

**SLOVENSKÁ TECHNICKÁ UNIVERZITA V BRATISLAVE**  
**FAKULTA ELEKTROTECHNIKY A INFORMATIKY**

AIS ID: 92320

**NEURÓNOVÉ SIETE - KATEGORIZÁCIA POČASIA**  
**ZADANIE**

Predmet: I-SUNS – Strojové učenie a neurónové siete  
Prednášajúci: prof. Dr. Ing. Miloš Oravec  
Cvičiaci: Ing. Marián Šebeňa

**Bratislava 2024**

**Bc. Lukáš Patrnčíak**

# Obsah

<b>Úvod</b>	<b>1</b>
<b>1 Vstupné dáta</b>	<b>2</b>
1.1 Predspracovanie dát . . . . .	2
<b>2 Analýza dát</b>	<b>6</b>
<b>3 Neurónová sieť</b>	<b>12</b>
3.1 Architektúra siete s pretrénovaním . . . . .	12
3.2 Architektúra siete bez pretrénovania . . . . .	15
3.3 Experimenty s trénovaním siete . . . . .	17
3.3.1 Dropout vrstvy . . . . .	22
<b>Záver</b>	<b>25</b>

# Zoznam obrázkov a tabuliek

Obrázok 1	Zobrazenie konfúznej matice pre číselné hodnoty v dátovej množine, po odstránení nulových hodnôt a outlierov . . . . .	5
Obrázok 2	Zobrazenie korelačnej matice pre EDA . . . . .	6
Obrázok 3	Zobrazenie grafu (scatterplot) pre pravdepodobnosť zrážok a vlhkosť . . . . .	7
Obrázok 4	Zobrazenie grafu (kdetplot) pre teplotnú hustotu vzhľadom na ročné obdobie . . . . .	8
Obrázok 5	Zobrazenie grafu (countplot) pre počet typov počasia podľa UV indexu . . . . .	9
Obrázok 6	Zobrazenie grafu (lineplot) pre zrážky podľa viditeľnosti a polohy	10
Obrázok 7	Zobrazenie grafu (violinplot) pre oblačnosť a atmosférický tlak .	11
Obrázok 8	Graf trénovacej a validačnej straty . . . . .	13
Obrázok 9	Graf trénovacej a validačnej úspešnosti . . . . .	13
Obrázok 10	Zobrazenie konfúznej matice pre trénovaciu úspešnosť . . . . .	14
Obrázok 11	Zobrazenie konfúznej matice pre testovaciu úspešnosť . . . . .	14
Obrázok 12	Graf trénovacej a validačnej straty s použitím Early Stopping .	15
Obrázok 13	Graf trénovacej a validačnej úspešnosti s použitím Early Stopping	15
Obrázok 14	Zobrazenie konfúznej matice pre trénovaciu úspešnosť s použitím Early Stopping . . . . .	16
Obrázok 15	Zobrazenie konfúznej matice pre testovaciu úspešnosť s použitím Early Stopping . . . . .	16
Obrázok 16	Graf trénovacej a validačnej straty pre najlepší experiment . . .	18
Obrázok 17	Graf trénovacej a validačnej úspešnosti pre najlepší experiment .	18
Obrázok 18	Konfúzna matica trénovacej množiny pre najlepší experiment . .	19
Obrázok 19	Konfúzna matica testovacej množiny pre najlepší experiment . .	19
Obrázok 20	Graf trénovacej a validačnej straty pre najhorší experiment . . .	20
Obrázok 21	Graf trénovacej a validačnej úspešnosti pre najhorší experiment	20
Obrázok 22	Konfúzna matica trénovacej množiny pre najhorší experiment . .	21
Obrázok 23	Konfúzna matica testovacej množiny pre najhorší experiment . .	21
Obrázok 24	Graf trénovacej a validačnej straty pre neurónovú sieť s dropout vrstvami . . . . .	22

Obrázok 25	Graf trénovacej a validačnej úspešnosti pre neurónovú sieť s droupout vrstvami . . . . .	23
Obrázok 26	Zobrazenie konfúznej matice pre trénovacu úspešnosť pre neurónovú sieť s droupout vrstvami . . . . .	23
Obrázok 27	Zobrazenie konfúznej matice pre testovacu úspešnosť pre neurónovú sieť s droupout vrstvami . . . . .	24
Tabuľka 2	Tabuľka vykonaných experimentov na základe pomeru trénovacích a testovacích dát s výslednou úspešnosťou . . . . .	4
Tabuľka 3	Konfigurácia jednotlivých experimentov a presnosť trénovania a testovania . . . . .	17

# Zoznam skratiek

<b>EDA</b>	Exploratory Data Analysis
<b>UV</b>	Ultrafialové žiarenie
<b>X</b>	Vstupná množina
<b>Y</b>	Výstupná množina

# Úvod

Cieľom tohto zadania je natréňovať neurónovú sieť pre kategorizáciu počasia, pričom pracujeme s dátami z AIS. Najprv bolo tieto dáta potrebné načítať z príslušného súboru a následne boli upravené tak, aby neobsahovali žiadne chýbajúce alebo neobvyklé hodnoty (outliery). Ak daná množina (resp. stĺpec v tabuľke dát) obsahoval príliš veľa chýbajúcich hodnôt (viac, ako 25 % dát chýbalo), bol tento stĺpec odstránený. Dáta boli následne vhodne zakódované, škálované a rozdelené do množín, aby bolo možné natréňovať jednoduchý klasifikačný model. Boli vykonané 4 experimenty, pričom pre najlepší z nich bola vytvorená tabuľka a konfúzna matica.

Následne boli dáta analyzované pomocou Exploratory Data Analysis (EDA), pričom boli použité dáta už upravené dáta, očistené od outlierov, duplícít a chýbajúcich hodnôt. Následne bolo vykreslených a popísaných 5 rôznych grafov, na základe korelačnej matice, ktorá ukázala, ako jednotlivé hodnoty medzi sebou vzájomne korelujú (zistoval sa ich vzájomný vzťah, ako spolu dané hodnoty súvisia).

Následne bola natréňovaná neurónová sieť, pričom bolo vykonaných 5 pokusov, kde boli zmenené rôzne parametre, aby sme zistili rozdiely v jednotlivých experimentoch pri rôznych hodnotách. Konfigurácie týchto experimentov boli zapísané v jednej tabuľke, pričom pre najlepšie a najhoršie tréňovanie bola vykreslená konfúzna matica.

# 1 Vstupné dáta

V tejto kapitole bude popísaný spôsob predspracovania a následného vyhodnotenia vstupných dát zo súboru *weather\_data.csv*. Pre spracovanie, tréňovanie a vykreslenie údajov boli použité knižnice:

- NumPy<sup>1</sup>: Pre matematické výpočty.
- Pandas<sup>2</sup>: Manipulácia s tabuľkovými dátami a ich analýza.
- Seaborn<sup>3</sup>: Vysokoúrovňová vizualizácia dát so zameraním na štatistiku.
- Matplotlib.pyplot<sup>4</sup>: Základná knižnica na tvorbu grafov s nízkoúrovňovou kontrolou nad vizualizáciami.
- Tensorflow (Keras)<sup>5</sup>: Tvorba a tréňovanie modelov strojového učenia a neurónových sietí. API Keras zjednodušuje tvorbu a tréňovanie neurónových sietí.
- sklearn<sup>6</sup>: Knižnica určená pre klasické metódy strojového učenia. Obsahuje množstvo algoritmov pre klasifikáciu, regresiu,...

## 1.1 Predspracovanie dát

Po načítaní dát z daného súboru nasledoval ich výpis, aby sme zistili, rôzne informácie o nich, predovšetkým koľko duplicit a chýbajúcich vzoriek obsahujú a aký je ich celkový počet. Najprv boli odstránené množiny<sup>7</sup>, v ktorých chýbalo viac ako 25% dát (v tomto prípade to bola množina *Irradiance*). Potom sme rozdelili dáta do dvoch skupín - tie, ktorý by mali obsahovať len číselné hodnoty a tie, ktoré by mali obsahovať len textové hodnoty.

Množiny dát, ktoré mali obsahovať len číselné hodnoty, boli prekontrolované funkciou *median\_replace()*, ktorá v prípade výskytu napr. textových hodnôt túto hodnotu zmazala a nahradila mediánom všetkých hodnôt v danej množine. Pre odstránenie outlierov (nezvyčajných hodnôt) bola navrhnutá funkcia *outliers()*, ktorá vrátila novú dátovú

---

<sup>1</sup><https://numpy.org/>

<sup>2</sup><https://pandas.pydata.org/>

<sup>3</sup><https://seaborn.pydata.org/>

<sup>4</sup>[https://matplotlib.org/stable/api/ pyplot\\_summary.html](https://matplotlib.org/stable/api/ pyplot_summary.html)

<sup>5</sup><https://www.tensorflow.org/guide/keras>

<sup>6</sup><https://scikit-learn.org/stable/>

<sup>7</sup>V tabuľke dát sú to stĺpce

množinu bez outlierov<sup>8</sup> a následne boli odstránené duplicitné hodnoty. Proces očistenia dát bol zakončený kontrolou, ktorá vrátila stĺpce (množiny) dát podľa nasledujúcich intervalov:

- Temperature:  $-25^{\circ}C$  do  $109^{\circ}C$
- Humidity: od 20% do 109%
- Wind Speed: 0 km/h do 48,5 km/h
- Precipitation (%): od 0% do 109%
- Cloud Cover: partly cloudy, clear, overcast, cloudy
- Atmospheric Perssure: od 984 hPa do 1067 hPa
- UV Index: od 0 do 14
- Season: Winter, Spring, Summer, Autumn
- Visibility (km): od 0 km do 20 km
- Location: inland, mountain, coastal
- Weather Type: Rainy, Cloudy, Sunny, Snowy
- Irradiance: od 200 do 800 W/m (**stĺpec(množina) odstránená z dôvodu veľkého počtu výskytov prázdnych hodnôt**)

Kódovanie nečíselných hodnôt prebehlo pomocou *Label Encoding*, ktorý priraduje textovým hodnotám práve jedno číslo (napríklad množine dát "Weather Type", ktorá obsahuje hodnoty "Sunny, Cloudy, Rainy, Snowy") boli pridelené hodnoty ("0, 1, 2, 3"). Keďže účelom implementácie bola kategorizácia počasia, výstupná množina Y obsahovala práve jednu množinu dát - Weather Type. Vstupná množina X obsahuje zvyšné množiny dát.

Následne bola navrhnutá funkcia *train\_model()*, ktorá vykonáva rozdelenie vstupných a výstupných dát (trénovacia<sup>9</sup>, validačná<sup>10</sup> a testovacia<sup>11</sup> množina podľa zvoleného deliaceho pomeru). Následne je vykonaná normalizácia dát pomocou **MinMaxScaler**, ktorý škáluje, resp. normalizuje hodnoty vstupných dát ( $X$ ) do intervalu (napríklad  $< 0, 1 >$ ). Trénovanie modelu prebieha logistickou regresiou so zadaným počtom iterácií (v tomto prípade 200). Následne prebehne predikcia a vyhodnotenia dát, pričom výstupom tejto funkcie sú nasledovné veličiny:

- model logistickej regresie,
- presnosť modelu vzhľadom na testovacie dáta,

---

<sup>8</sup>Hodnoty, ktoré sa výrazne líšia od ostatných hodnôt v danej množine

<sup>9</sup>Model sa z týchto dát učí, ako predikovať cieľové hodnoty

<sup>10</sup>Slúži na optimalizáciu modelu počas tréningu (ladenie hyperparametrov - napr. rýchlosť učenia, počet vrstiev v neurónovej sieti, ...)

<sup>11</sup>Slúži na finálne vyhodnotenie modelu po tréningu a ladení hyperparametrov



- konfúzna matica.

Výstup tvoria taktiež aj normalizované dáta (tréningové, testovacie a validačné) a cieľové (výstupné) dáta.

Číslo experimentu	Deliaci pomer 1	Deliaci pomer 2	Úspešnosť
1	0.3	0.4	0.933941
2	0.3	0.1	0.936364
3	0.2	0.4	0.943052
4	0.1	0.2	0.963636

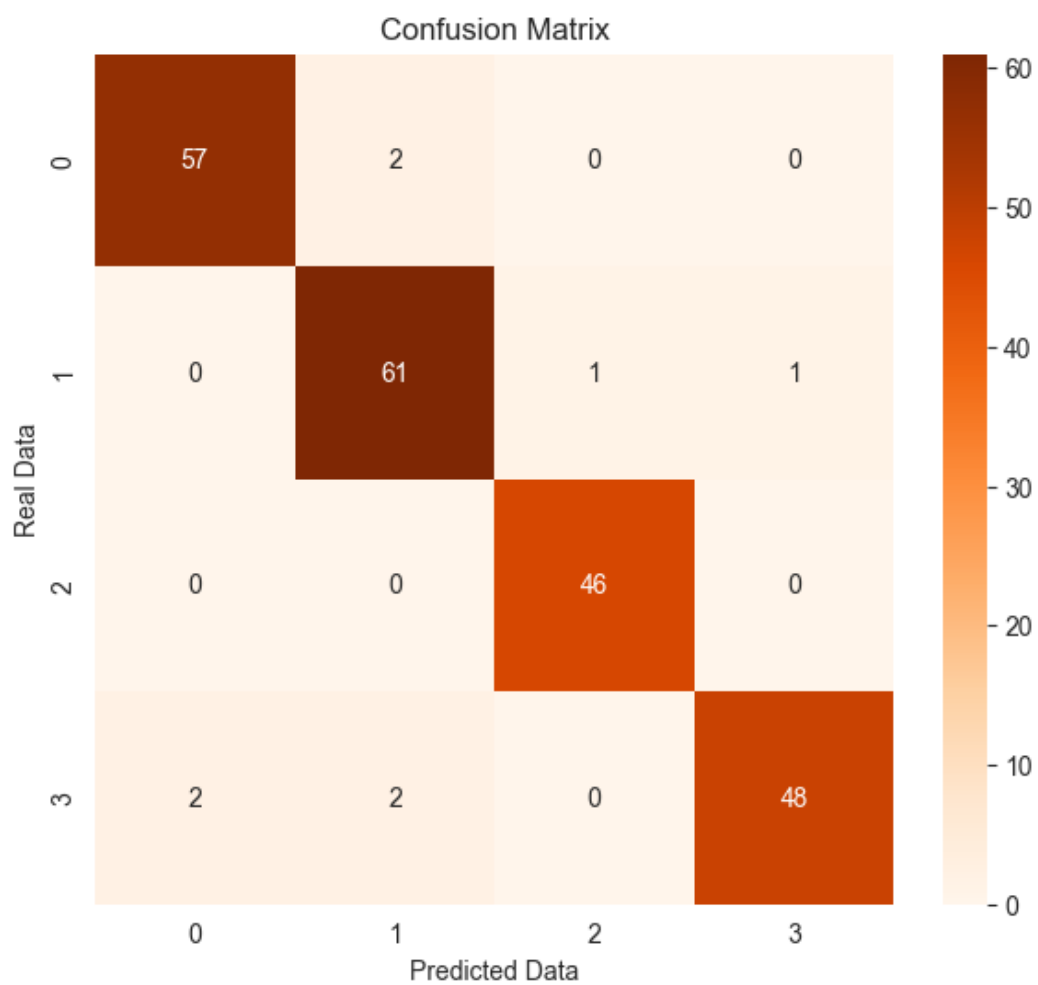
Tabuľka 2: Tabuľka vykonaných experimentov na základe pomeru tréningových a testovacích dát s výslednou úspešnosťou

Ako je možné z tabuľky vidieť, najlepší (najúspešnejší) experiment je ten, ktorého deliaci pomer<sup>12</sup> je experiment číslo 3. Táto presnosť bola vygenerovaná funkciou *accuracy\_score()* na základe skutočných hodnôt (ktoré obsahuje náš dataset) a predikovaných hodnôt (výstupné hodnoty - predpovede, vytvorené na základe vstupných dát).

Na obrázku nižšie je možné vidieť vygenerovanú konfúznú maticu<sup>13</sup>. Matica bola graficky vykreslená funkciou *generate\_confusion\_matrix()* na základe vygenerovanej matice *confusion\_matrix* vo funkcii *train\_model()*.

<sup>12</sup>Pomer tréningových a testovacích dát, zvolená hodnota uvádza, aká časť dát bude použitá ako testovacie dáta (napr. 0.1 – 10% pre testovacie dáta)

<sup>13</sup>Matica, ktorá znázorňuje výkonnosť klasifikačného modelu tým, že porovnáva predikované výsledky modelu s reálnymi (skutočnými) hodnotami. Hlavná diagonála obsahuje počet správne klasifikovaných predikcií (model správne identifikoval triedy)

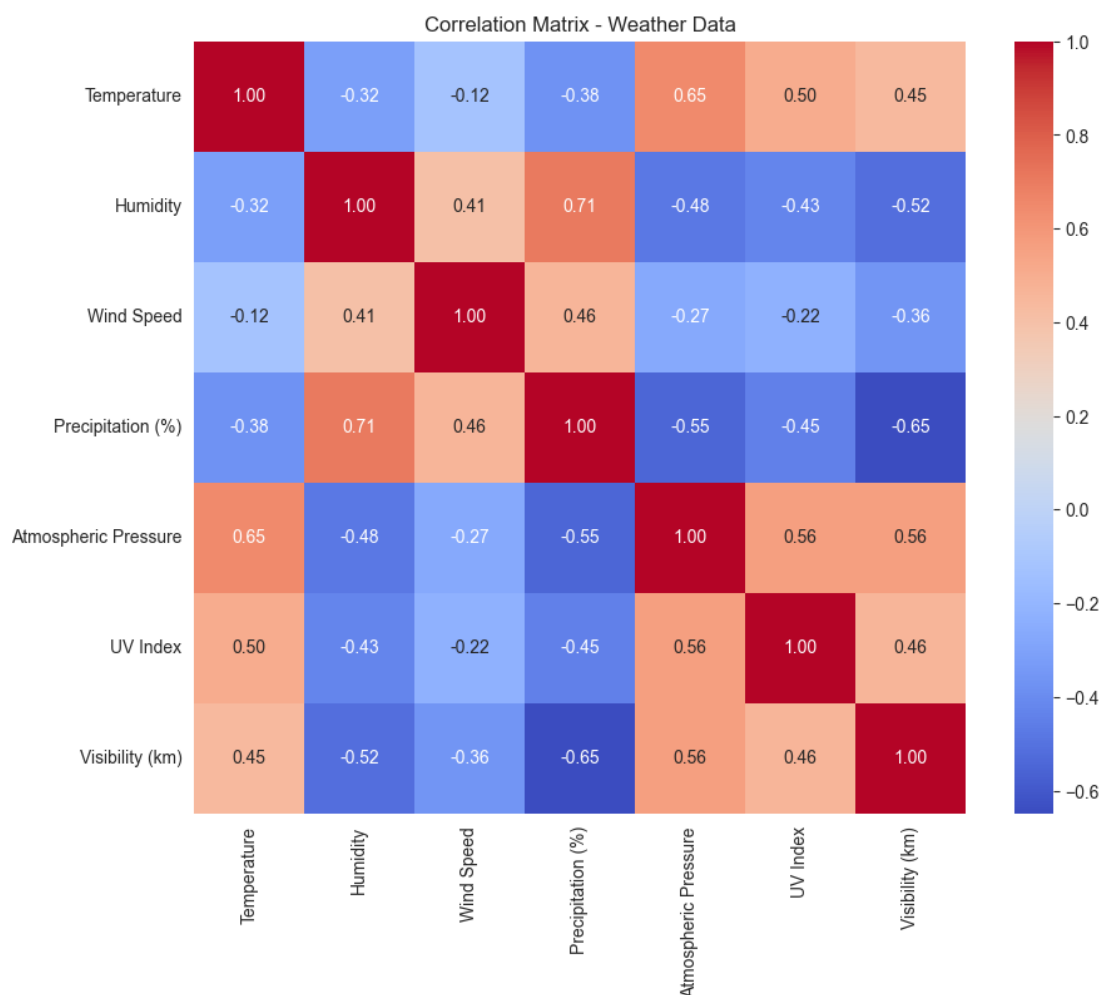


Obrázok 1: Zobrazenie konfúznej matice pre číselné hodnoty v dátovej množine, po odstránení nulových hodnôt a outlierov

## 2 Analýza dát

Výstupné dáta boli analyzované pomocou **Exploratory Data Analysis (EDA)**. V tomto prípade pracujeme s dátami (pred zakódovaním - aby bolo možné použiť slovné hodnoty), ktoré neobsahujú outliery alebo chýbajúce hodnoty. Zaujímavé hodnoty môžeme nájsť aj pomocou **korelačnej matice**<sup>14</sup>. Táto matica obsahuje zvyčajne hodnoty v intervale  $[-1, 1]$ :

- 1: Perfektná pozitívna korelácia (keď hodnota jednej premennej stúpa, druhá stúpa rovnako).
- -1: Perfektná negatívna korelácia (keď hodnota jednej premennej stúpa, druhá klesá).
- 0: Žiadna lineárna korelácia (premenné nie sú lineárne závislé).

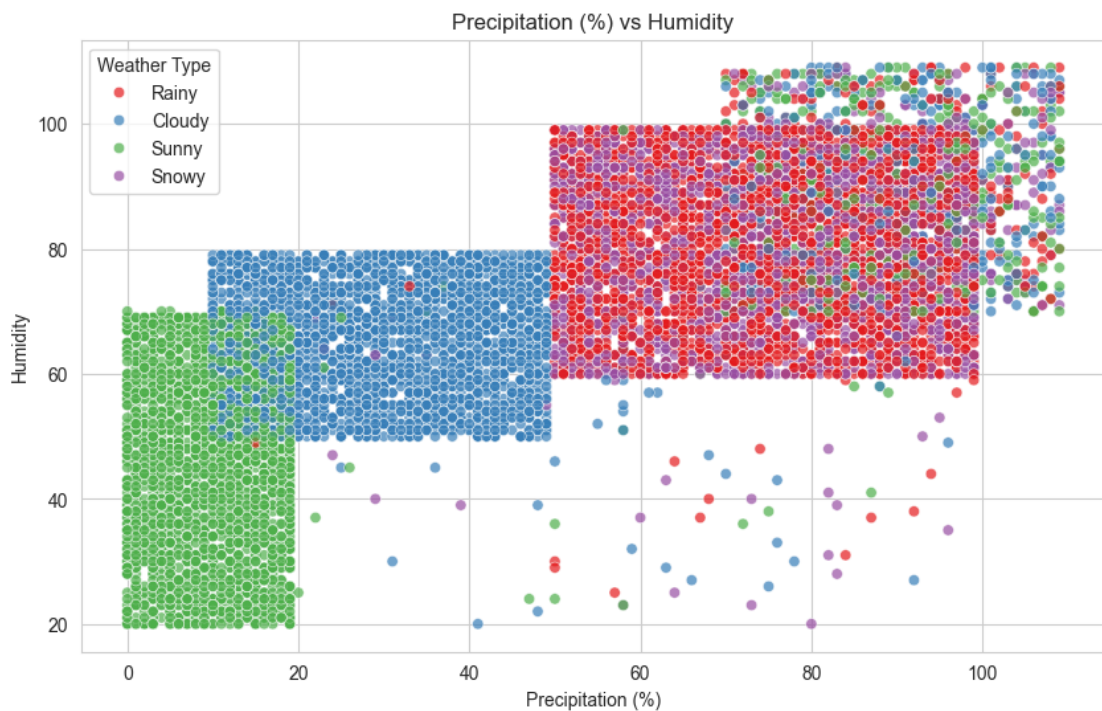


Obrázok 2: Zobrazenie korelačnej matice pre EDA

<sup>14</sup>Matica, ktorá ukazuje vzťahy medzi rôznymi premennými v dátových množinách

Táto matica bola vykreslená funkciou `generate_correlation_matrix()` pre množiny dát, v ktorých sa nachádzajú číselné hodnoty (keďže korelačná matica dokáže pracovať len s takými). Čím sa hodnota jednotlivých vzťahov približuje čo najbližšie ku 1 alebo  $-1$ , tým je korelácia lepšia.

1. **Pravdepodobnosť zrážok a vlhkosť:** Na tomto grafe znázorňuje os X pravdepodobnosť zrážok a os y percentuálny podiel vlhkosti v ovzduší. Každý bod v grafe predstavuje jedinečný záznam (riadok) v súbore údajov, kde je záznamom kombinácia percenta zrážok a vlhkosti pre daný typ počasia.

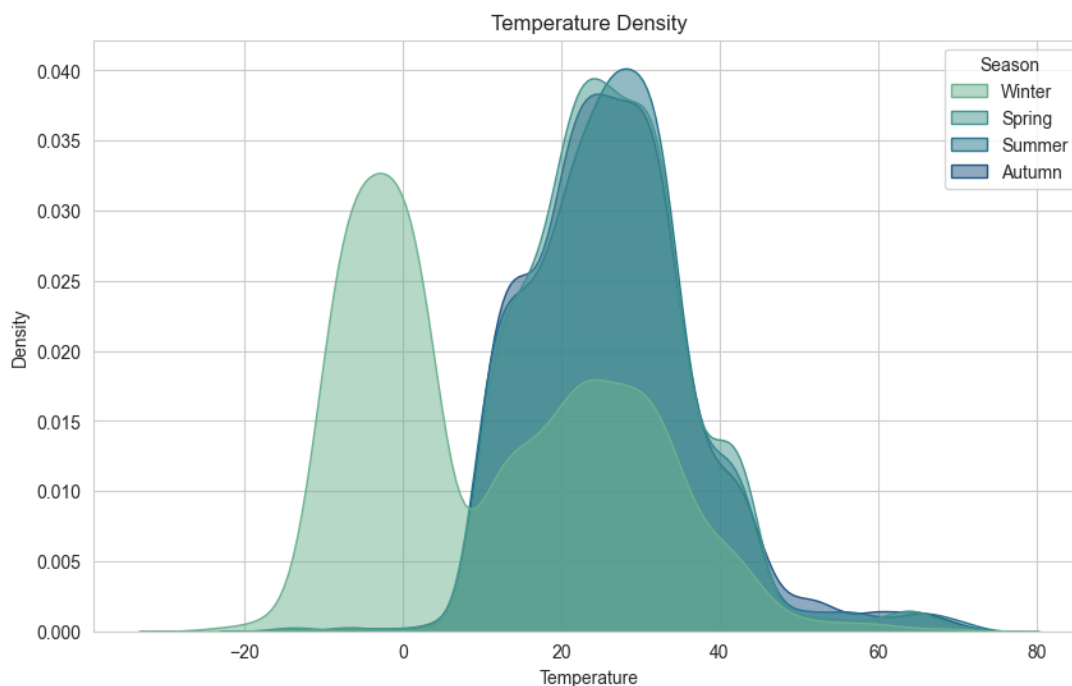


Obrázok 3: Zobrazenie grafu (scatterplot) pre pravdepodobnosť zrážok a vlhkosť

- Sunny: Tieto body môžu byť zvyčajne umiestnené v dolnej časti grafu, kde sú zrážky nízke (0-20 %) a vlhkosť tiež môže byť relatívne nízka.
- Rainy: Tieto body sa pravdepodobne nachádzajú v oblasti vyšších percentuálnych zrážok (prevažne 50-100 %) a môžu mať variabilnú vlhkosť, s niektorými bodmi na vyššej vlhkosti.
- Snowy: Podobne ako rainy, ale môžu mať nižšie teploty, a teda aj variabilitu vo vlhkosti.
- Cloudy: Tieto body môžu ukazovať stredné úrovne zrážok a vlhkosti, s možným výskytom nízkych až stredných hodnôt.

2. **Teplotná hustota vzhľadom na ročné obdobie:** Tento graf zobrazuje hustotu pravdepodobnosti teploty v závislosti od ročného obdobia. KDE plot ukazuje, ako je rozloženie teplôt rozšírené v rôznych obdobiach roka.

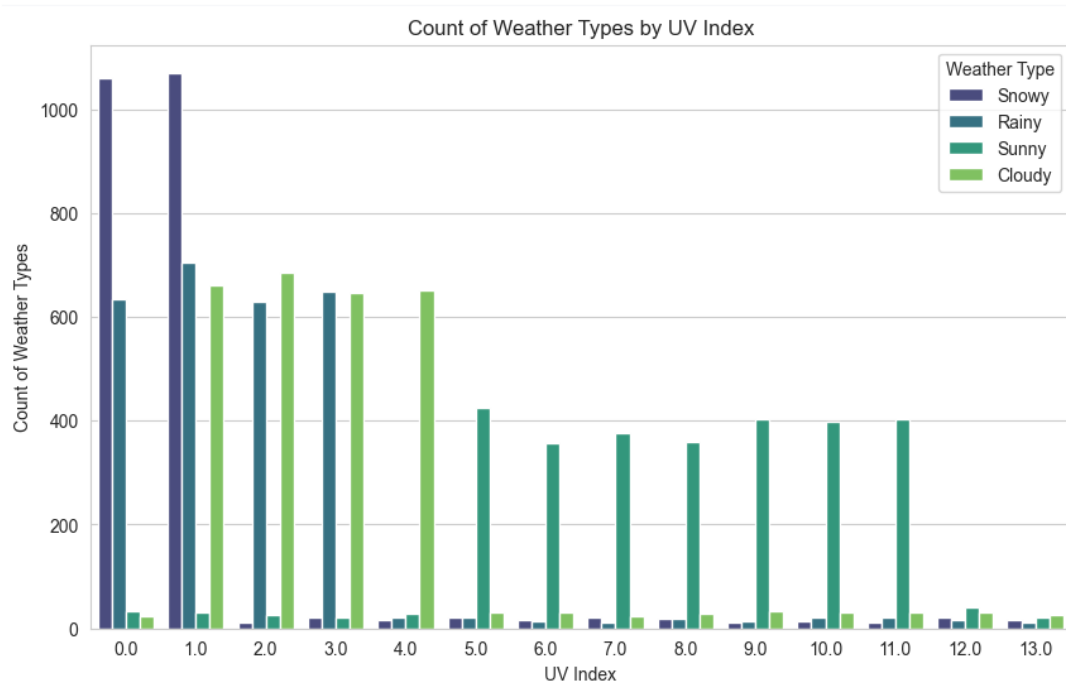
Vyššie oblasti na grafe naznačujú vyššiu hustotu pozorovaní, zatiaľ čo nižšie oblasti ukazujú nižšiu hustotu. Pomocou tohto grafu môžete identifikovať trendy, ako sú teplejšie alebo chladnejšie obdobia v priebehu roka.



Obrázok 4: Zobrazenie grafu (kdeplot) pre teplotnú hustotu vzhľadom na ročné obdobie

3. **Počet typov počasia podľa UV indexu:** Tento graf ukazuje, ako sa mení počet jednotlivých typov počasia s rastúcim UV indexom. Môžeme si všimnúť, že s vyššími hodnotami UV indexu (napríklad nad 6) sú častejšie výskyty slnečného počasia, zatiaľ čo pri nižších hodnotách je bežnejšie daždivé počasie. Os x reprezentuje hodnoty UV indexu. Na osi y je zobrazený počet výskytov každého typu počasia pre konkrétnu hodnotu UV indexu.

Tento počet indikuje, koľkokrát sa daný typ počasia objavil pri určitej hodnote UV indexu v množine dát.

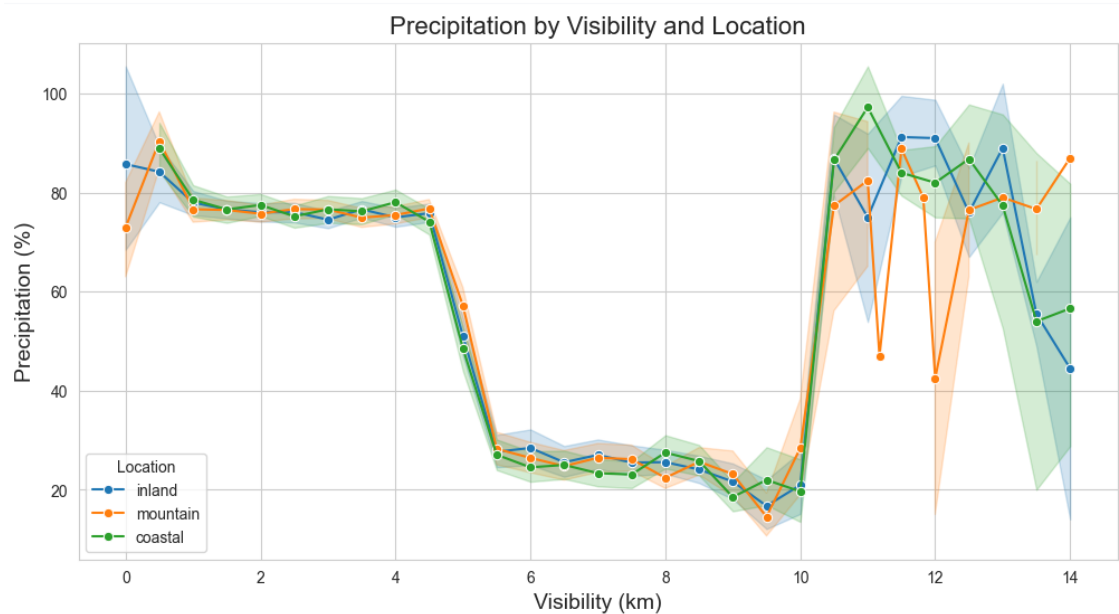


Obrázok 5: Zobrazenie grafu (countplot) pre počet typov počasia podľa UV indexu

- Sunny (slnečno): Obyčajne najvyšší počet výskytov pri vyšších hodnotách UV indexu.
- Rainy (dážď): Často sa vyskytuje pri nižších hodnotách UV indexu.
- Cloudy (oblačno) a Snowy (sneženie): Môžu mať variabilný počet výskytov v závislosti od podmienok.

4. **Zrážky podľa viditeľnosti a polohy:** Tento graf zobrazuje zrážky v závislosti od viditeľnosti a polohy. Line plot zobrazuje trend zrážok ako funkciu viditeľnosti a geografickej polohy. Os x reprezentuje viditeľnosť a os y zobrazuje percento zrážok. I Môžete vidieť, ako sa percento zrážok mení s rôznymi úrovňami viditeľnosti. Napríklad, ak je viditeľnosť vysoká, môže to naznačovať, že zrážky sú minimálne (nízke percento), zatiaľ čo pri nízkej viditeľnosti môže percento zrážok rásť v závislosti od geografickej polohy (vnútrozemská, pobrežná, horská).

Ak máme napríklad v horách viditeľnosť 12km, pravdepodobnosť zrážok je cca 41 %.

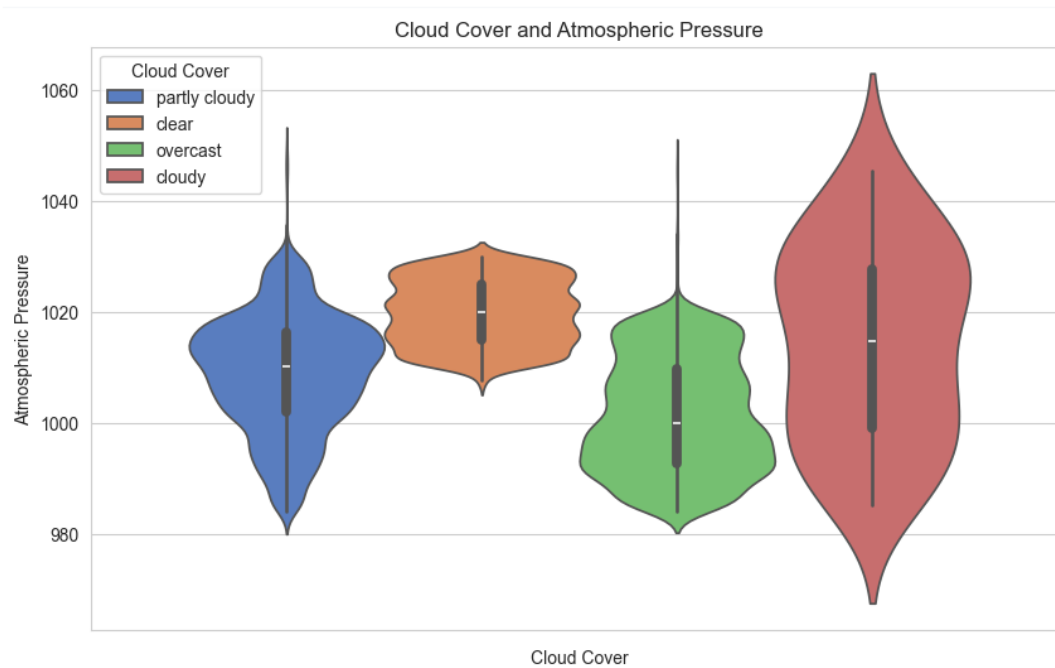


Obrázok 6: Zobrazenie grafu (lineplot) pre zrážky podľa viditeľnosti a polohy

5. **Zrážky podľa oblačnosti a atmosférického tlaku:** Tento graf zobrazuje rozloženie oblačnosti na osi x a atmosférického tlaku na osi y.

Šírka každého tvaru v grafe zobrazuje hustotu dát v rámci kategórie oblačnosti. Čím je tvar širší, tým viac hodnôt sa v tejto oblasti nachádza. Vnútorň box plot (čierny obdĺžnik s bielou čiarou):

- Biely bod - medián atmosférického tlaku pre danú kategóriu Cloud Cover.
- Obdĺžnikové pole - ukazuje interkvartilový rozsah (rozpätie medzi prvým a tretím kvartilom) - kde sa nachádza 50 % dát.
- Zvislé čiary - ukazujú rozsah dát.



Obrázok 7: Zobrazenie grafu (violinplot) pre oblačnosť a atmosférický tlak



## 3 Neurónová sieť

Základ tvorí funkcia *build\_model()*, v ktorej je možné nastaviť počet neurónov v 1. a 2. skrytej vrstve, pričom počet prvkov vstupnej vrstvy je daný ako počet stĺpcov trénovanej množiny X (ktorá je zakódovaná a škálovaná). Ako aktivačná funkcia v skrytých vrstvách neurónu je použitá ReLU<sup>15</sup> a vo výstupnej vrstve je to Softmax<sup>16</sup>. Taktiež je tu možné nastaviť rýchlosť učenia pre Adam Optimizer<sup>17</sup>. V tomto prípade bude neurónová sieť používať rovnaký batch<sup>18</sup> - 16.

Ako funkcia straty bola použitá *sparse\_categorical\_crossentropy*, ktorá sa používa pri reprezentovaní údajov ako celé čísla (bol použitý Label Encoding). Pre meranie úspešnosti bola použitá *sparse\_categorical\_data*, kde sú údaje taktiež reprezentované ako celá čísla (bol použitý Label Encoding), na rozdiel od One-Hot Encoding.

### 3.1 Architektúra siete s pretrénovaním

Pretrénovanie nastáva vtedy, keď sa model príliš prispôsobí trénovacím údajom a nedokáže generalizovať na nové, neznáme dáta. Aby bolo dosiahnuté, bolo pridaných viac neurónov do 1. skrytej (128 neurónov) a 2. skrytej (64 neurónov) vrstvy a taktiež bol zvýšený počet epoch<sup>19</sup> (500). Vstupná konfigurácia siete je:

- **rýchlosť učenia:** 0.001
- **1. skrytá vrstva:** 128
- **2. skrytý vrstva:** 64
- **počet epoch:** 500

Ak sa strata na tréningových dátach neustále znižuje a dosahuje veľmi nízke hodnoty, zatiaľ čo strata na validačných dátach najskôr klesá, ale potom začne stúpať alebo oscilovať, znamená to, že model sa príliš prispôbil trénovacím dátam a začína strácať schopnosť generalizovať na nové dáta.

---

<sup>15</sup>Aktivačná funkcia používaná v neurónových sieťach, ktorá transformuje vstupy tak, že všetky záporné hodnoty sú nahradené nulou a kladné hodnoty zostávajú nezmenené.

<sup>16</sup>Aktivačná funkcia, ktorá transformuje výstupy modelu na pravdepodobnosti, pričom ich súčet je 1. Pre každú triedu vypočíta pravdepodobnosť, že vstup patrí do tejto triedy.

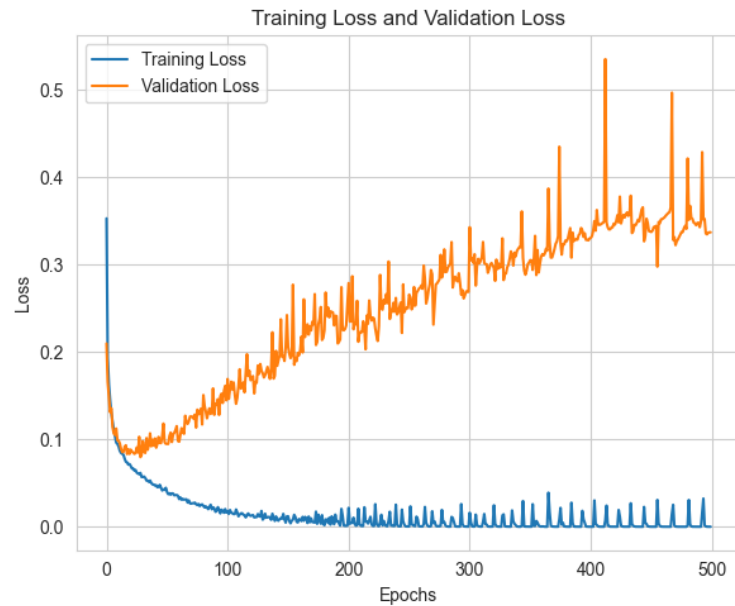
<sup>17</sup>adaptívne nastavuje rýchlosť učenia

<sup>18</sup>Počet vzoriek (riadkov) z trénovacích dát, ktoré model spracováva naraz predtým, než aktualizuje váhy

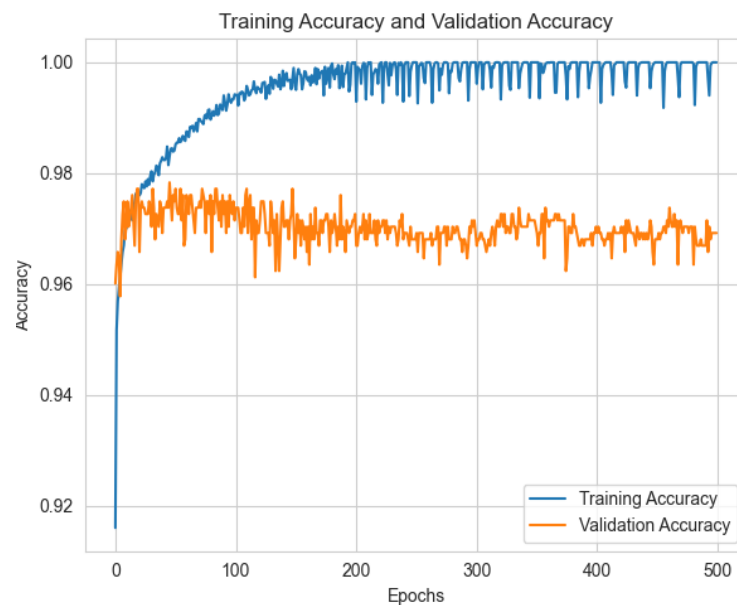
<sup>19</sup>Jeden kompletný prechod všetkými trénovacími dátami

Ak presnosť na tréningových dátach neustále stúpa a dosahuje veľmi vysoké hodnoty (blíži sa k 1), zatiaľ čo presnosť na validačných dátach dosiahne určitý bod a potom začne klesať alebo stagnovať, znamená to, že model sa príliš učí detaily tréningových dát (jeho schopnosť generalizovať sa zhoršuje).

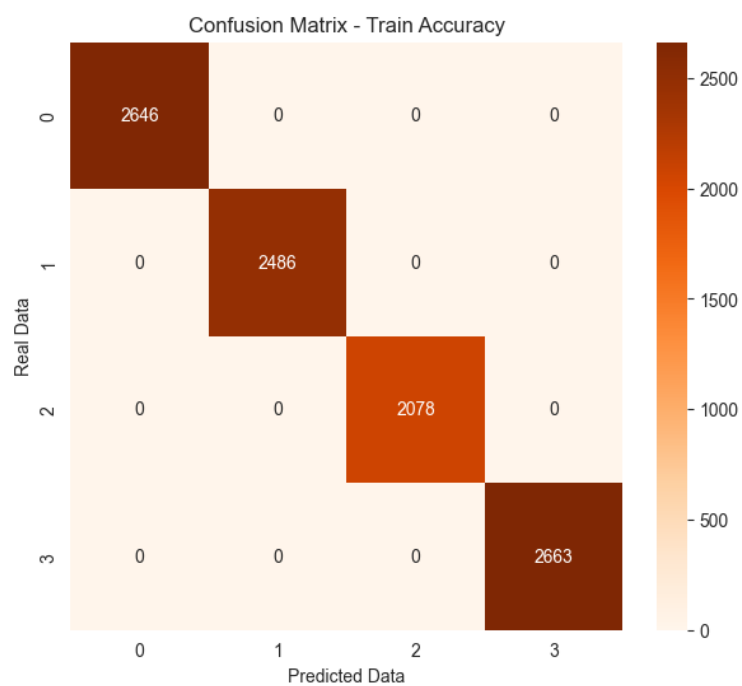
V tomto prípade dosiahla tréningová úspešnosť hodnotu **1** a testovacia úspešnosť hodnotu **0.977272**.



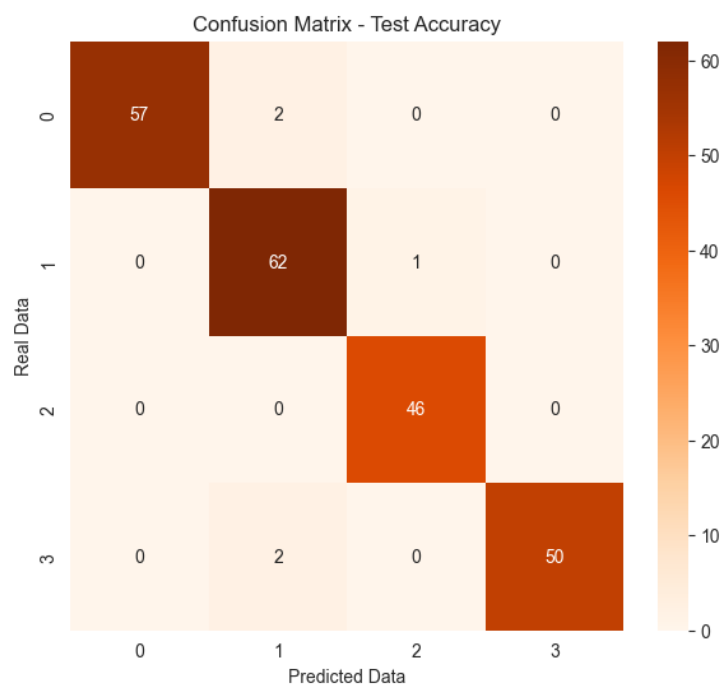
Obrázok 8: Graf tréningovej a validačnej straty



Obrázok 9: Graf tréningovej a validačnej úspešnosti



Obrázok 10: Zobrazenie konfúznej matice pre trénovaciu úspešnosť



Obrázok 11: Zobrazenie konfúznej matice pre testovaciu úspešnosť

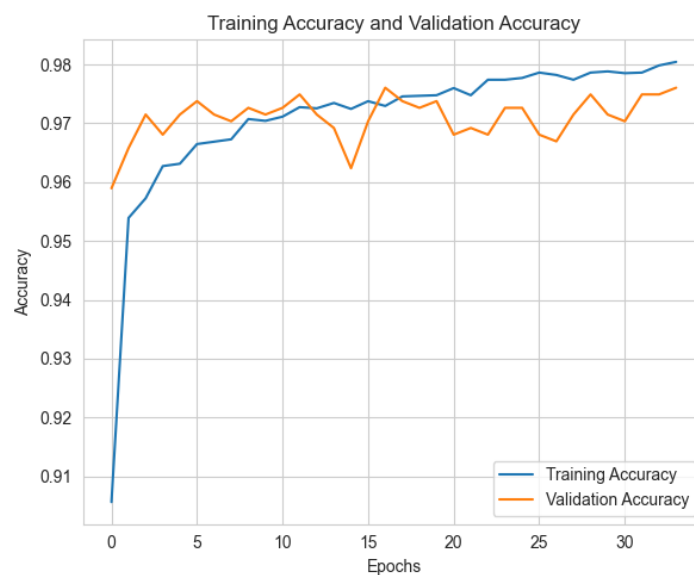
## 3.2 Architektúra siete bez pretrénovania

V tomto prípade pracujeme s konfiguráciou uvedenou v sekcii 3.1, avšak aby sme odstránili pretrénovanie, pridáme **Early Stopping**, pre skoré zastavenie tréningu v prípade, že sa vyskytne náznak pretrénovania. To môže ovplyvniť počet epôch, ktorý sa v tomto prípade zníži (ako je možné vidieť na grafoch).

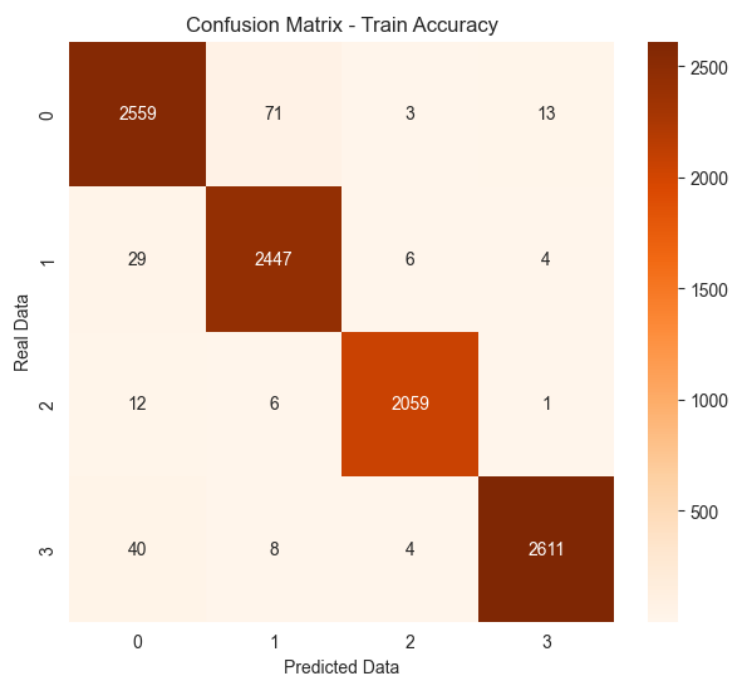
Ako je možné na grafoch vidieť, táto neurónová sieť nemá vďaka použitiu Early Stoppingu vyššie uvedené známky pretrénovania.



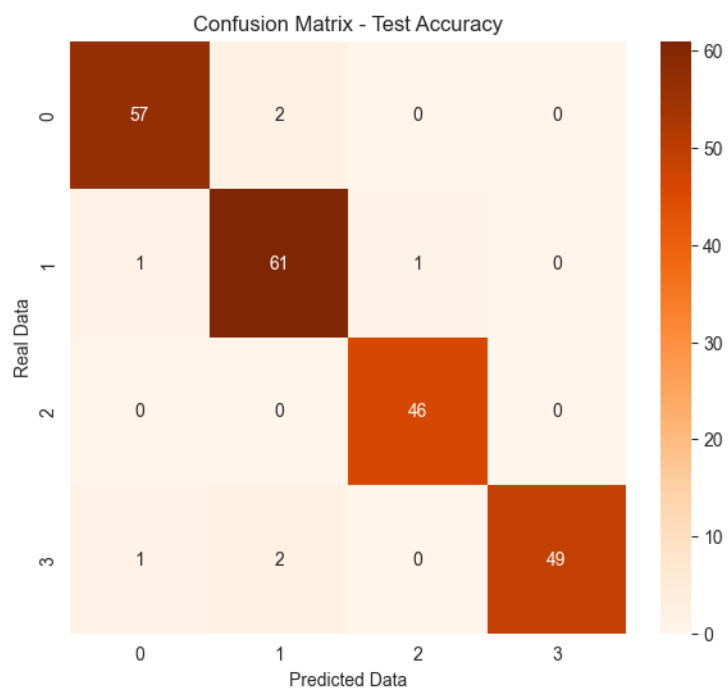
Obrázok 12: Graf tréningovej a validačnej straty s použitím Early Stopping



Obrázok 13: Graf tréningovej a validačnej úspešnosti s použitím Early Stopping



Obrázok 14: Zobrazenie konfúznej matice pre trénovaciu úspešnosť s použitím Early Stopping



Obrázok 15: Zobrazenie konfúznej matice pre testovaciu úspešnosť s použitím Early Stopping

V tomto prípade dosiahla trénovacia úspešnosť hodnotu **0.98004** a testovacia úspešnosť hodnotu **0.968181**, čo indikuje zlepšenie, oproti predchádzajúcej sieti (bez Early Stoppingu).

### 3.3 Experimenty s trénovaním siete

Pre zobrazenie rozdielov v trénovaniach a výsledných úspešnostiach na základe rôznych vstupných hyperparametrov bolo vykonaných 5 experimentov, pričom v prvých troch experimentoch bola zmenená hodnota parametra rýchlosti učenia a v ďalších dvoch bol zmenený počet neurónov v 1. a 2. skrytej vrstve.

Všetky tieto zmeny, vrátane výpisu úspešnosti jednotlivých experimentov sú zapísané v tabuľke nižšie.

Číslo	Rýchlosť učenia	1. vrstva	2. vrstva	Trénovacia úspešnosť	Testovacia úspešnosť
1	0.01	64	32	0.996354	0.981818
2	0.0001	64	32	0.976299	0.972727
3	0.01	64	32	0.986934	0.972727
4	0.01	128	16	0.980958	0.968182
5	0.01	32	32	0.980654	0.968182

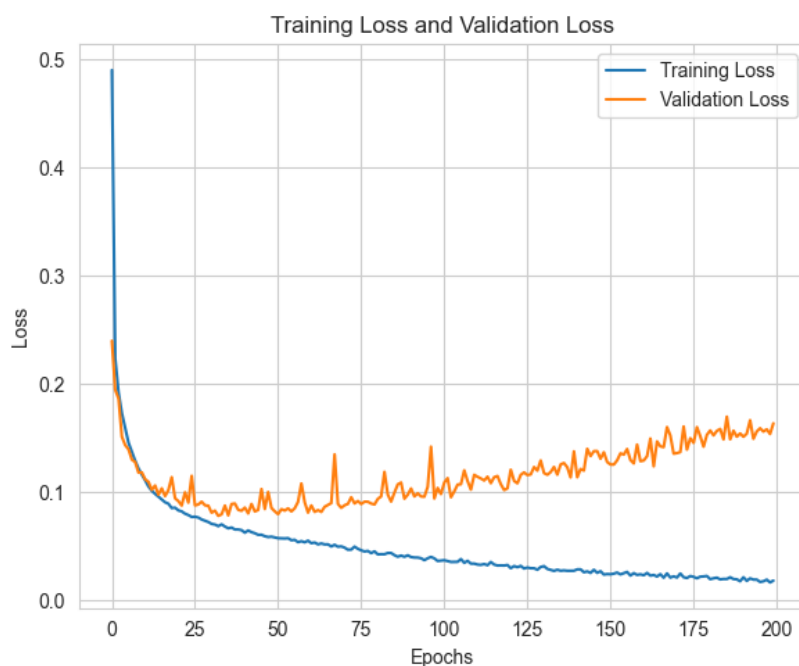
Tabuľka 3: Konfigurácia jednotlivých experimentov a presnosť trénovania a testovania

Výber najlepšieho a najhoršieho testu  $T$  bol zvolený na základe priemeru trénovacej  $t_1$  a testovacej  $t_2$  úspešnosti:

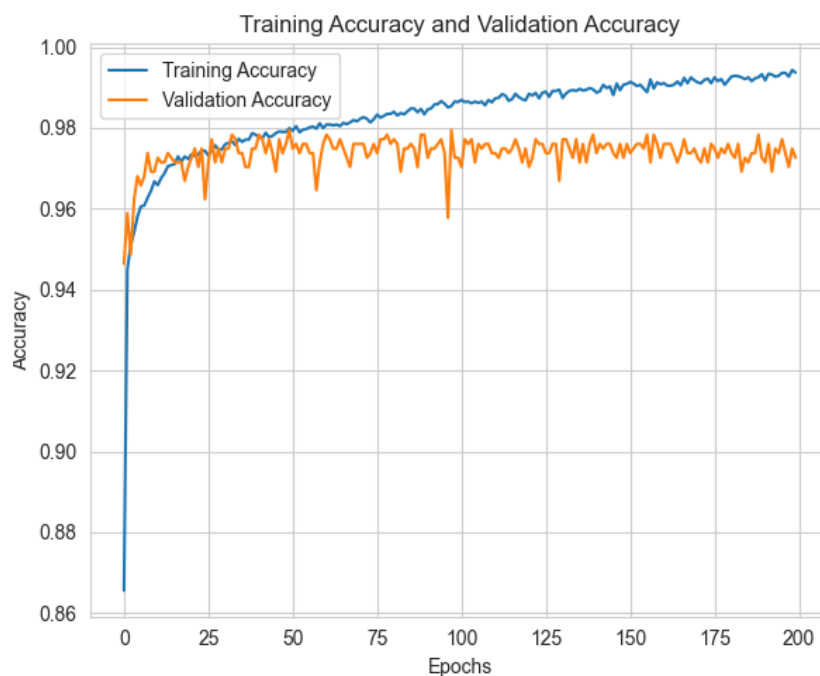
$$T = \frac{t_1 + t_2}{2},$$

pričom tento priemer bol vykonaný pre každý experiment v tabuľke. Najvyššia hodnota  $T$  predstavuje najlepší experiment (č. 1 s priemerom 0.989085) a najmenšia hodnota predstavuje najhorší (č. 5 s priemerom 0.974418) experiment.

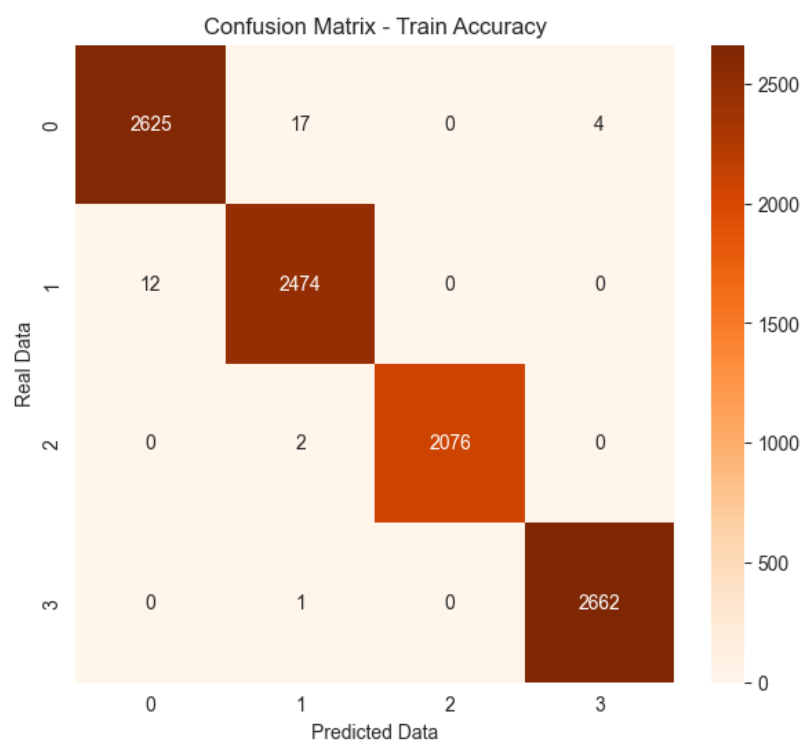
- Grafy úspešnosti a strát a konfúzne matice pre najlepší experiment



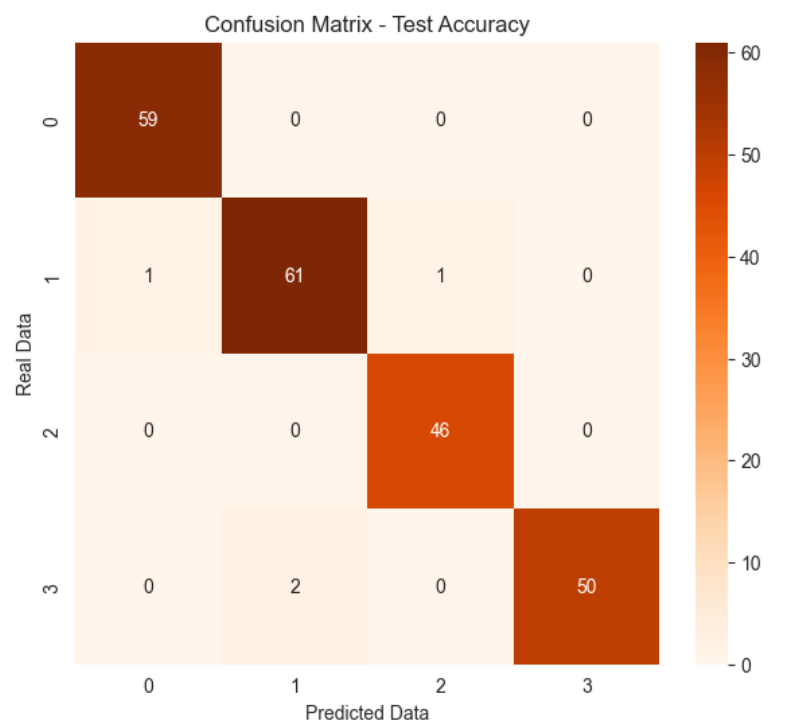
Obrázok 16: Graf trénovacej a validačnej straty pre najlepší experiment



Obrázok 17: Graf trénovacej a validačnej úspešnosti pre najlepší experiment



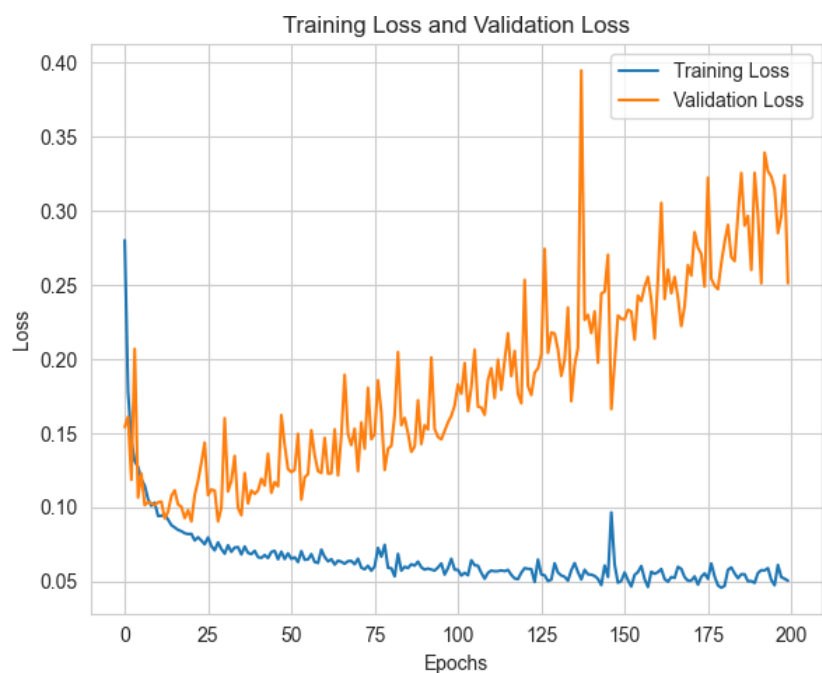
Obrázok 18: Konfúzna matica trénovacej množiny pre najlepší experiment



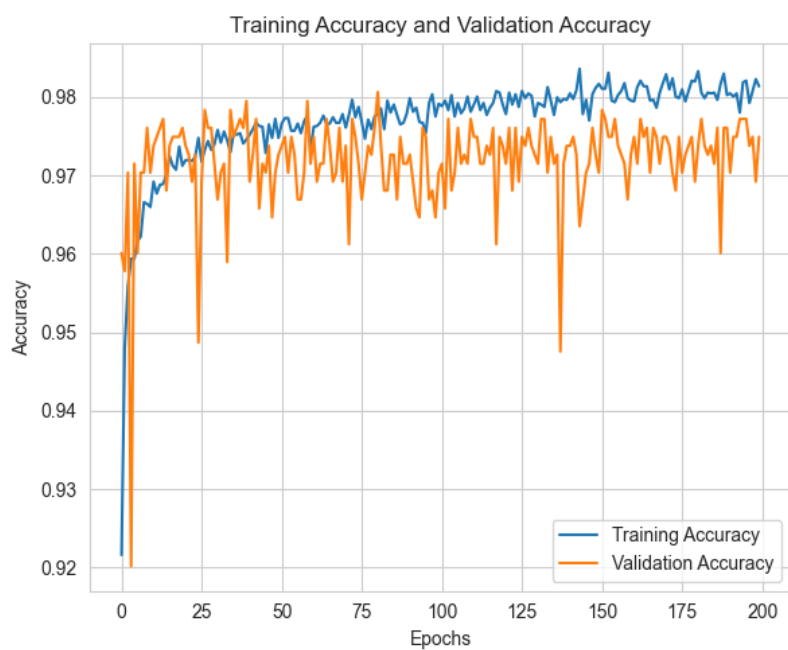
Obrázok 19: Konfúzna matica testovacej množiny pre najlepší experiment



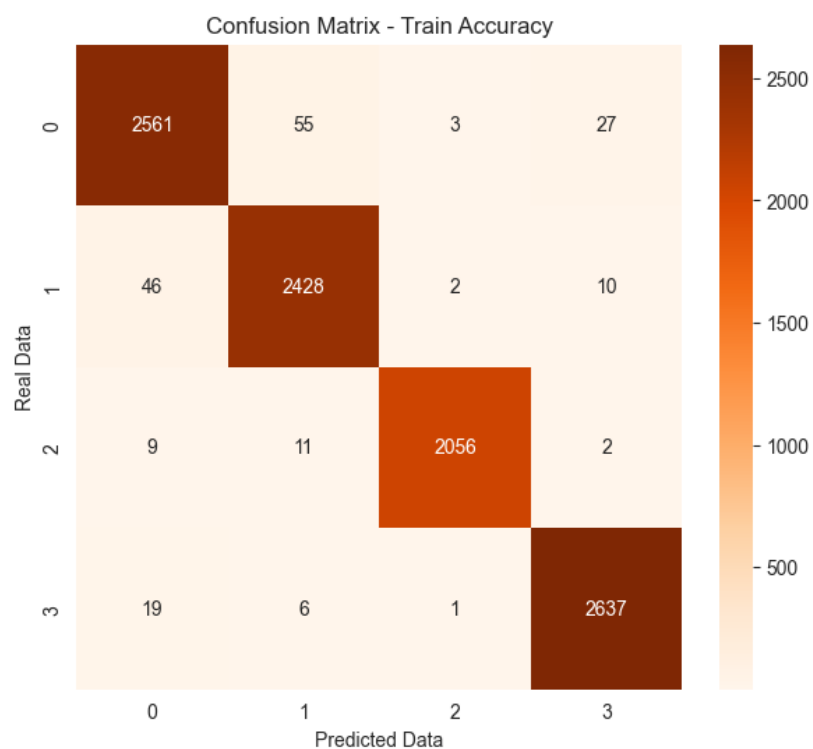
- Grafy úspešnosti a strát a konfúzne matice pre najhorší experiment



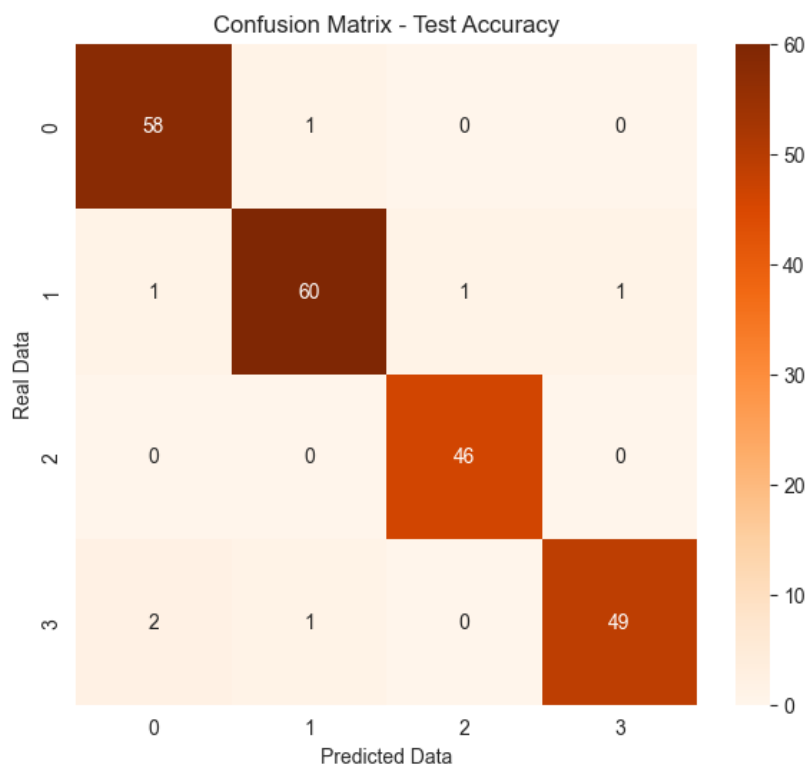
Obrázok 20: Graf trénovacej a validačnej straty pre najhorší experiment



Obrázok 21: Graf trénovacej a validačnej úspešnosti pre najhorší experiment



Obrázok 22: Konfúzna matica trénovacej množiny pre najhorší experiment

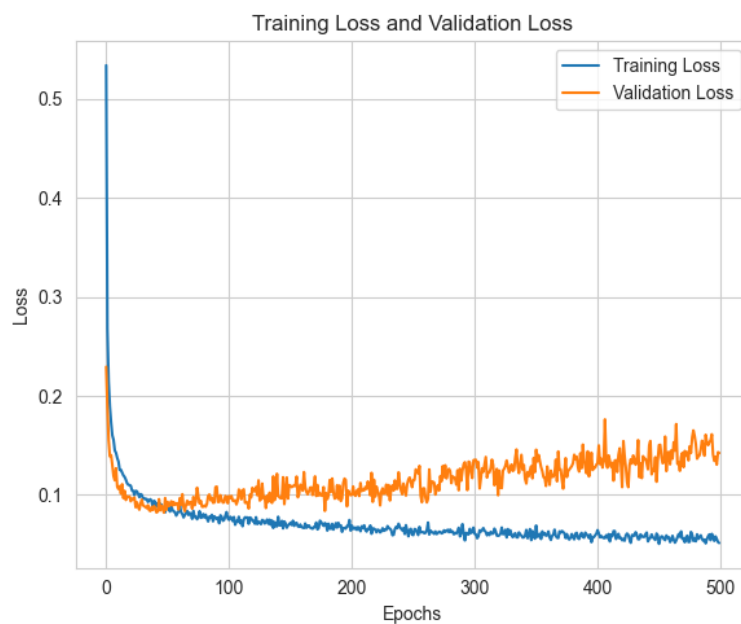


Obrázok 23: Konfúzna matica testovacej množiny pre najhorší experiment

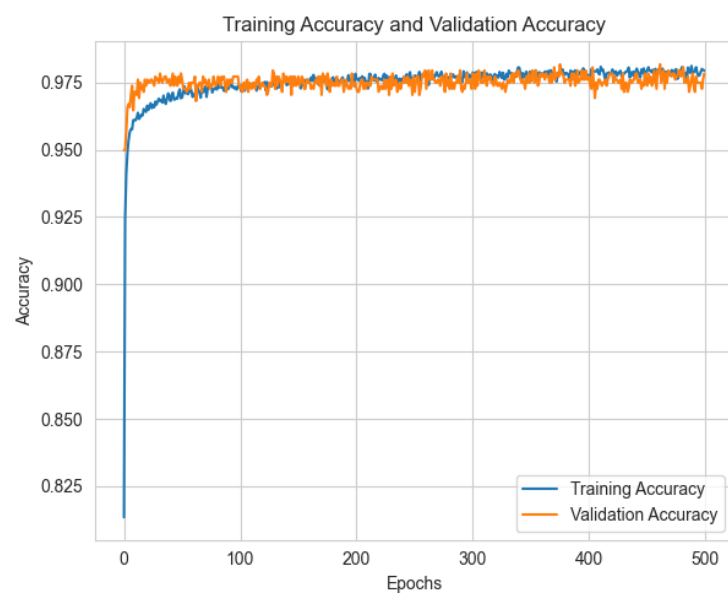
### 3.3.1 Dropout vrstvy

Dropout vrstvy môžu pomôcť znížiť pretrénovanie pridaním náhodného vyradovania neurónov počas tréningu. Túto vrstvu pridávame vždy s pravdepodobnosťou vypnutia neurónov (napríklad 0.2, resp. 20 % pravdepodobnosť, že sa neuróny vypnú).

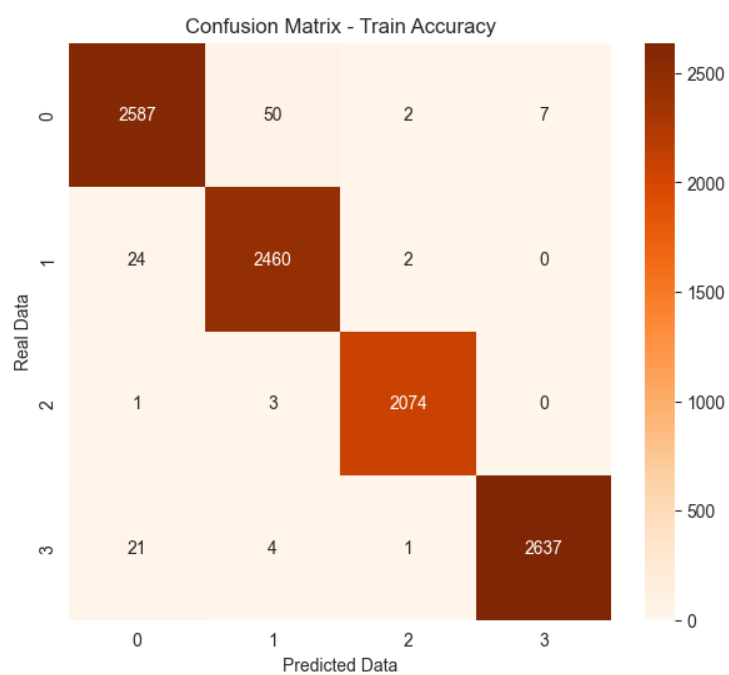
Dropout vrstvu môžeme pridať do každej skrytej vrstvy a tak povypínať neuróny s rôznou pravdepodobnosťou. V tomto prípade pracujeme s konfiguráciou uvedenou v sekcii 3.1, pričom boli pridané 2 dropout vrstvy do skrytej vrstvy neurónovej siete - obidve s pravdepodobnosťami 40%.



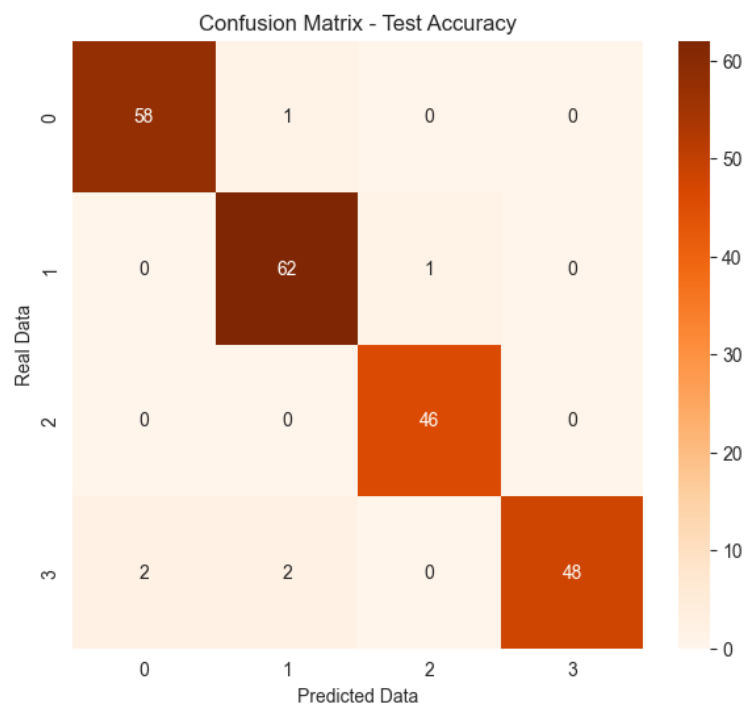
Obrázok 24: Graf trénovacej a validačnej straty pre neurónovú sieť s dropout vrstvami



Obrázok 25: Graf trénovacej a validačnej úspešnosti pre neurónovú sieť s dropout vrstvami



Obrázok 26: Zobrazenie konfúznej matice pre trénováciu úspešnosť pre neurónovú sieť s dropout vrstvami



Obrázok 27: Zobrazenie konfúznej matice pre testovaciu úspešnosť pre neurónovú sieť s dropout vrstvami

Uvedené parametre neurónovej siete boli vybrané zámerne - aby bolo možné demonštrovať na rovnakej pretrénovanej sieti rôzne úspešnosti. Ako je možné vidieť, úspešnosť sa oproti pôvodnej sieti zlepšila (trénovacia úspešnosť má hodnotu 0.98835 a testovacia úspešnosť má hodnotu 0.972727).

# Záver

Cieľom tohto zadania bolo kategorizovať počasie do štyroch kategórií (Rainy, Cloudy, Sunny, Snowy), pričom ako vstupné parameter boli použité množiny Cloud Cover, Season, Location, Temperature, Humidity, Wind Speed, Precipitation (%), Atmospheric Pressure, UV Index, Visibility (km).

Dáta boli očistené od outlierov, chýbajúcich hodnôt alebo odchýlok, boli zakódované Label Encoderom a následne normalizované MinMaxScalerom. S ďalej pracovali pri tréovaní neurónovej siete, pričom bolo vykonaných množstvo experimentov s cieľom zistiť, ako ovplyvňuje úspešnosť siete zmena jednotlivých hyperparametrov.

Taktiež bol pomocou Exploratory Data Analysis (EDA) analyzovaný dataset pred zakódovaním a následne bolo (aj) pomocou korelačnej matice nájdených 5 zaujímavých vzťahov, pre ktoré boli vykreslené a popísané príslušné grafy.