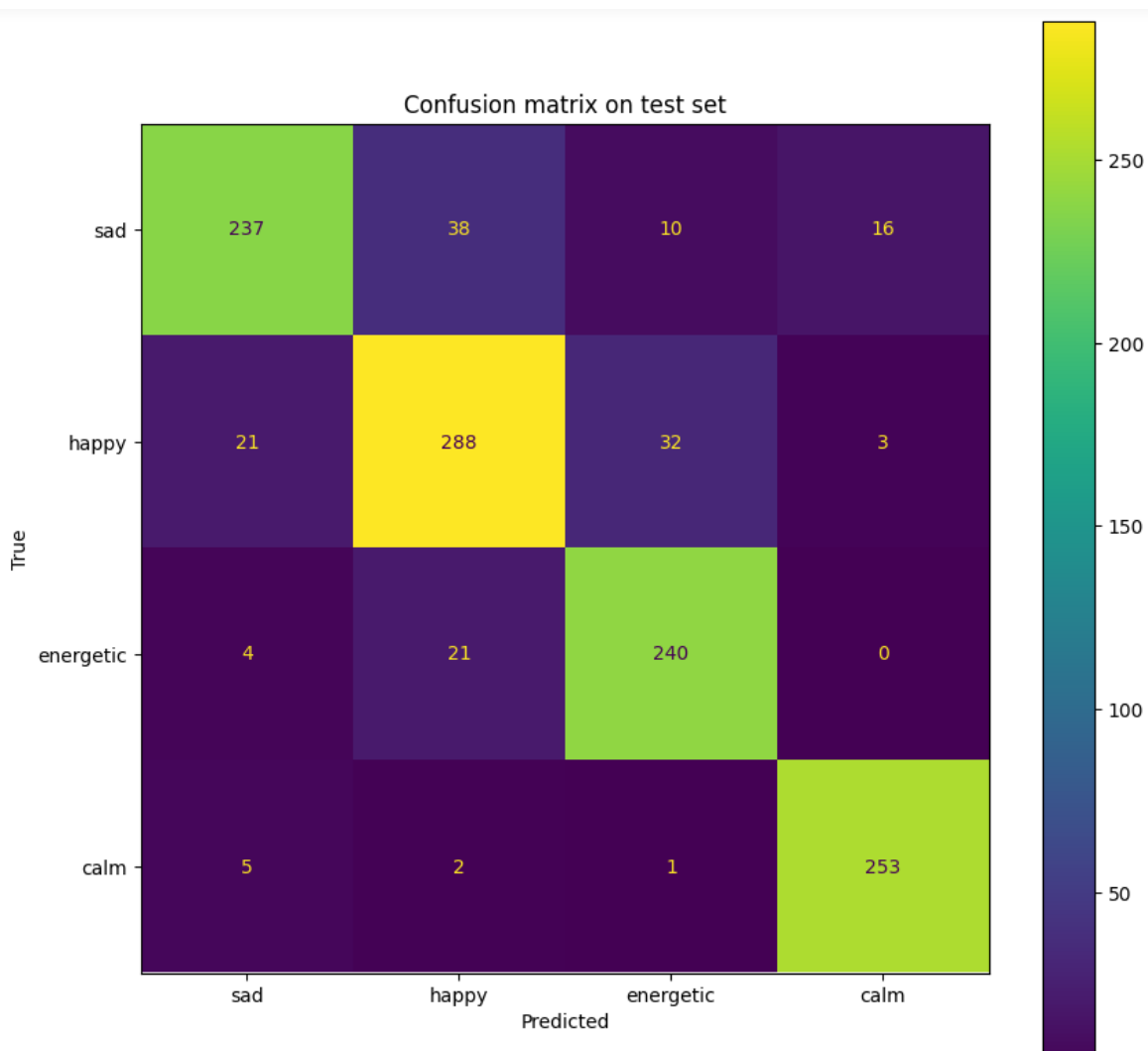
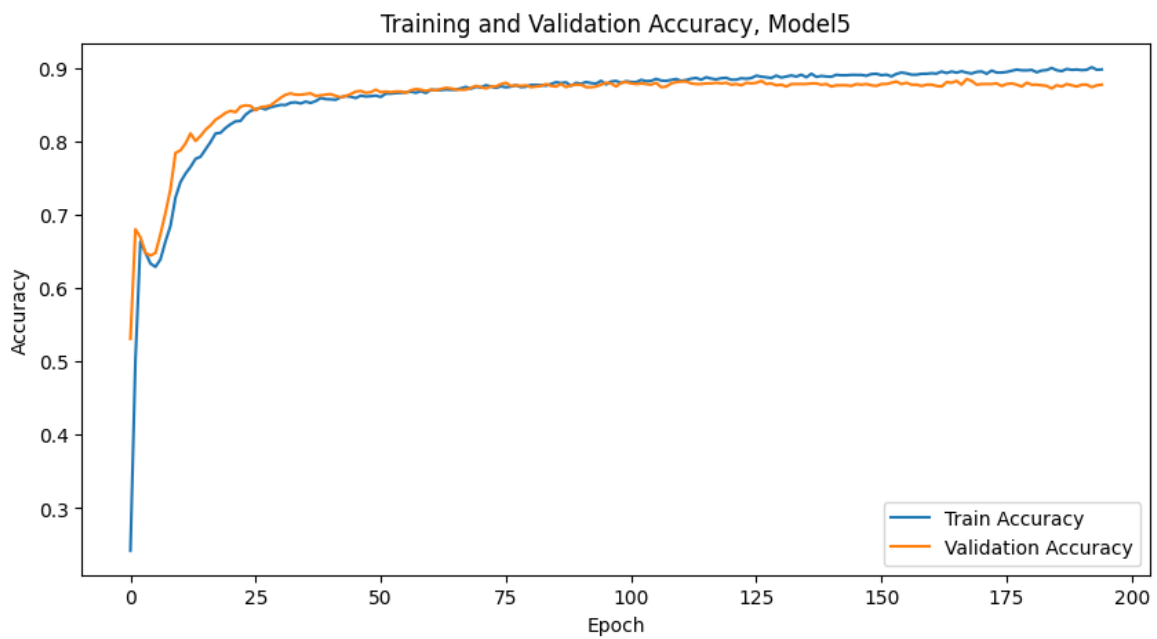


c. Zmena parametrov siete

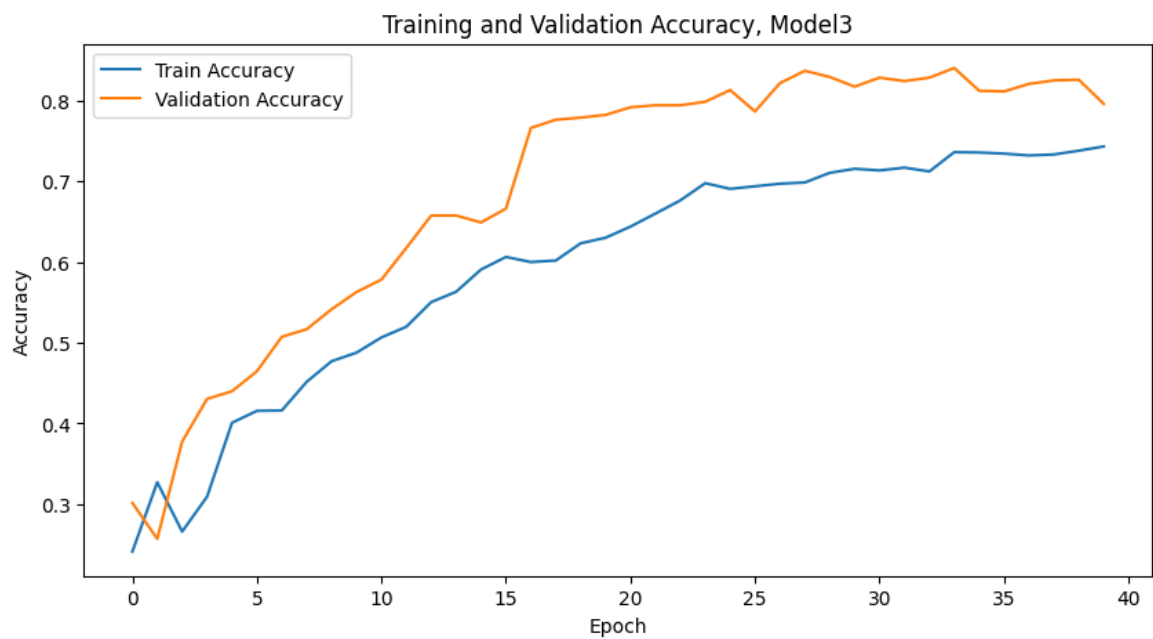
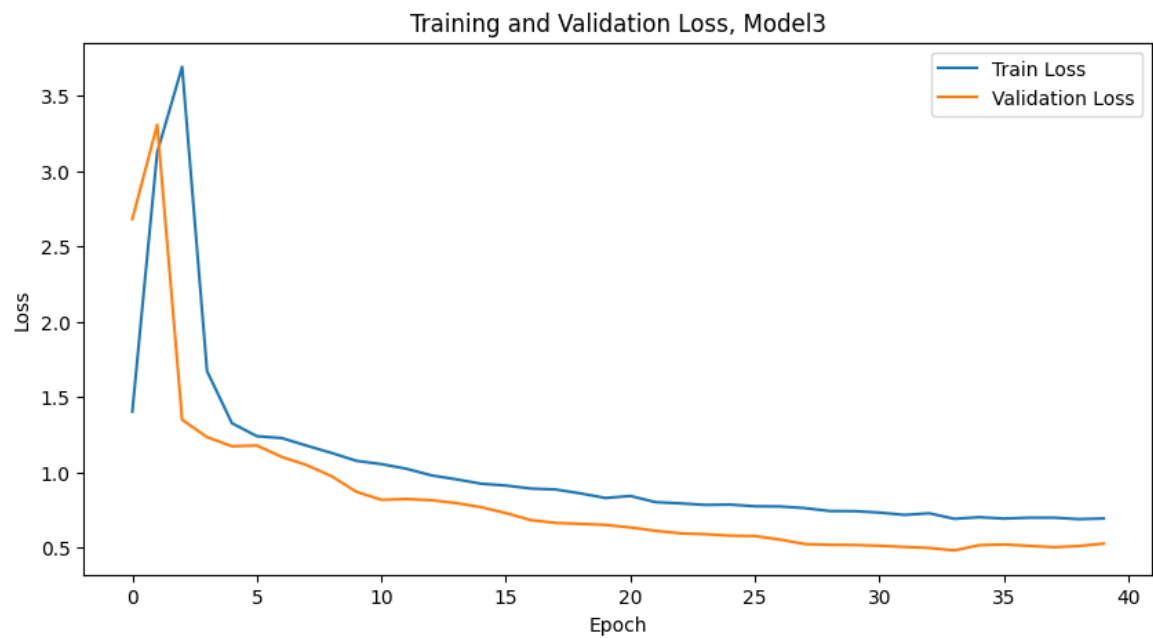
Model	Skryté vrstvy	Počet neurónov	Param rýchlosti učenia	Early stopping patience	Dropout	Epoch	Úspešnosť train/test
1.	4	[128,64,64,32]	0.001	10	0.5	1000	87.31%/ 85.91%
2.	2	[64,32]	0.01		0.3	300	89.89%/ 86.34%
3.	3	[128,64,32]	0.1	6	0.5	1000	80.40%/ 79.50%
4.	2	[64, 32]		10	0.3	300	87.61%/ 85.65%
5.			0.01	20	0.2	400	90.57%/ 86.93%

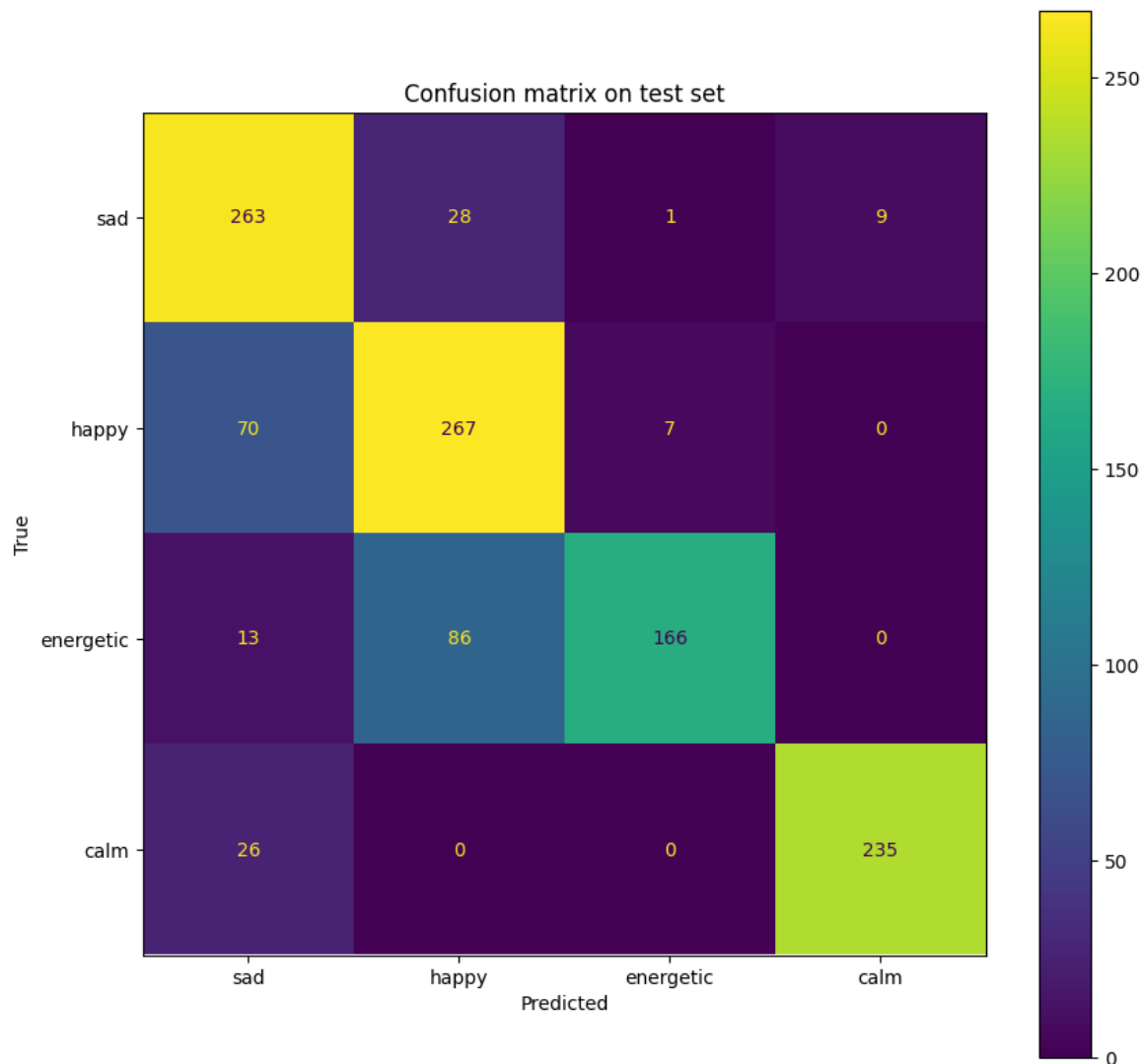
Grafy a konfúzna matica na najlepšie vyhodnotenom testovaní (model 5):





Grafy a konfúzna matica na najhoršie vyhodnotenom testovaní (model 5):





4. časť: Bonus

a. Používanie grid search

V ďalšej úlohe sme boli poverení využitím metódy Grid Search s cieľom identifikovať optimálne hyperparametre pre náš model. Ako nástroj pre túto úlohu som, v kontexte môjho modelu založeného na Pytorch, využil knižnicu 'tune' z balíka 'ray'. Základný model mal nastavený early stopping na true s patience 5, monitorovaná hodnota bola val_loss.

Kombinacia hyperparametrov:

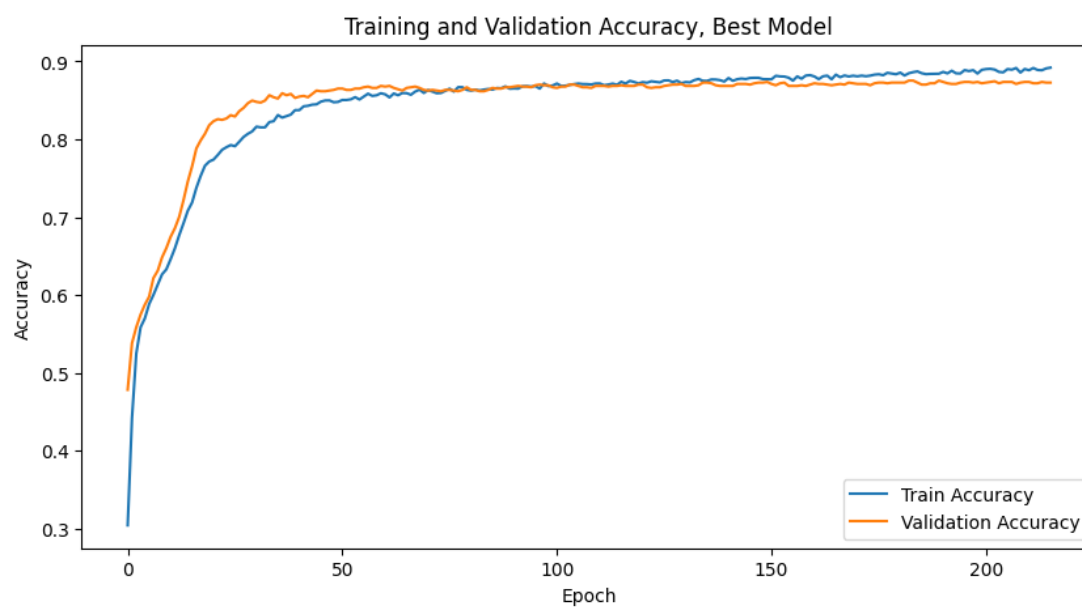
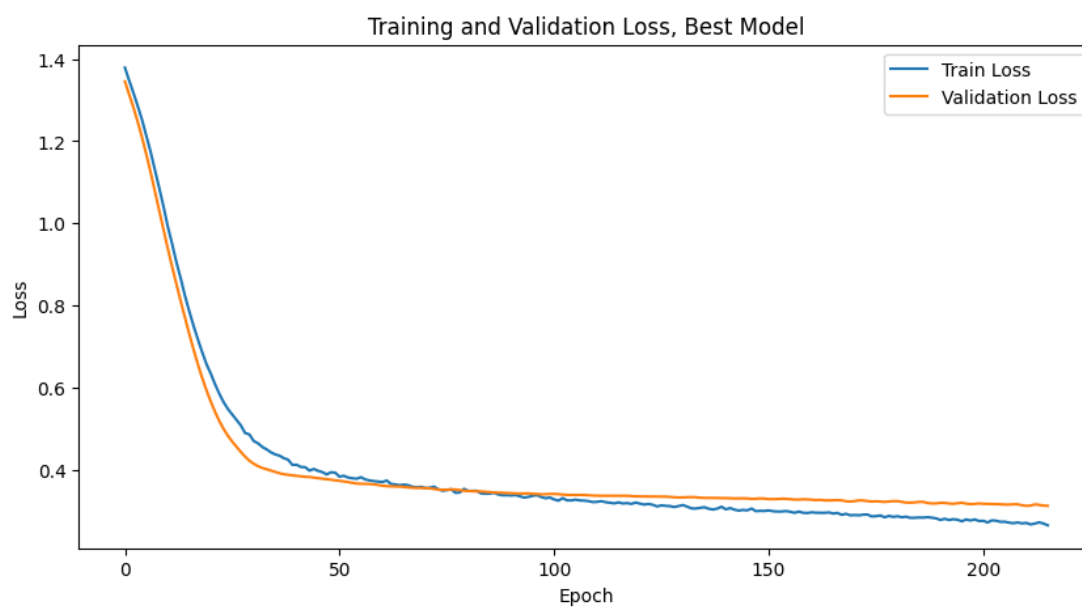
- learning rate: [0.001, 0.01, 0.1]
- skryté vrstvy: [128, 64, 32], [256, 128, 64], [512, 256, 128]
- dropout: [0.3, 0.2, 0.5]
- epochy: 500

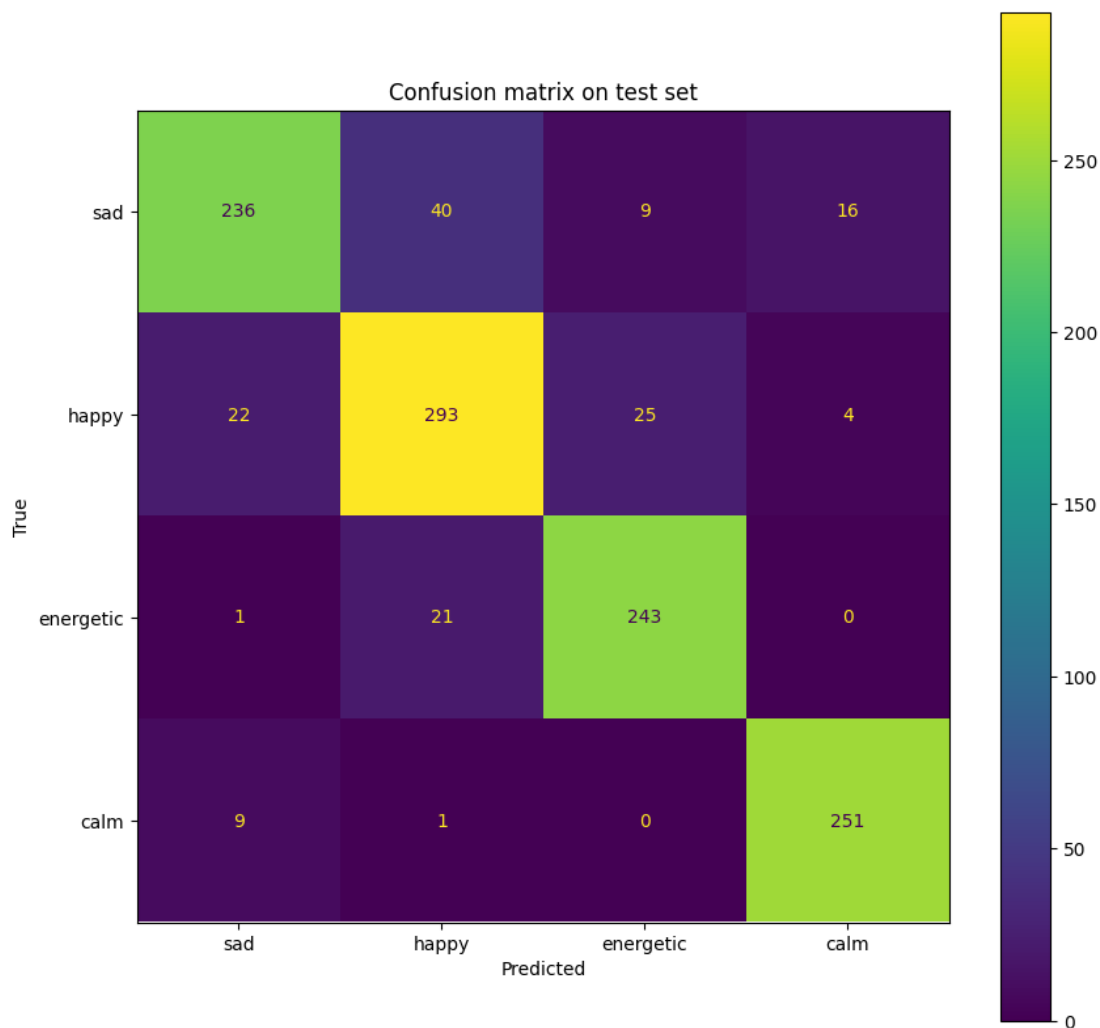
V nasledujúcom obrázku sú prezentované hodnoty 'val_loss' pre všetky skúšané kombinácie hyperparametrov. Dáta sú zotriedené podľa hodnoty 'val_loss', čo nám umožňuje identifikovať, ktoré kombinácie hyperparametrov vedú k najlepším a najhorším modelom.

	config/lr	config/hidden_layer_sizes	config/dropout_prob	val_loss
trial_id				
6522b_00004	0.001	[256, 128, 64]	0.2	0.296244
6522b_00015	0.010	[512, 256, 128]	0.3	0.308377
6522b_00014	0.010	[256, 128, 64]	0.5	0.311617
6522b_00006	0.001	[512, 256, 128]	0.3	0.313592
6522b_00009	0.010	[128, 64, 32]	0.3	0.315158
6522b_00013	0.010	[256, 128, 64]	0.2	0.318332
6522b_00010	0.010	[128, 64, 32]	0.2	0.318377
6522b_00007	0.001	[512, 256, 128]	0.2	0.319893
6522b_00016	0.010	[512, 256, 128]	0.2	0.321054
6522b_00012	0.010	[256, 128, 64]	0.3	0.321404
6522b_00011	0.010	[128, 64, 32]	0.5	0.322975
6522b_00008	0.001	[512, 256, 128]	0.5	0.327663
6522b_00003	0.001	[256, 128, 64]	0.3	0.329009
6522b_00000	0.001	[128, 64, 32]	0.3	0.334247
6522b_00001	0.001	[128, 64, 32]	0.2	0.338100
6522b_00005	0.001	[256, 128, 64]	0.5	0.338960
6522b_00002	0.001	[128, 64, 32]	0.5	0.343926
6522b_00018	0.100	[128, 64, 32]	0.3	0.381548
6522b_00021	0.100	[256, 128, 64]	0.3	0.453785
6522b_00019	0.100	[128, 64, 32]	0.2	0.456178
6522b_00020	0.100	[128, 64, 32]	0.5	0.491153
6522b_00017	0.010	[512, 256, 128]	0.5	0.535887
6522b_00022	0.100	[256, 128, 64]	0.2	0.720243
6522b_00023	0.100	[256, 128, 64]	0.5	0.875835
6522b_00025	0.100	[512, 256, 128]	0.2	1.053185
6522b_00024	0.100	[512, 256, 128]	0.3	1.265706

Úspešnosť vyhodnotenia na najlepší model je 89.66% na tréningových dát a 87.36% na testovacích dát.

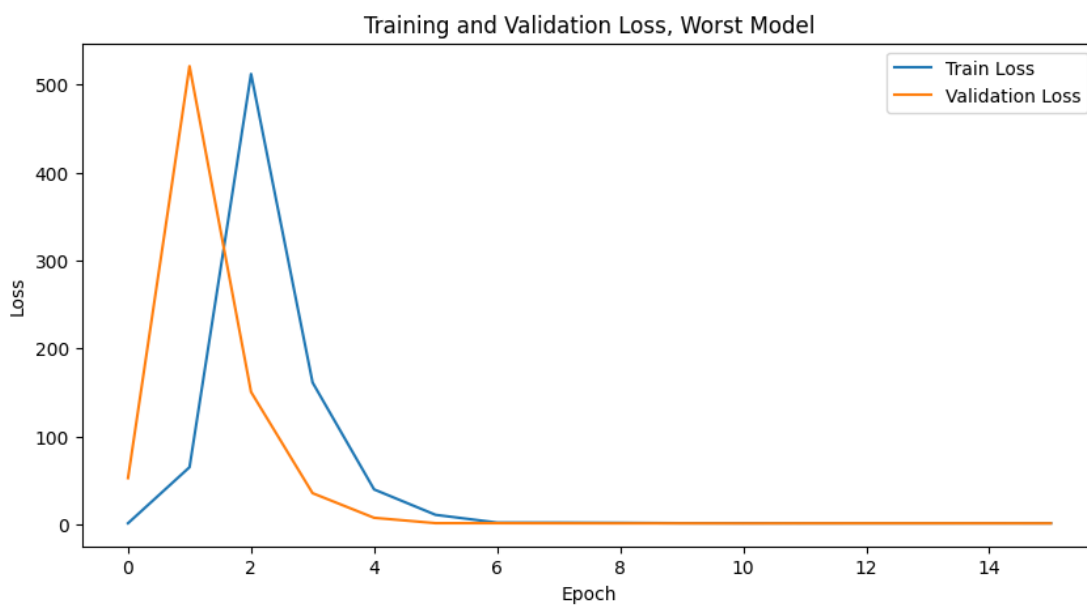
Grafy a konfúzna matica:

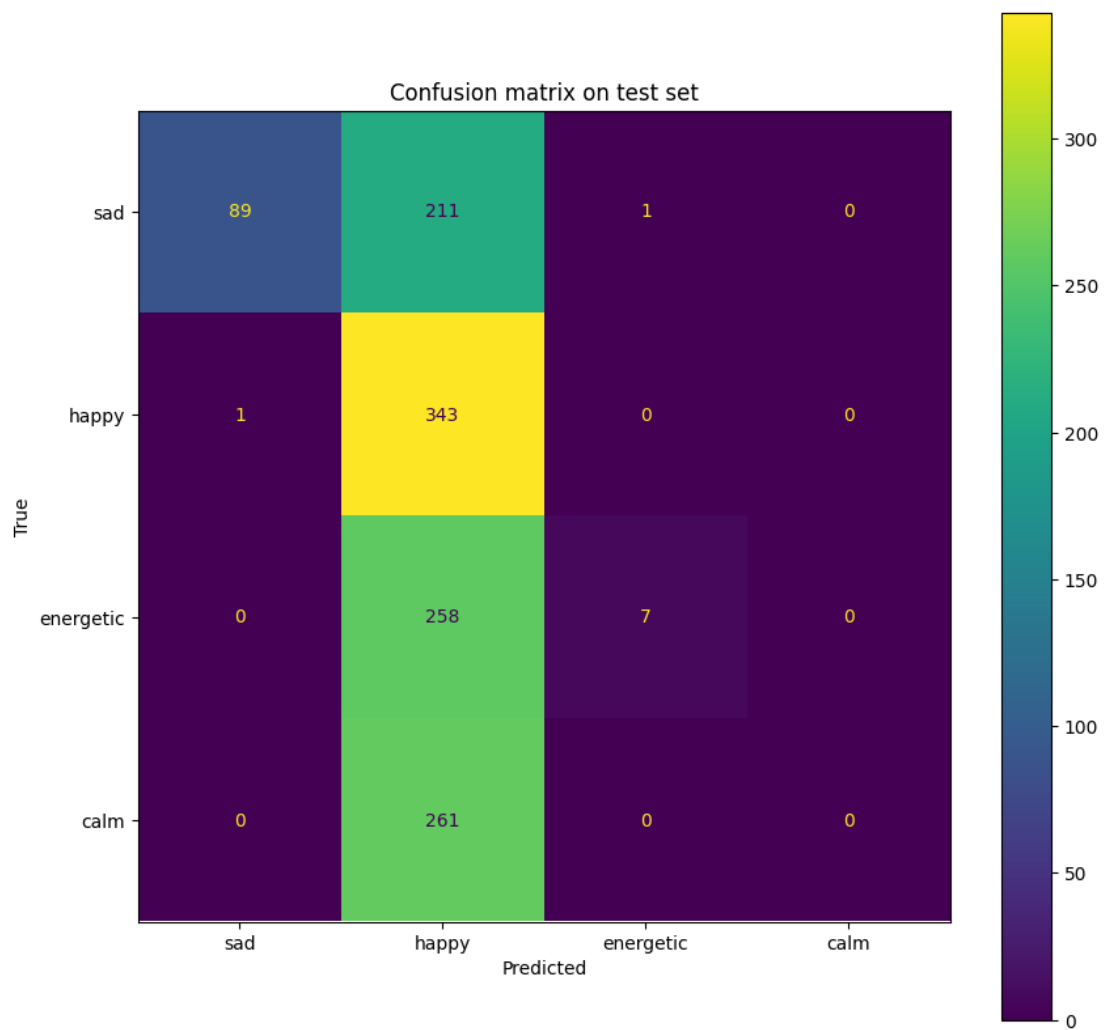




Úspešnosť vyhodnotenia na najhorší model je 36.29% na tréningových dát a 37.49% na testovacích dát.

Grafy a konfúzna matica:





b. Klasifikácia popularity

Ako tému našej problematiky sme si vybrali klasifikáciu obľúbenosti piesní z datasetu Spotify. Prvým krokom bolo adaptovať stĺpec "popularity" tak, aby reprezentoval kategorické hodnoty. Riešenie spočívalo v definovaní piatich kategórií: "Very Low", "Low", "Medium", "High" a "Very High". Hodnoty týchto kategórií boli pridelené na základe pôvodných numerických hodnôt s cieľom dosiahnuť vyvážené zastúpenie jednotlivých kategórií.

V nasledujúcej fáze spracovania dát sme aplikovali metódu One Hot Encoding, známu aj ako Dummy Encoding, na atribút 'emotion'. Po tejto transformácii sme získali nové stĺpce, ktoré budeme využívať v procese trénovania modelu.

#	Column	Non-Null	Count	Dtype
0	danceability	11825	non-null	float64
1	energy	11825	non-null	float64
2	loudness	11825	non-null	float64
3	speechiness	11825	non-null	float64
4	acousticness	11825	non-null	float64
5	instrumentalness	11825	non-null	float64
6	liveness	11825	non-null	float64
7	valence	11825	non-null	float64
8	tempo	11825	non-null	float64
9	duration_ms	11825	non-null	float64
10	popularity	11825	non-null	int64
11	number_of_artists	11825	non-null	float64
12	explicit	11825	non-null	int64
13	genre_ambient	11825	non-null	int64
14	genre_anime	11825	non-null	int64
15	genre_bluegrass	11825	non-null	int64
16	genre_blues	11825	non-null	int64
17	genre_classical	11825	non-null	int64
18	genre_comedy	11825	non-null	int64
19	genre_country	11825	non-null	int64
20	genre_dancehall	11825	non-null	int64
21	genre_disco	11825	non-null	int64
22	genre_edm	11825	non-null	int64
23	genre_emo	11825	non-null	int64
24	genre_folk	11825	non-null	int64
25	genre_forro	11825	non-null	int64
26	genre_funk	11825	non-null	int64
27	genre_grunge	11825	non-null	int64
28	genre_hardcore	11825	non-null	int64
29	genre_house	11825	non-null	int64
30	genre_industrial	11825	non-null	int64
31	genre_j-pop	11825	non-null	int64
32	genre_j-rock	11825	non-null	int64
33	genre_jazz	11825	non-null	int64
34	genre_metal	11825	non-null	int64
35	genre_metalcore	11825	non-null	int64
36	genre_opera	11825	non-null	int64
37	genre_pop	11825	non-null	int64
38	genre_punk	11825	non-null	int64
39	genre_reggaeton	11825	non-null	int64
40	genre_rock	11825	non-null	int64
41	genre_rockabilly	11825	non-null	int64
42	genre_ska	11825	non-null	int64
43	genre_sleep	11825	non-null	int64
44	genre_soul	11825	non-null	int64
45	emotion_calm	11825	non-null	int64
46	emotion_energetic	11825	non-null	int64
47	emotion_happy	11825	non-null	int64
48	emotion_sad	11825	non-null	int64

Využil som pôvodnú neurónovú sieť, ktorú som predtým konštruoval. Aby som identifikoval optimálny model, aplikoval som metódu Grid Search s rovnakými kombináciami hyperparametrov, aké boli použité v predchádzajúcej úlohe. Po dokončení tejto optimalizácie som identifikoval nasledujúci najlepší model s hyperparametrami:

Best configuration is: {'lr': 0.001, 'hidden_layer_sizes': [256, 128, 64], 'dropout_prob': 0.3, 'epochs': 500}

Úspešnosť vyhodnotenia na model je 70.34% na tréningových dát a 68.22% na testovacích dát.

Grafy a konfúzna matica:

