

VILNIAUS UNIVERSITETAS
MATEMATIKOS IR INFORMATIKOS FAKULTETAS
PROGRAMŲ SISTEMŲ BAKALAURO STUDIJŲ PROGRAMA

**Vieno neurono mokymas sprendžiant
klasifikavimo uždavinį**

Skaitmeninis intelektas ir sprendimų priėmimas

Užduoties ataskaita

Atliko: 4 kurso 4 grupės studentas
Vilius Puškunalis

Vilnius – 2023

TURINYS

1. ĮVADAS	2
2. DUOMENYS	3
3. PROGRAMOS KODAS SU KOMENTARAIŠ	4
4. NEURONO MOKYMO TAISYKLĖS	8
4.1. Pradinių svorių reikšmių parinkimas	8
4.2. Stochastinis gradientinis nusileidimas.....	8
4.3. ADALINE taisyklė	8
5. TYRIMO REZULTATAI	9
5.1. Paklaidos reikšmių priklausomybė nuo epochų skaičiaus	9
5.2. Klasifikavimo tikslumo priklausomybė nuo epochų skaičiaus.....	10
5.3. Rezultatų priklausomybė nuo skirtingų mokymosi greičio reikšmių.....	11
6. GAUTI DUOMENYS	13
7. REZULTATAI IR IŠVADOS	14
ŠALTINIAI	15

1. Įvadas

Užduoties tikslas

Apmokyti vieną neuroną spręsti dviejų klasių uždavinį ir atlikti tyrimą su dviem duomenų aibėmis.

Užduoties variantas

Studento numeris – 2016025.

Paskutinio skaitmens 5 dalybos iš 3 liekana – 2.

2-as variantas – neurono mokymui naudoti stochastinį gradientinį nusileidimą ir ADALINE mokymo taisyklę.

Uždaviniai

- Parsisiųsti ir paruošti irisų bei krūties vėžio navikų duomenų aibes.
- Sukurti programą, kuri įgyvendintų vieno neurono mokymo ir testavimo procesą, sprendžiant klasifikavimo užduotį.
- Irisų ir krūties vėžio navikų duomenų aibėms ištirti:
 - Paklaidos reikšmės priklausomybę nuo epochų skaičiaus.
 - Klasifikavimo tikslumo priklausomybę nuo epochų skaičiaus.
 - Rezultatų priklausymą nuo skirtingų mokymosi greičio reikšmių.

2. Duomenys

Dirbtinio neurono mokymui ir testavimui buvo naudojamos irisų [Fis36] bei krūties vėžio [Wo192] duomenų aibės.

Irisų duomenų aibėje yra 3 rūšys (Setosa, Versicolor, Virginica) po 50 duomenų įrašų, tačiau darbe naudojamos Versicolor ir Virginica rūšys. Kiekvienas įrašas turi po 4 požymius. Taigi, apdorotoje duomenų aibėje yra 100 įrašų.

Krūties vėžio duomenų aibėje yra 699 duomenų įrašai. Kiekvienas įrašas turi po 9 požymius. Atmetus eilutes, kur kai kurie duomenys trūksta, lieka 683 įrašai. Duomenys klasifikuojami į nepiktybinius ir piktybinius navikus.

Duomenys buvo padalinti į mokymo ir testavimo aibes santykiu 70:30.

3. Programos kodas su komentrais

Dirbtinis neuronas įgyvendintas Python programavimo kalba. Gali būti modifikuojami šie hiperparametrai: mokymo greitis, epochų skaičius, santykis tarp mokymo ir testavimo aibių.

```
1 import pandas as pd
2 import numpy as np
3
4 # Nuskaityti ir paruošti iris duomenų aibę
5 def fetch_and_prepare_iris_data():
6     # Skaitome CSV formato failą
7     df = pd.read_csv('iris/iris.data', header=None)
8
9     # Filtruojame duomenis
10    filtered_indices = df[df.iloc[:, -1].isin(['Iris-versicolor', 'Iris-
        -virginica'])].index
11    df = df.loc[filtered_indices]
12
13    # Pakeičiame klasių pavadinimus į numerius
14    df.iloc[:, -1] = df.iloc[:, -1].map({'Iris-versicolor': 0, 'Iris-
        -virginica': 1})
15
16    return df
17
18 # Nuskaityti ir paruošti kruties vezio duomenų aibę
19 def fetch_and_prepare_breast_cancer_data():
20     # Skaitome CSV formato failą
21     df = pd.read_csv('breast-cancer/breast-cancer-wisconsin.data',
        header=None)
22
23     # Pasaliname pirmą stulpelį
24     df = df.iloc[:, 1:]
25
26     # Išmetame eilutes su klaustukais
27     indices_to_keep = ~df.isin(['?']).any(axis=1)
28     df = df[indices_to_keep]
29
30     # Konvertuojame visus stulpelius į skaičių formata (problema dėl
        buvusių klaustukų)
31     for col in df.columns:
32         df[col] = pd.to_numeric(df[col], errors='coerce')
33
34     # Pakeičiame klasių numerius
35     df = df.loc[indices_to_keep].copy()
```

```

36     df.iloc[:, -1] = df.iloc[:, -1].map({2: 0, 4: 1})
37
38     return df
39
40 # Skaido duomenis i mokymo ir testavimo rinkinius
41 def split_data(train_test_split, df):
42     # Maisome duomenis
43     df = df.sample(frac=1, random_state=0)
44
45     # Apskaiciuojame indeksa, kuriuo padalinsime duomenis
46     split_index = int(train_test_split * len(df))
47
48     # Daliname duomenis
49     train_data = df[:split_index]
50     test_data = df[split_index:]
51
52     return train_data, test_data
53
54 # Neuronu mokymas naudojant ADALINE taisykle ir stochastini gradientini
    nusileidima
55 def adaline_SGD(train_data, epochs, learning_rate):
56     features = train_data.iloc[:, :-1].values
57     classes = train_data.iloc[:, -1].values
58
59     # Inicializuojame svorius ir bias
60     weights = np.zeros(features.shape[1])
61     bias = 0
62
63     # Mokymo ciklas
64     error_per_epoch = []
65     accuracy_per_epoch = []
66     for epoch in range(epochs):
67         total_error = 0
68         for xi, t in zip(features, classes):
69             # Aktyvacijos funkcija yra tiesine
70             y = np.dot(xi, weights) + bias
71             weights += learning_rate * (t - y) * xi
72             bias += learning_rate * (t - y)
73
74             error = (t - y) ** 2
75             total_error += error
76
77     error_per_epoch.append(total_error)

```

```

78         accuracy_per_epoch.append(accuracy(train_data, weights, bias,
79                                             False))
80
81     return weights, bias, error_per_epoch, accuracy_per_epoch
82
83 # Skaiciuoja klasifikavimo tiksluma
84 def accuracy(data, weights, bias, should_print):
85     features = data.iloc[:, :-1].values
86     classes = data.iloc[:, -1].values
87
88     y = np.dot(features, weights) + bias
89
90     # Slenkstine funkcija
91     y = np.where(y > 0.5, 1, 0)
92
93     if should_print:
94         for real_class, predicted_class in zip(classes, y):
95             print(f"Real class: {real_class}, predicted class: {
96                   predicted_class}, prediction correct: {real_class ==
97                   predicted_class}")
98
99     return np.mean(y == classes)
100
101 # Skaiciuoja vidutine kvadratine paklaida
102 def mean_squared_error(data, weights, bias):
103     features = data.iloc[:, :-1].values
104     classes = data.iloc[:, -1].values
105
106     # Aktyvacijos funkcija yra tiesine
107     y = np.dot(features, weights) + bias
108
109     # MSE skaiciavimas
110     mse = np.mean((classes - y)**2)
111
112     return mse
113
114 if __name__ == "__main__":
115     train_test_split = 0.7 # Dalinimas i mokymo/testavimo aibes 70:30
116     epochs = 100 # Epochu skaicius
117     learning_rate = 0.0001 # Mokymosi greitis
118
119     # Pasirinkti iris arba krutu vezio duomenų aibes
120     df = fetch_and_prepare_iris_data()

```

```

118 #df = fetch_and_prepare_breast_cancer_data()
119
120 # Duomenų dalinimas į mokymo ir testavimo aibes
121 train_data, test_data = split_data(train_test_split, df)
122
123 # Neuronų mokymas
124 weights, bias, error_per_epoch, accuracy_per_epoch = adaline_SGD(
125     train_data, epochs, learning_rate)
126
127 # Tikslumo skaičiavimas
128 test_accuracy = accuracy(test_data, weights, bias, True)
129 test_mse = mean_squared_error(test_data, weights, bias)
130
131 print(f"Gauti svoriai: {np.round([bias, *weights], decimals=4)}")
132 print(f"Gautos paklaidos po kiekvienos epochos mokymo duomenims: {
133     np.round(error_per_epoch, decimals=4)}")
134 print(f"Gauta paklaida testavimo duomenims: {test_mse:.4f}")
135 print(f"Gautas klasifikavimo tikslumas po kiekvienos epochos mokymo
136     duomenims: {np.round(accuracy_per_epoch, decimals=4)}")
137 print(f"Gautas klasifikavimo tikslumas testavimo duomenims: {
138     test_accuracy:.4f}")

```


4. Neuronų mokymo taisyklės

4.1. Pradinių svorių reikšmių parinkimas

Pradiniai svoriai bei poslinkis visi nustatyti su reikšme lygia 0.

4.2. Stochastinis gradientinis nusileidimas

Stochastinis gradientinis nusileidimas (angl. Stochastic Gradient Descent; SGD) yra optimizavimo algoritmas, skirtas rasti funkcijos lokalų minimumą arba maksimumą. Stochastinis gradientinis nusileidimas yra vienas iš gradientinio nusileidimo algoritmo variantų.

Stochastinio gradientinio nusileidimo atveju viena epocha atitinka m iteracijų, čia m yra mokymo duomenų kiekis [Kur21].

4.3. ADALINE taisyklė

ADALINE (angl. Adaptive Linear Neuron) yra dirbtinio neuronų mokymo taisyklė, sukurta 1959-aisiais. ADALINE esminis dalykas – naudojama tiesinė aktyvacijos funkcija. ADALINE mokymosi taisyklių formulė yra ta pati kaip ir perceptrono. Svoriai keičiami (atnaujinami) pagal šią mokymo taisyklę (formulę) [Kur21]:

$$w_k := w_k - \eta(t_i - y_i)(-x_{ik})$$

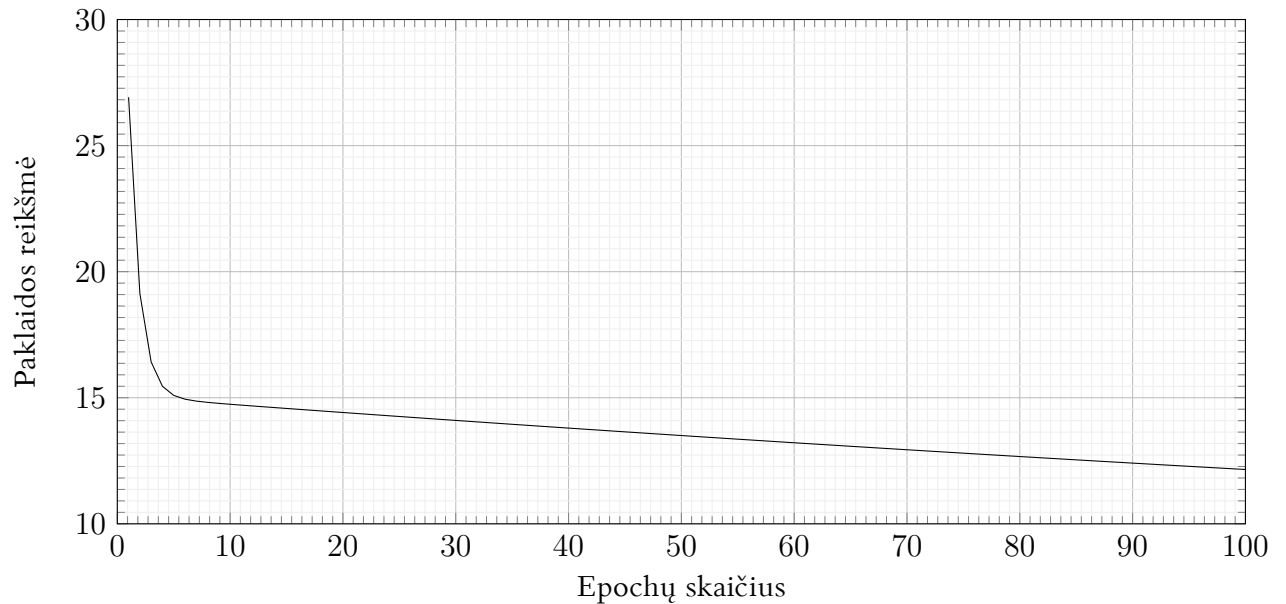
arba

$$w_k := w_k + \eta(t_i - y_i)(-x_{ik})$$

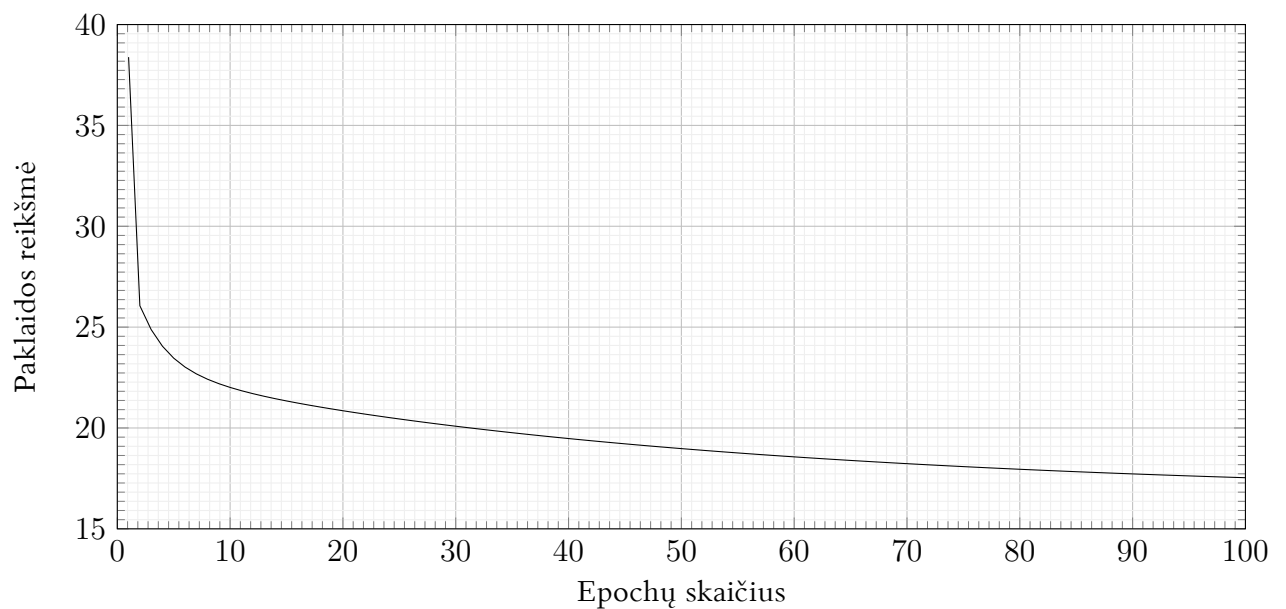
5. Tyrimo rezultatai

5.1. Paklaidos reikšmių priklausomybė nuo epochų skaičiaus

Ir irisų, ir krūties vėžio navikų paklaidos reikšmė greit nukrenta, o po to ir toliau lėtai, logaritmiškai mažėja didėjant epochų skaičiui.



1 pav. Paklaidos reikšmės priklausomybė nuo epochų skaičiaus irisus klasifikuojančiam neuronui

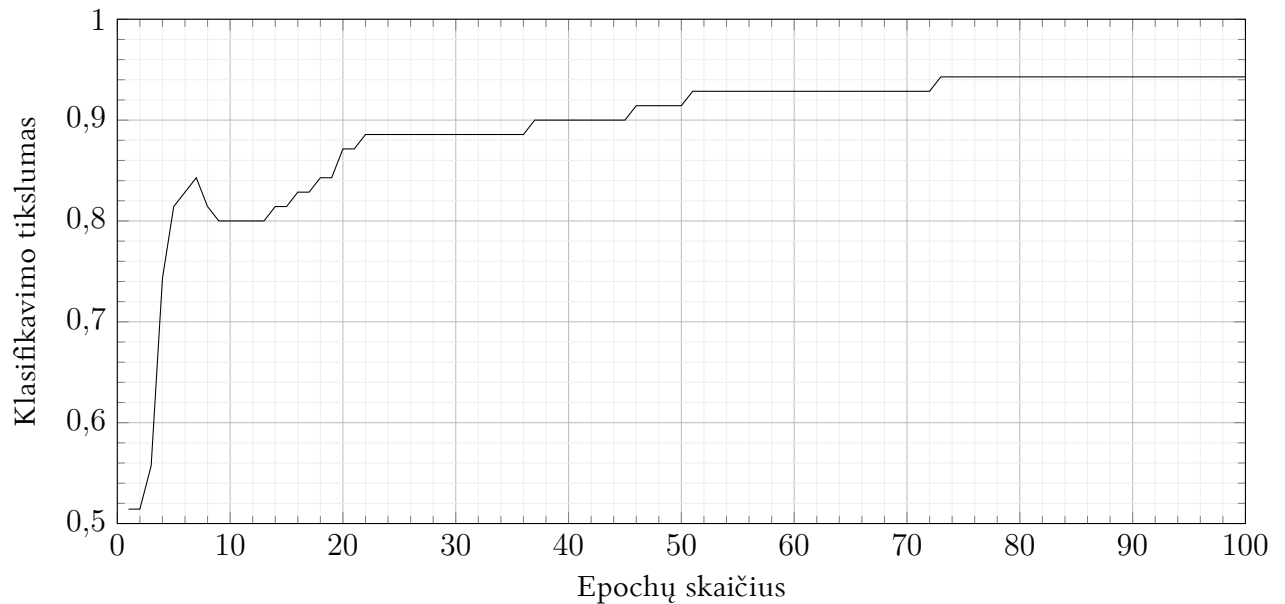


2 pav. Paklaidos reikšmės priklausomybė nuo epochų skaičiaus krūties vėžio navikus klasifikuojančiam neuronui

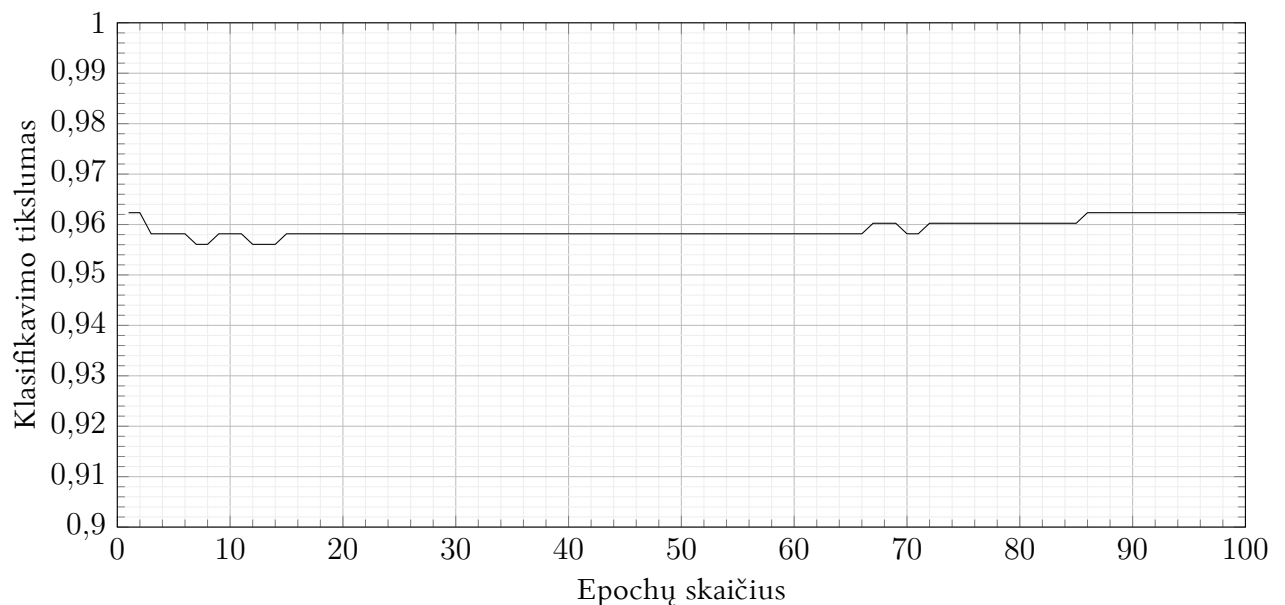
5.2. Klasifikavimo tikslumo priklausomybė nuo epochų skaičiaus

Irisų klasifikavimo tikslumas labiausiai auga pirmomis epochomis, vėlesnėmis epochomis auga ir toliau, tačiau lėčiau.

Krūties vėžio navikų klasifikavimo tikslumas visada lieka aukštas ir beveik nepriklauso nuo epochų skaičiaus.



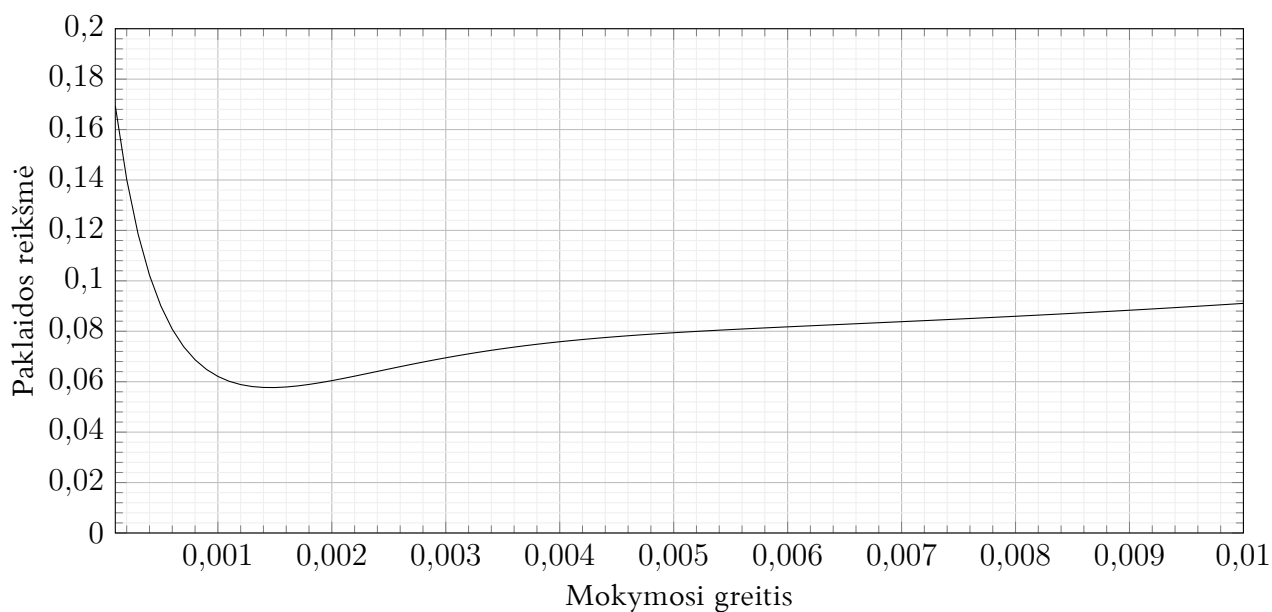
3 pav. Klasifikavimo tikslumo priklausomybė nuo epochų skaičiaus irisus klasifikuojančiam neuronui



4 pav. Klasifikavimo tikslumo priklausomybė nuo epochų skaičiaus krūties vėžio navikus klasifikuojančiam neuronui

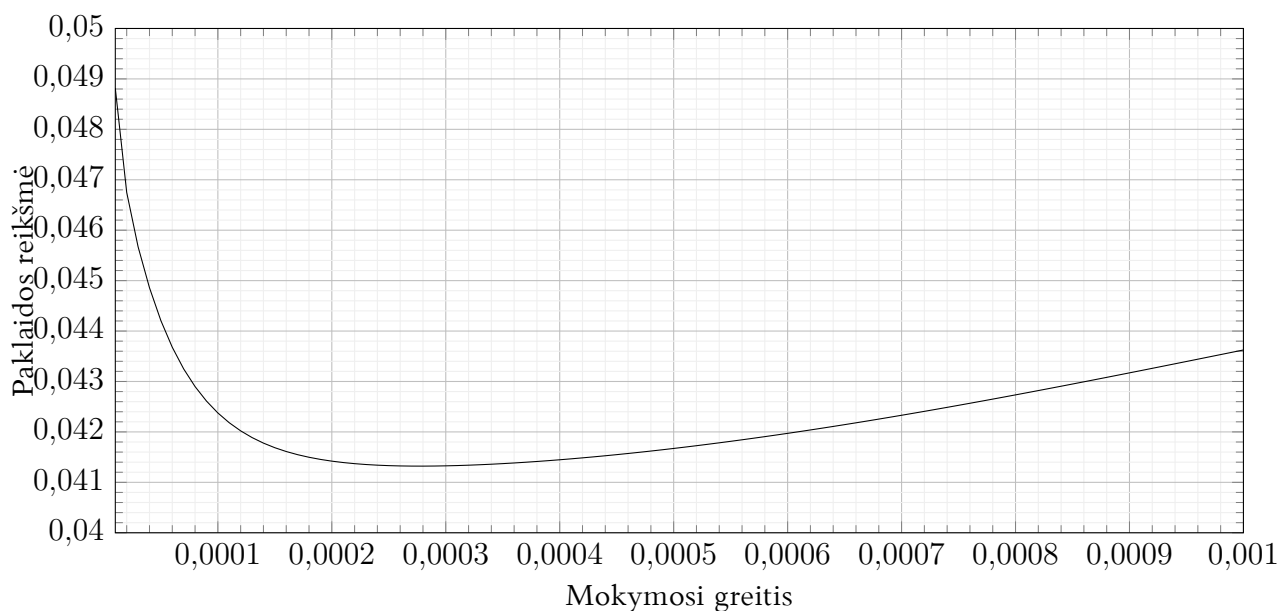
5.3. Rezultatų priklausomybė nuo skirtingų mokymosi greičio reikšmių

Irisus klasifikuojančio neurono paklaidos reikšmė esant mažam mokymosi greičiui yra aukšta, tada nukrenta iki minimumo, kai mokymosi greitis yra 0,0015.



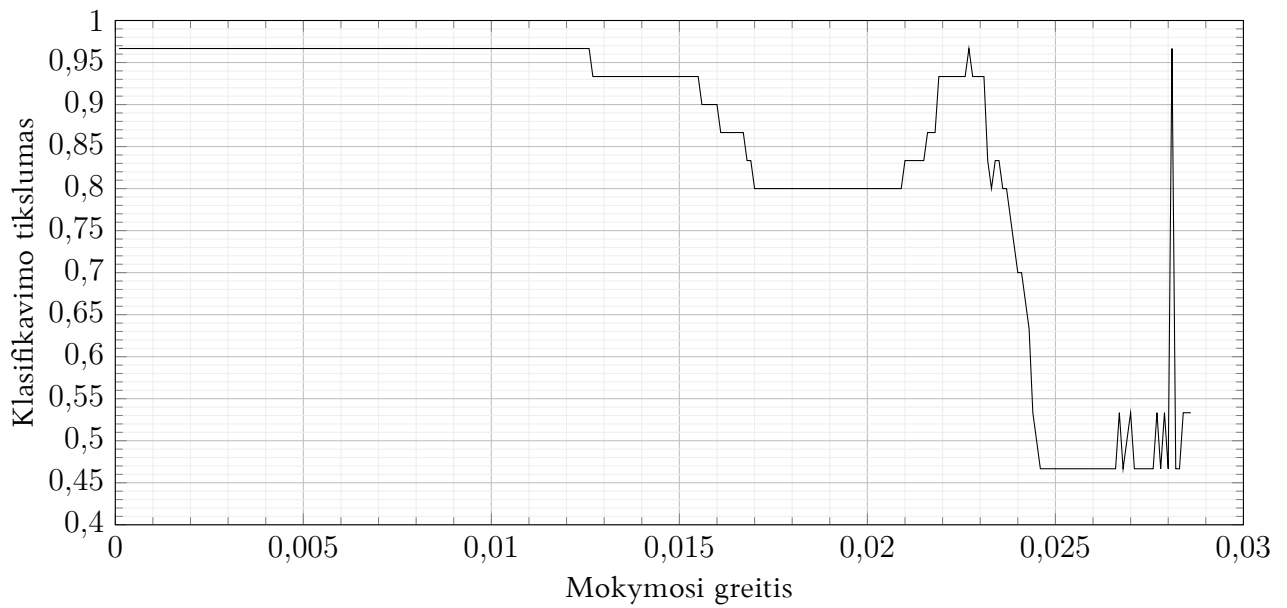
5 pav. Paklaidos reikšmės priklausomybė nuo mokymosi greičio irisus klasifikuojančiam neuronui

Krūties vėžio navikus klasifikuojančio neurono paklaidos reikšmė esant mažam mokymosi greičiui yra aukšta, tačiau nukrenta iki minimumo, kai mokymosi greitis yra 0,00028.



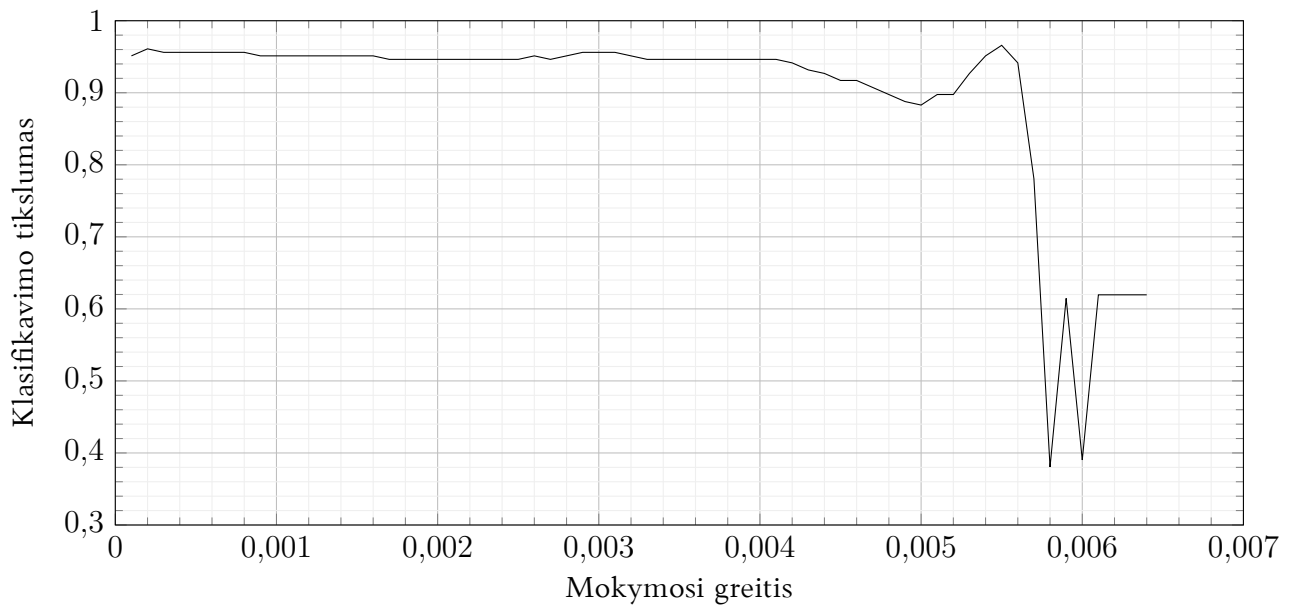
6 pav. Paklaidos reikšmės priklausomybė nuo mokymosi greičio krūties vėžio navikus klasifikuojančiam neuronui

Aukščiausias irisų klasifikavimo tikslumas gaunamas, kai mokymosi greitis yra nuo 0,0001 iki 0,0126.



7 pav. Klasifikavimo tikslumo priklausomybė nuo mokymosi greičio irisus klasifikuojančiam neuronui

Aukščiausias krūties vėžio navikų klasifikavimo tikslumas pasiekiamas, kai mokymosi greitis lygus 0,0055.



8 pav. Klasifikavimo tikslumo priklausomybė nuo mokymosi greičio krūties vėžio navikus klasifikuojančiam neuronui

6. Gauti duomenys

Ir irisus, ir krūties vėžio navikus klasifikuojančių dirbtinių neuronų mokymui naudota 100 epochų, kadangi tada klasifikavimo tikslumas jau būna pasiekęs aukščiausią galimą reikšmę.

Irisus klasifikuojančiam neurono klasifikavimo tikslumas aukščiausias, o paklaidos reikšmė mažiausia, kai mokymosi greitis lygus 0,0015.

Krūties vėžio navikus klasifikuojančiam neurono paklaidos reikšmė mažiausia, kai mokymosi greitis lygus 0,0055. Tačiau, paklaidos reikšmė mažiausia, kai mokymosi greitis yra 0,00028. Atsižvelgiant į tai, kad klasifikavimo tikslumas su mokymosi greičiu lygiu 0,00028 mažesnis tik per 0,009756, optimalus santykis tarp mažiausios paklaidos reikšmės ir aukščiausio klasifikavimo tikslo gaunamas su mokymosi greičiu, kuris lygus 0,00028.

1 lentelė. Irisus klasifikuojančio dirbtinio neurono rezultatai

Gauti svoriai	-0,0956; -0,2692; -0,196; 0,3945; 0,6613
Epochų skaičius	100
Paklaida paskutinėje epochoje mokymo duomenims	4,3894
Klasifikavimo tikslumas paskutinėje epochoje mokymo duomenims	0,9143
Paklaida testavimo duomenims	0,0807
Klasifikavimo tikslumas testavimo duomenims	0,9667

2 lentelė. Krūtų vėžio navikus klasifikuojančio dirbtinio neurono rezultatai

Gauti svoriai	-0,2409; 0,032; 0,0247; 0,0158; 0,011; 0,0113; 0,0423; 0,0171; 0,009; 0,0093
Epochų skaičius	100
Paklaida paskutinėje epochoje mokymo duomenims	17,0885
Klasifikavimo tikslumas paskutinėje epochoje mokymo duomenims	0,9623
Paklaida testavimo duomenims	0,0413
Klasifikavimo tikslumas testavimo duomenims	0,9561

7. Rezultatai ir išvados

Apmokius neuroną klasifikuoti irisus ir krūties vėžio navikus pasiektas 95 % klasifikavimo tikslumas, apskaičiuoti svoriai, su kuriais pasiekiamas aukščiausias tikslumo ir žemiausios paklaidos rezultatas. Dirbtinis neuronas naudoja stochastinį gradientinį nusileidimą ir ADALINE mokymosi taisyklę. Paklaidos reikšmė krenta didėjant epochų skaičiui. Paklaidos reikšmė įgija minimumą tik su tam tikru mokymosi greičiu, o toliau didėja. Klasifikavimo tikslumas taip pat priklauso nuo mokymosi greičio ir krenta mokymosi greičio reikšmei kylant.

Šaltiniai

- [Fis36] R. A. Fisher. *Iris* [UCI Machine Learning Repository]. 1936. Pasiekiamas per DOI: 10.24432/C56C76.
- [Kur21] O. Kurasova. *Skaitmeninis intelektas ir sprendimų priėmimas*. 2021.
- [Wo192] William Wolberg. *Breast Cancer Wisconsin (Original)* [UCI Machine Learning Repository]. 1992. Pasiekiamas per DOI: 10.24432/C5HP4Z.