

# I-SUNS: Zadanie č.1

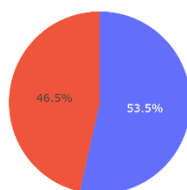
Neurónové siete

Tomáš Minárik

## Príprava dát

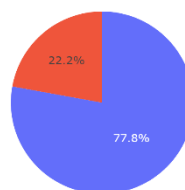
- Najprv sme si načítali dáta pomocou funkcie `read_csv`. Po načítaní dát sme ako prvé skontrolovali riadky kde sa nachádzajú nulové hodnoty. Zistili sme že aspoň jedna nulová hodnota sa nachádza v 15 841 z 62 869 riadkov.
- Rozhodli sme sa tieto riadky odstrániť keďže neskôr použijeme undersampling.
- Ďalej sme sa rozhodli odstrániť stĺpce s id a menami - `D_appid`, `D_name`, `D_developer`, `D_publisher`.
- Z hodnotení sme sa rozhodli ponechať len stĺpec `score` a odstrániť stĺpce `positive`, `negative`, `D_reviews`
- Zo stĺpca `D_owners` sme dostali počet majiteľov tak že sme spravili priemer z hornej a dolnej hranice.
- Zo stĺpca `D_release_date` sme vybrali rok.
- Stĺpec `D_genre` sme odstránili keďže žánre budeme ďalej vyberať z tagov.
- Ďalej sme sa pozreli na stĺpce ktoré obsahovali hodnoty `True` a `False`.
- Odstránili sme stĺpec `coming_soon` lebo obsahoval len hodnoty `False`.

has\_website\_linked



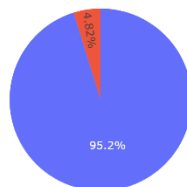
has\_controller\_support

■ true  
■ false



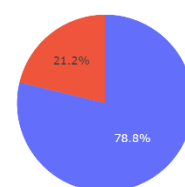
■ false  
■ true

is\_single\_player



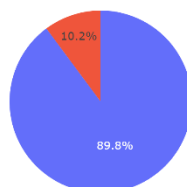
is\_multi\_player

■ true  
■ false



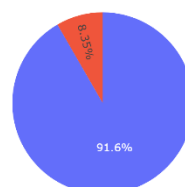
■ false  
■ true

is\_early\_access



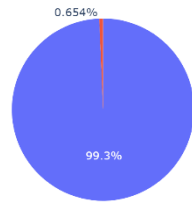
mature\_content

■ false  
■ true



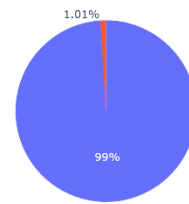
■ false  
■ true

Addictive



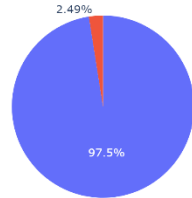
false  
true

Beautiful



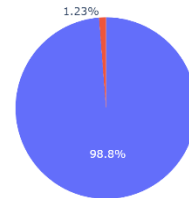
false  
true

Classic



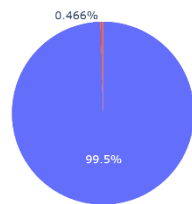
false  
true

Competitive



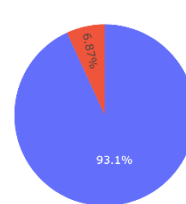
false  
true

Cult Classic



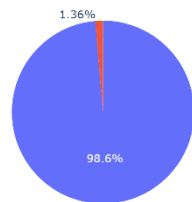
false  
true

Difficult



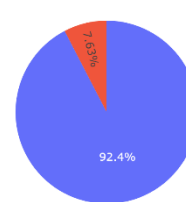
false  
true

Emotional



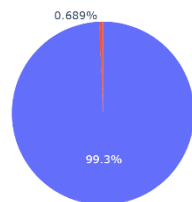
false  
true

Funny



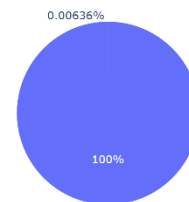
false  
true

Lore-Rich



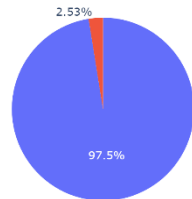
false  
true

Masterpiece



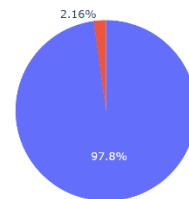
false  
true

Replay Value

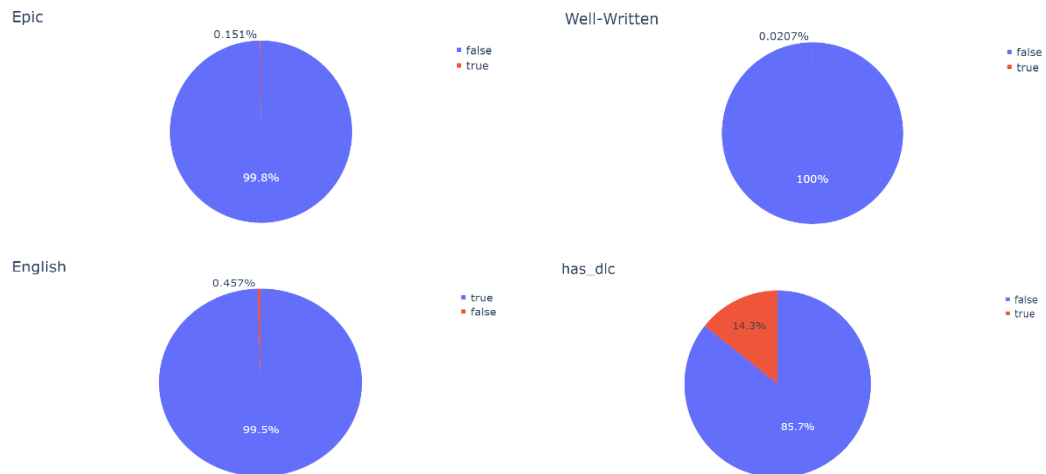


false  
true

Short



false  
true



- Z grafov môžeme vidieť že väčšina stĺpcov obsahuje viac ako 90% len hodnotu False a preto sme sa rozhodli ponechať len stĺpce has\_website\_linked, is\_multi\_player, has\_controller\_support, has\_dlc.
- Ďalej sme zo stĺpca D\_tags získali tagy, ktoré sa vyskytujú najčastejšie.

Tags	
Indie	32319
Singleplayer	21884
Action	21703
Casual	20747
Adventure	20158
2D	11606
Strategy	10290
Simulation	10003
RPG	8693
Puzzle	8134
Atmospheric	7586
Early Access	6366
Story Rich	6038
Pixel Graphics	5990
Multiplayer	5774
Colorful	5310
Arcade	5288
3D	5263
Cute	5103
First-Person	4940

- Z vyššie uvedeného zoznamu sme vybrali 8 najčastejšie sa vyskytujúcich žánrov a pre každý sme vytvorili vlastný stĺpec (Indie, Singleplayer, Action, Casual, Adventure, 2D, Strategy, Simulation, RPG).
- Zo stĺpcov, ktoré informovali o vydavateľovi sme sa rozhodli ponechať stĺpce publisher\_est, self\_published a odstrániť stĺpec developer\_est.

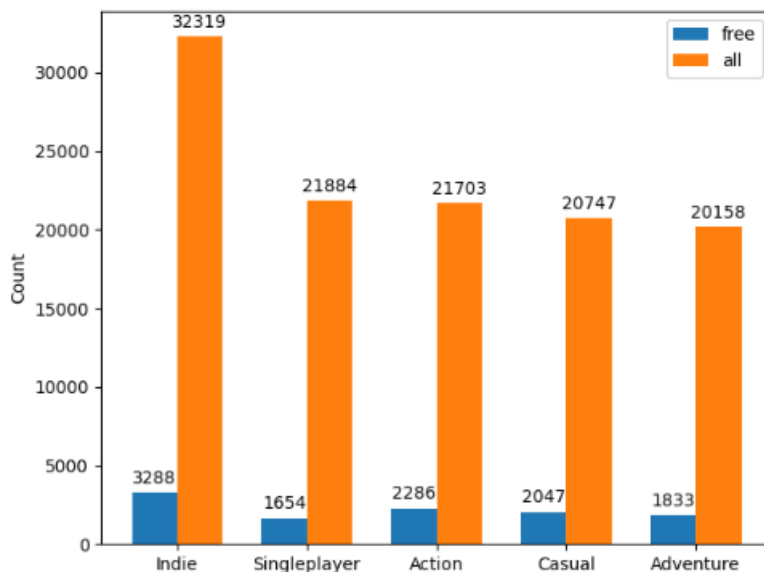
- Stĺpec ccu sme sa rozhodli odstrániť keďže máme stĺpec s počtom majiteľov a zároveň 99.97% nemá viac ako 100 hráčov.
- Na záver sme normalizovali stĺpce s číselnými hodnotami (D\_owners, D\_release\_date, languages, publisher\_est). V nasledujúcom obrázku som zobrazené priemery pred a po použití MinMaxScalera.

```
Before
Owners: 122628.41355860163
Date: 2018.2569121241102
languages: 3.6380405907980022
publishers: 17.478588885346934
After
Owners: 0.0015019124357727915
Date: 0.8502764849644099
languages: 0.09421573538564294
publishers: 0.04088979872294524
```

- Rovnaké úpravy boli spravené aj na testovacích dátach

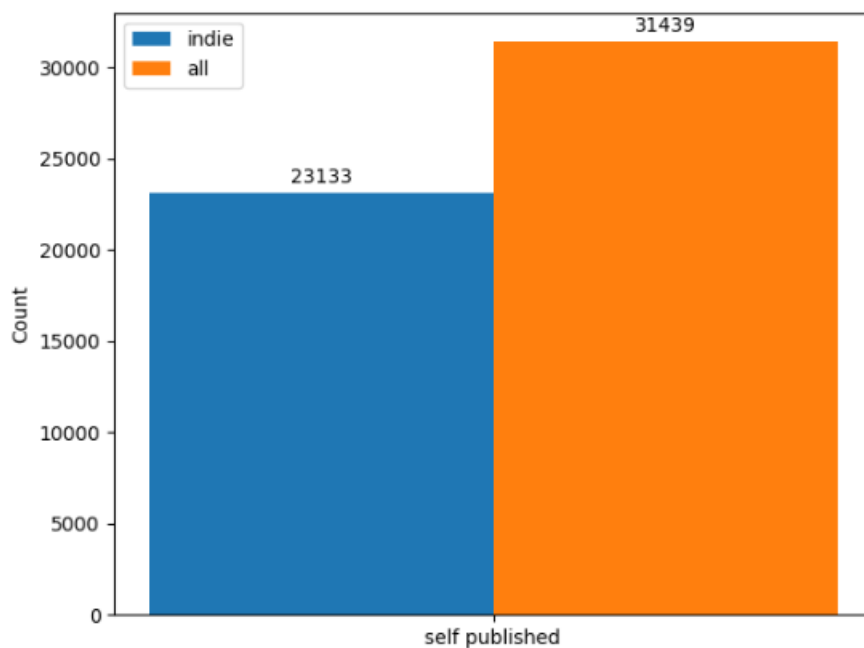
EDA

### Top žánre a is\_free



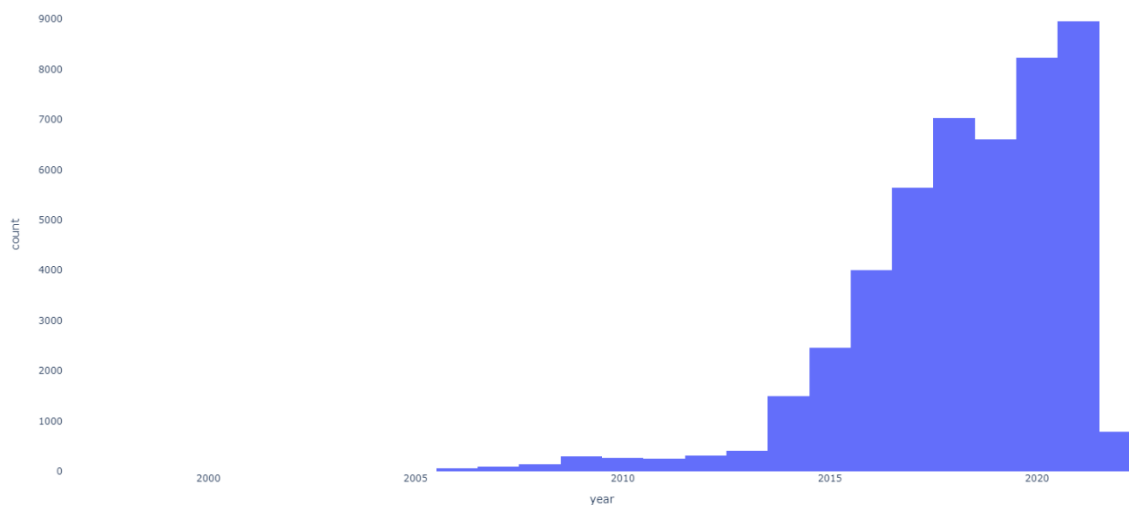
Z grafu môžeme vidieť že najlepší pomer hier zadarmo k plateným hrám má žáner action pri ktorom je to približne 10.5%. Najhoršie je na tom žáner singleplayer z ktorého je približne 7.6% hier zadarmo.

### Self published a Indie



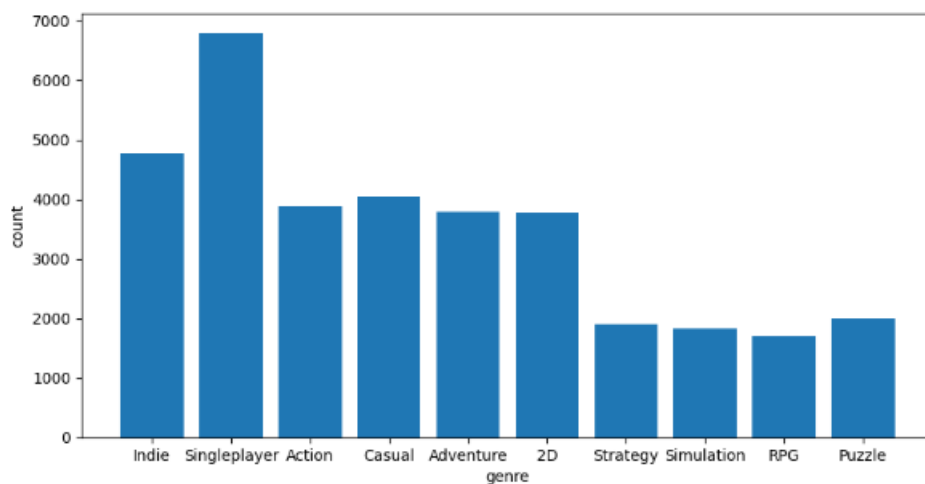
Z grafu môžeme vidieť že väčšina hier, ktoré sú self published patria pod žáner indie. Čo sa dalo predpokladať keďže indie znamená nezávislá a ide o hry tvorené jednotlivcom alebo malým tímom.

#### Počet vydaných hier podľa rokov



Z grafu môžeme vidieť že takmer každý rok je vydaných viac hier ako v predchádzajúcom roku. Zatiaľ najviac vydaných hier bolo v roku 2021 a to 8954.

#### Top herné žánre v roku 2021

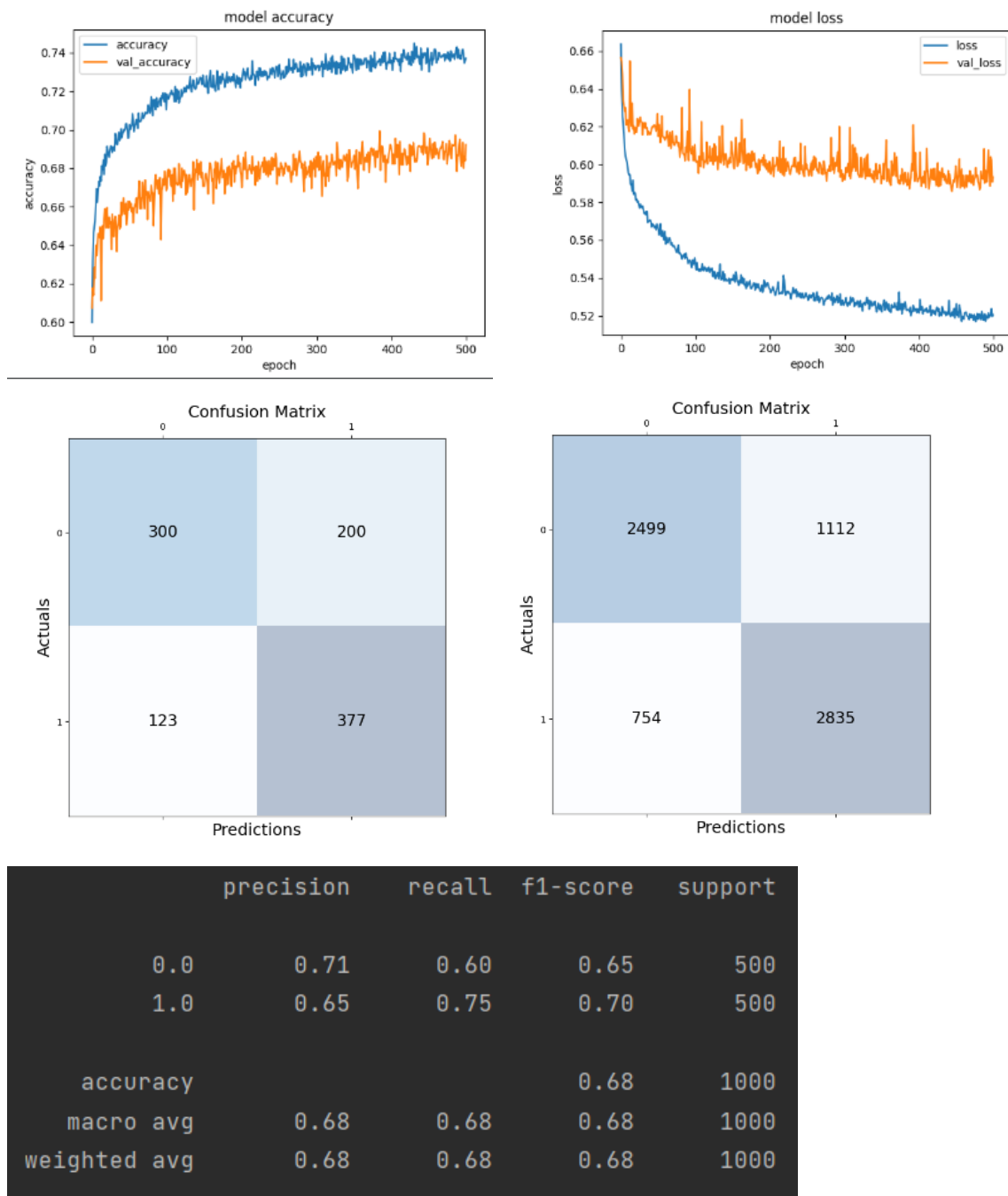


V roku 2021 bolo najviac hier vzdaných ako singleplayer na druhom mieste boli indie a na treťom casual.

## Trénovanie

### Sieť 1

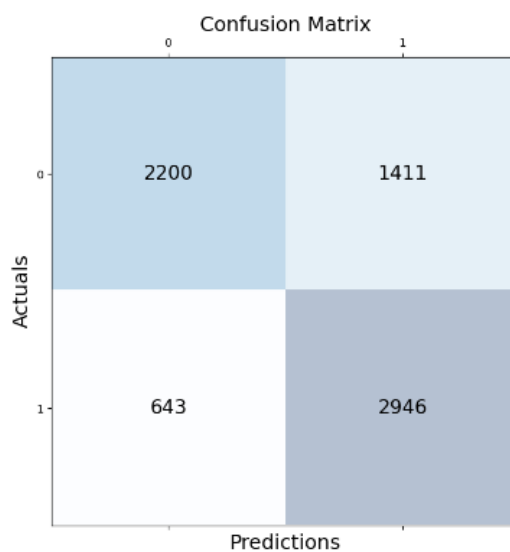
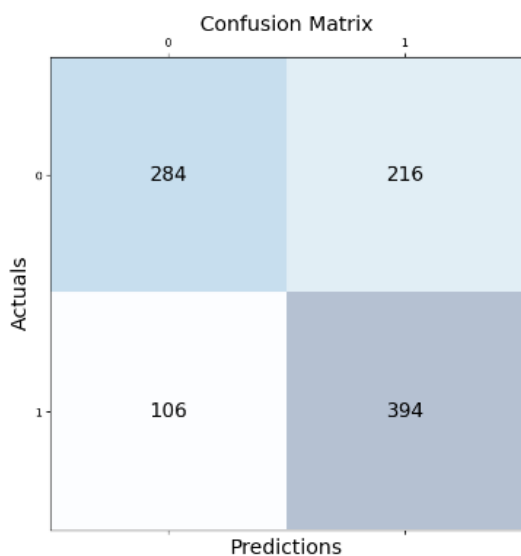
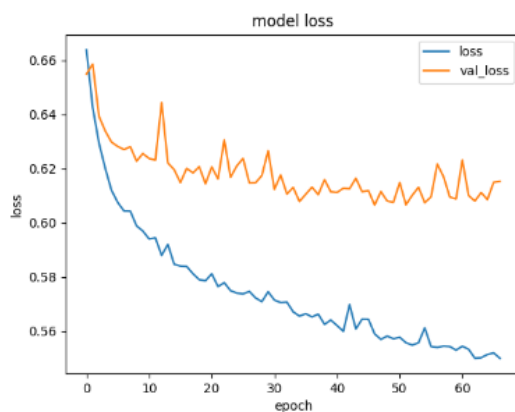
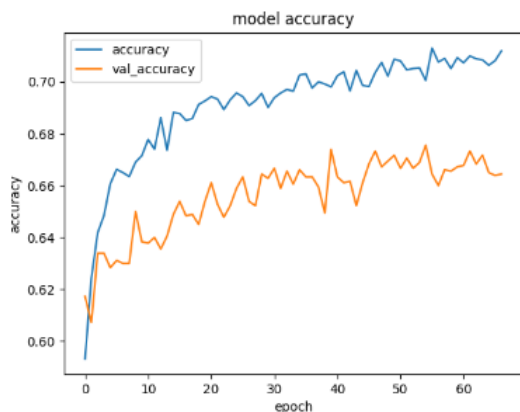
- Vstupné parametre(20) -  
score,D\_owners,D\_release\_date,languages,publisher\_est,self\_published,has\_dlc,has\_website\_linked,has\_controller\_support,is\_multi\_player,is\_early\_access,Indie,2D,Action,Casual,Adventure,Singleplayer,Strategy,Simulation,RPG
- Rozdelenie:
  - 4500 zadarmo
  - 4500 platené
- Validáčné dáta – 20% z trénovacích
- 1 skrytá vrstva s 20 neurónmi a aktivačnou funkciou relu
- 1 výstupný neurón, aktivačná funkcia sigmoid
- Kriteiálna funkcia Binary Cross Entropy
- Solver – Adam s learning rate 0.01
- Batch size = 100
- Počet epoch = 500



- Úspešnosť na tréningových dátach je 73% a validačných je 69%
- Úspešnosť na testovacích dátach je 68%.

## Sieť 2

- Rovnaká sieť ako sieť 1
- Pridaný EarlyStopping s patience 20



	precision	recall	f1-score	support
0.0	0.73	0.57	0.64	500
1.0	0.65	0.79	0.71	500
accuracy			0.68	1000
macro avg	0.69	0.68	0.67	1000
weighted avg	0.69	0.68	0.67	1000

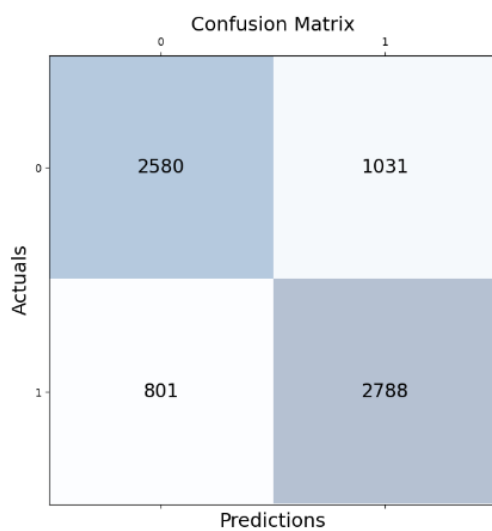
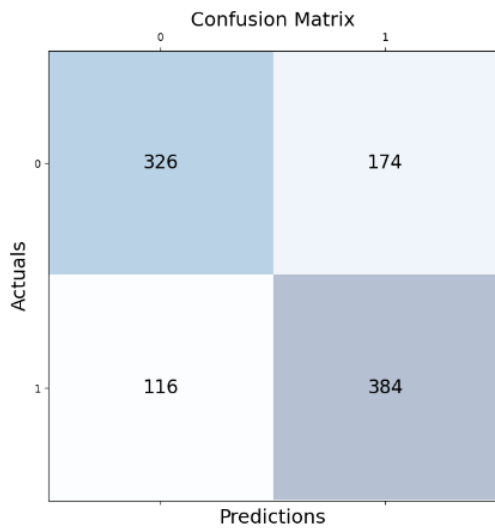
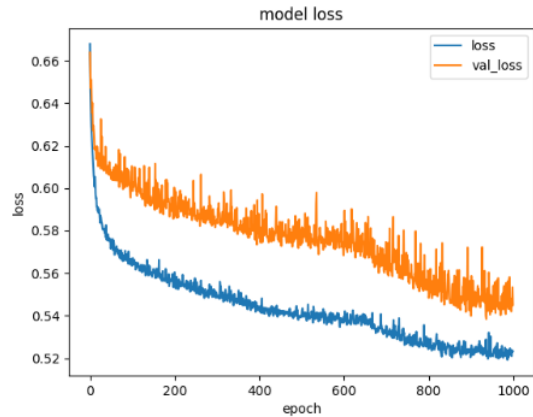
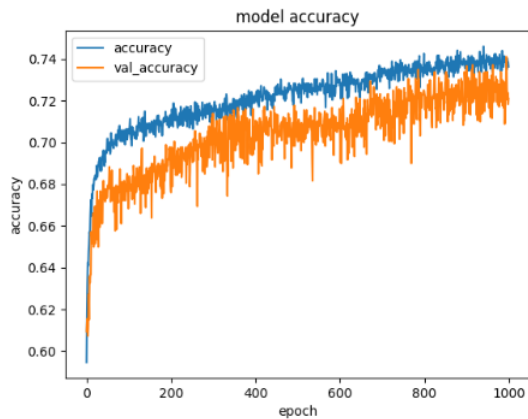
- Úspešnosť na tréningových dátach je 71% a na validačných 66%.
- Úspešnosť na testovacích dátach je 69%.
- Trénovanie bolo ukončené po 67. epoche.

### Sieť 3

- Vstupné parametre(16) - score,D\_owners,D\_release\_date,is\_free,languages,publisher\_est,self\_published,has\_dlc,has\_website\_linked,has\_controller\_support,is\_multi\_player,is\_early\_access,Indie>Action,Casual,Adventure,Singleplayer
- Použitie len Top 5 žánrov namiesto 9
- 1 skrytá vrstva s 16 neurónmi a aktivačnou funkciou relu



- Počet epoch = 1000
- EarlyStopping s patience=100



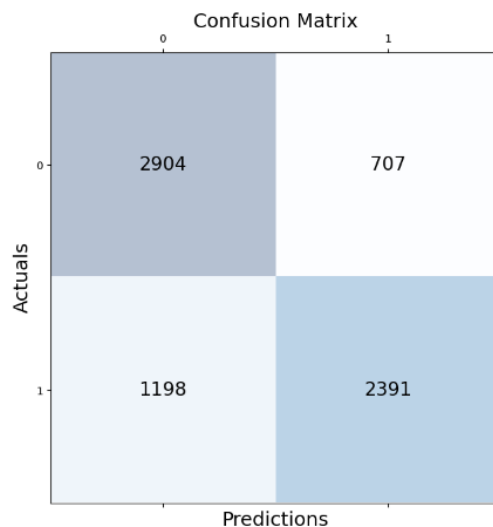
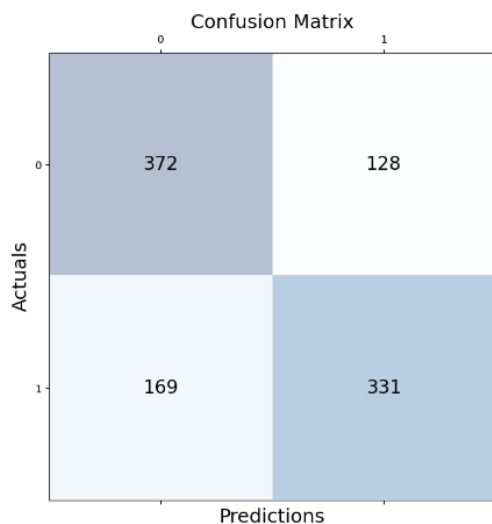
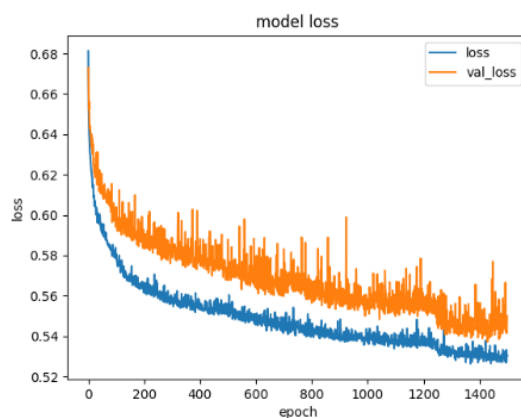
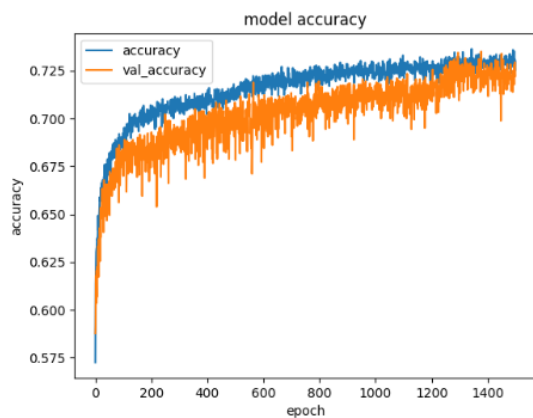
	precision	recall	f1-score	support
0.0	0.74	0.65	0.69	500
1.0	0.69	0.77	0.73	500
accuracy			0.71	1000
macro avg	0.71	0.71	0.71	1000
weighted avg	0.71	0.71	0.71	1000

- Trénovanie ukončené po 1000. epoche.
- Úspešnosť na tréningových dátach je 74% a na validačných 72%
- Úspešnosť na testovacích dátach je 71%

#### Sieť 4

- Vstupné parametre(14) -  
score,D\_owners,D\_release\_date,is\_free,languages,publisher\_est,self\_published,has\_dlc,has\_website\_linked,is\_multi\_player,Indie,Action,Casual,Adventure,Singleplayer
- Odstránili sme is\_early\_access a has\_controller\_support

- 1 skrytá vrstva s 14 neurónmi a aktivačnou funkciou relu
- Počet epoch = 1500
- EarlyStopping s patience=100



	precision	recall	f1-score	support
0.0	0.69	0.74	0.71	500
1.0	0.72	0.66	0.69	500
accuracy			0.70	1000
macro avg	0.70	0.70	0.70	1000
weighted avg	0.70	0.70	0.70	1000

- Trénovanie ukončené po 1500. epoche.
- Úspešnosť na tréningových dátach je 73% a na validačných 72%.
- Úspešnosť na testovacích dátach je 70%.

## Záver

Dosiahli sme výsledky v rozsahu 68-71%. Pričom najlepšia sieť bola sieť 3. Mala 16 vstupných parametrov a 1 skrytú vrstvu so 16 neurónmi. Pri prvej sieti bol rozdiel medzi tréningovými

a testovacími dátami 5%. Po přidání early-stoppingu v druhé síti jsme dosáhli rozdílu 2%, který byl nejmenší ze všech sítí.