

OPPIMISTEHTÄVÄ

Koneoppiminen

Hanna Rantanen Jenna Räty

DATA-ANALYYSI JA TEKOÄLYN PERUSTEET Joulukuu 2022

Tietotekniikka Tietoliikennetekniikka ja tietoverkot

1 AIHEEN JA DATAN ESITTELY JA TYÖN TAVOITE

1.1 Kuvaus datasta

Oppimistehtävän datasetti löydettiin Kaggle-sivustolta. Kyseinen datasetti on nimeltään 'Suicidal Behaviours Among Adolescents', ja se tarkastelee tekijöitä, jotka ennustavat nuorten itsetuhoista käyttäytymistä. Dataa on kerätty 27 maasta käyttäen Global School-based Student Health Survey-kyselylomaketta (GSHS).

Datasetti on csv-tiedosto. Rivejä tiedostossa on 107kpl ja sarakkeita 17kpl. Country-sarakkeeseen on tallennettu maa (esim. 'Argentina'), Year-sarakkeessa on vuosiluku väliltä 2010–2018, Age Group-sarakkeessa on ikäryhmä väliltä 13–15 tai 16–17, sekä Sex-sarakkeessa sukupuoli ('Male' tai 'Female'). Lopuissa sarakkeista, jotka kuvaavat itsetuhoisiin yrityksiin vaikuttavia tekijöitä sekä itsetuhoisten yritysten esiintymistä, arvot on esitetty prosenttilukuina kerätystä aineistosta. Numeroarvot eivät siis esitä ilmiöiden absoluuttista määrää, vaan prosentuaalista osuutta kyseisessä maassa kyseisenä vuonna tietyssä ikäryhmässä ja sukupuolessa esiintyneistä arvoista.

Datasetissä on prosenttilukuja, jotka saavat jatkuvia arvoja. Näin ollen koneoppimisen menetelmäksi oli valittava regressiomenetelmä, luokittelumenetelmiä ei tällä datasetillä voida käyttää. Koneoppimisen menetelmiksi valittiin lineaarinen regressio, Decision Tree-regressio sekä Random Forest-regressio.

1.2 Kuvaus datan esikäsittelystä

Ensiksi haluttiin tarkistaa, onko datasetissä NaN-arvoja. Tämä katsottiin isnull()-funktiolla. Tästä saatiin tieto, että aineistossa on 6 NaN-arvoa. Nämä arvot muutettiin nollaksi fillna()-funktiolla. Tämän jälkeen tarkistettiin jälleen isnull():n avulla, että NaN-arvoja ei enää ole.

Describe()-funktiolla aineistoa tarkasteltaessa havaittiin, että Currently_Drink_Alcoholsarakkeessa on virheellinen prosenttiluku 548%. Tässä on varmaankin tarkoitus olla arvo 54.8%. Arvo korjattiin loc-toiminnon avulla.

Ennen koneoppimisen menetelmien käyttöä tuli päättää, mitkä sarakkeet datasetistä valitaan opetusvaiheeseen, eli talletettiin X-muuttujaan. Datasetissä oli 17 saraketta, mikä on melko paljon, joten oletuksena oli, ettei niitä kaikkia välttämättä tarvitse käyttää parhaiden oppimistulosten saavuttamiseksi. Tässä vaiheessa kuitenkin jätettiin kaikki sarakkeet X-muuttujaan.

Datasetissä on kategorisia muuttujia, jotka täytyi muuttaa koneoppimisen ymmärtämään muotoon. Tämä tarkoittaa numeerista muotoa, eli luotiin niistä dummymuuttujat. Datasetin kategoriset muuttujat olivat sarakkeissa 'Country', 'Age Group' sekä 'Sex'. Dummy-muuttujien luomiseen käytettiin ColumnTransformer- sekä OneHotEncoder-tekniikoita.

2 DATAN KÄSITTELYN KUVAUS KONEOPPIMISEN MENETELMÄLLÄ

2.1 Datan käsittely lineaarisen regression menetelmällä

X- ja y-muuttujien tallettamisen sekä dummy-muuttujien luomisen jälkeen data jaettiin opetus- ja testidataan. Tässä käytettiin suhdetta 80%-20%, sillä tämä oli kurssilla yleisimmin käytössä ollut tapa jakaa data. Ennen lineaarisen regression mallin opettamista data skaalattiin StandardScaler()-funktion avulla. Malli opetettiin LinearRegression()-funktion avulla, ja ennustus testidatalla tehtiin model.predict()-funktiolla.

Tämän jälkeen mallia voitiin arvioida metriikoiden avulla. Metriikoiksi valittiin R2, MAE sekä RMSE. Opetus- ja ennustusvaiheen jälkeen lineaarisen regression malli sai seuraavat metriikat:

Linear Regression: R2: 0.40400815310106764 MAE: 2.883220826427646 RMSE: 3.9026459653405694

R2-scoresta nähdään, että mallin avulla voidaan selittää vain noin 40% havainnoista. Mitä lähempänä R2-score olisi arvoa 1, sen paremmin malli sopii. Tästä päätellen voisi olla hyödyllistä pyrkiä muokkaamaan opetusprosessia R2-scoren parantamiseksi.

Yrityksen ja erehdyksen taktiikalla sekä koodia useaan kertaan ajamalla saatiin optimoitua X-muuttujaan talletetut sarakkeet parhaan R2-scoren tuottaviin valintoihin. Tällöin X-muuttujasta jätettiin pois sarakkeet 'No_close_friends' sekä 'Really_Get_Drunk'

Linear Regression: R2: 0.45368968517499664 MAE: 2.773112323256568 RMSE: 3.736445880749801

R2-score parani hieman, mutta on edelleen melko vaatimaton. Seuraava keino sen parantamiseen oli opetus- ja testidatan jakosuhteen muuttaminen. Internet-lähteiden perusteella yleisimmät datan jakosuhteet ovat 70%-30%, 80%-20% sekä 90%-10%. Jakosuhteella 70%-30% R2-score heikkeni huomattavasti, arvoon 0.057. Jakosuhtee 90%-10% taas R2-scoreksi saatiin 0.65, mihin voitiin olla tyytyväisiä. Myös MAE- ja RMSE-metriikoiden arvot paranivat näiden muutosten myötä.

Linear Regression: R2: 0.6501257557107862 MAE: 2.5284690323253254 RMSE: 2.964539430303372 Mallin onnistumista voitiin tarkastella myös vertaamalla y_test-muuttujan alkuperäisiä arvoja sekä mallin ennustamia y_pred-muuttujan arvoja. Alla olevasta kuvasta nähdään, että joidenkin ennustusten kohdalla oli enemmän vaihtelua, mutta suuruusluokaltaan ennustukset osuivat kutakuinkin lähelle totuutta.



Lineaarisen regression mallia testattiin myös ajamalla sen läpi itse luotua testiaineistoa. Tätä varten luotiin kolme testitapausta: ensimmäisessä tapauksessa selittävien muuttujien arvot saivat keskiluokkaisia arvoja, toisessa tapauksessa selittävät muuttujat saivat suuria arvoja (poislukien 'Have_Understanding_Parents'-sarake, jolle annettiin pieni arvo), ja kolmannessa tapauksessa selittävät muuttujat saivat pieniä arvoja (poislukien 'Have_Understanding_Parents'-sarake, jolle annettiin suuri arvo). Luodut testitapaukset on esitelty vielä alla.

```
#Tehdään kolme testitapausta

new_cases = [{'Country':'Argentina', 'Year':2017, 'Age Group':'16-17', 'Sex':'Female', 'Currently_Drink_Alcohol':20, 'Overwieght':10,

'Use_Marijuana':30, 'Have_Understanding_Parents':40, 'Missed_classes_without_permssion':25, 'Had_sexual_relation':55,

'Smoke_cig_currently':25, 'Had_fights':15, 'Bullied':30, 'Got_Seriously_injured':43},

{'Country':'Samoa', 'Year':2011, 'Age Group':'13-15', 'Sex':'Male', 'Currently_Drink_Alcohol':90, 'Overwieght':80,

'Use_Marijuana':70, 'Have_Understanding_Parents':20, 'Missed_classes_without_permssion':60, 'Had_sexual_relation':50,

'Smoke_cig_currently':80, 'Had_fights':80, 'Bullied':90, 'Got_Seriously_injured':88},

{'Country':'Indonesta', 'Year':2015, 'Age Group':'16-17', 'Sex':'Female', 'Currently_Drink_Alcohol':0, 'Overwieght':0,

'Use_Marijuana':0, 'Have_Understanding_Parents':90, 'Missed_classes_without_permssion':0, 'Had_sexual_relation':0,

'Smoke_cig_currently':0, 'Had_fights':0, 'Bullied':0, 'Got_Seriously_injured':0}]
```

Kun näiden testitapausten perusteella tehtiin ennustus luodulla mallilla, tuloksista nähtiin, että keskivertotapauksen todennäköisyys itsetuhoiselle yritykselle olisi noin 34%, suuren riskin tapauksen todennäköisyys olisi noin 66% ja pienen riskin tapauksen noin -18%. Tämä tulos on järkevä ja tukee käsitystä mallin toimivuudesta.



2.2 Datan käsittely DecisionTree- ja RandomForest-menetelmillä

DecisionTree- ja RandomForest-menetelmien osalta käytettiin regressiota. Aiemmassa vaiheessa oli jo toteutettu X- ja y-muuttujien luominen, opetus- ja testidataan jako sekä dummy-muuttujien luonti. Myös DecisionTree- ja RandomForest-menetelmiä varten data skaalattiin StandardScaler()-funktion avulla. Koska tässä vaiheessa käytettiin kahta menetelmää, päätettiin ne ajaa aliohjelmien avulla. DecisionTree-menetelmässä malli opetettiin DecisionTreeRegressor()-funktiolla ja RandomForest-menetelmässä RandomForestRegressor()-funktiolla. Aliohjelmat suorittamalla ne palauttivat alla esitellyt metriikat.

Decision Tree:
R2: 0.05280850865491982
MAE: 3.801002164502164
RMSE: 4.91991829015039

Random Forest:
R2: 0.49639177438497895
MAE: 2.819909090909091
RMSE: 3.5874462555366082

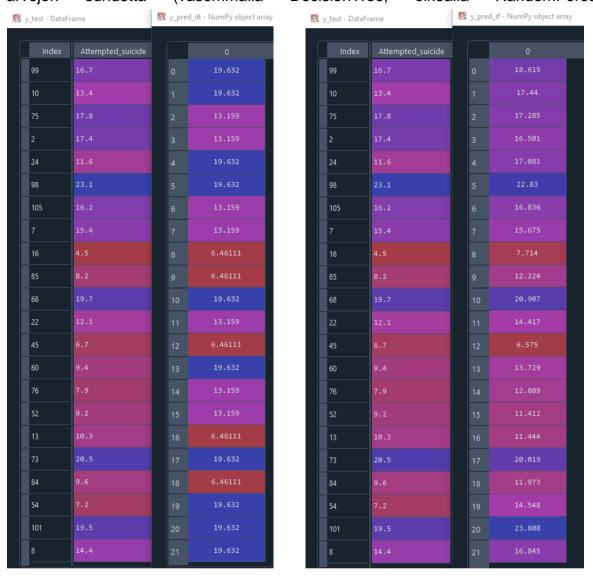
Ylläolevasta kuvasta nähdään, että DecisionTree-malli antoi hyvin heikon R2-scoren, 0.05. RandomForest-malli antoi R2-scoreksi 0.496, mikä on kohtalainen mutta ei vielä kovin vahva todiste mallin toimivuudesta. Tämän parantamiseksi päätettiin yrittää muuttaa X-muuttujan arvoja sekä opetus- ja testidatan jakosuhdetta. X-muuttujan arvoihin lisättiin takaisin ne sarakkeet, jotka lineaarisen regression mallia varten poistettiin. Datan jakosuhde vaihdettiin alkuperäiseen 80%-20%. Tämän myötä metriikat olivat seuraavat:

Decision Tree:
R2: 0.03748913764323869
MAE: 3.862986013986013
RMSE: 4.959544778208524

Random Forest:
R2: 0.6114412883185596
MAE: 2.4757727272727266
RMSE: 3.151137262236373

DecisionTree-mallin R2-score on edelleen heikko, mutta RandomForest-mallin R2-scorea saatiin selkeästi parannettua (0.496 → 0.611).

Alla olevissa kuvissa on esitelty molempien menetelmien ennustusten ja todellisten y:n arvojen suhdetta (vasemmalla DecisionTree, oikealla RandomForest).



Silmämääräisesti vertailemalla voidaan todeta, että RandomForest-mallin ennustuksien virheet ovat pienempiä kuin DecisionTree-mallilla.

2.3 Saatujen tulosten pohdinta ja vertailu

Työssä käytettiin kolmea koneoppimisen menetelmää: lineaarinen regressio, DecisionTree-regressio sekä RandomForest-regressio.

Linear Regression: R2: 0.6501257557107862 MAE: 2.5284690323253254 RMSE: 2.964539430303372 Decision Tree: R2: 0.03748913764323869 MAE: 3.862986013986013 RMSE: 4.959544778208524

Random Forest:

R2: 0.6114412883185596 MAE: 2.4757727272727266 RMSE: 3.151137262236373 Yllä olevissa kuvissa on esitelty eri koneoppimisen mallien metriikat. Ensimmäiseksi tarkasteltiin R2-scorea. Se kertoo, kuinka hyvin malli sopii dataan, ja sen arvo tulisi olla mahdollisimman lähellä arvoa 1. R2-scoren arvoja tarkastelemalla voidaan todeta, että lineaarisen regression menetelmä toimii parhaiten (R2=0.65). Seuraavaksi parhaiten toimiva malli saatiin RandomForest-menetelmällä (R2=0.61). DecisionTreemenetelmä toimi tässä yhteydessä heikosti (R2=0.04).

Seuraavaksi voitiin tarkastella MAE- ja RMSE-arvojen suhteita. Mitä suurempi on näiden ero, sitä suurempaa vaihtelua on virheissä. Näin ollen RMSE- ja MAE-arvon erotus tulisi olla mahdollisimman pieni. Lineaarisen regression menetelmässä arvojen erotus on 0.436, DecisionTree-menetelmällä 1.0965 ja RandomForest-menetelmässä 0.6753. Näin ollen voidaan todeta, että lineaarisen regression menetelmässä virheiden vaihtelu on pienempää kuin muilla menetelmillä.

Tulosten perusteella voitiin vetää yhteen, että tässä oppimistehtävässä parhaiten koneoppimisen menetelmistä suoriutui lineaarinen regressio. RandomForestmenetelmä taas toimi paremmin kuin DecisionTree, mikä on järkevää, sillä se hyödyntää toiminnassaan useita DecisionTree-menetelmän luomia päätöspuita. Lineaarinen regressio saattoi tässä tilanteessa toimia RandomForestia paremmin siksi, että RandomForest hyödyntää toiminnassaan myös luokittelun menetelmiä, jotka eivät välttämättä sopineet valittuun datasettiin ja sen avulla opettamiseen. Tämän lisäksi RandomForest saattaa olla yksinkertaisesti liian monimutkainen ja raskas menetelmä koneoppimiseen tämän aineiston kohdalla, jolloin "kevyempi" lineaarinen regressio tuo paremmat tulokset.

3 TYÖN ARVIOINTI

Tekijöiden mielestä tässä oppimistehtävässä on käytetty laajasti kurssin koneoppimisosuudessa opittuja regressioon soveltuvia tekniikoita ja menetelmiä. Kurssin harjoitustehtävien lisäksi työssä käytettiin lähteenä internetiä.

Kehityskohteena tämän työn myötä voisi olla sellaisen datasetin valitseminen, joka soveltuisi luokittelevien koneoppimismenetelmien käyttöön. Tällöin tulisi harjoiteltua myös luokittelevien menetelmien käyttöä sekä niihin liittyviä visualisointimenetelmiä, kuten confusion matrixia. Tämän lisäksi valittu datasetti olisi voinut olla hieman laajempi, jotta metriikka-arvot olisivat mahdollisesti olleet paremmat ja luodut mallit siten luotettavampia. Suuremman datamäärän avulla olisi mahdollisesti voitu hvödyntää myös neuroverkkoien menetelmiä.

Onnistumisena tämän työn osalta oli kolmen tarkoituksenmukaisen koneoppimismenetelmän käyttö. Menetelmiä käyttämällä saatiin selkeästi vertailtua niiden toimivuutta hyödyntäen muun muassa erilaisia metriikoita. Myös koneoppimista edeltävä datan ja sen esikäsittelyn kuvaus oli kattava ja perusteellinen. Oppimistehtävässä pyrittiin hiomaan lopputulos mahdollisimman hyväksi muokkaamalla mallille syötettäviä arvoja, kuten X-muuttujaa sekä datan jakosuhdetta. Oppimistehtävän raportti on huolellisesti ja johdonmukaisesti laadittu, ja se etenee selkeästi havainnollistavien kuvien ja sanallisen kuvauksen myötä. Oppimistehtävä koettiin tekijöiden toimesta onnistuneeksi ja se sisältää 2 pisteen työlle vaadittavat ominaisuudet.

4 LIITTEET

Liite 1. Otos datasta

Index ▼	Country	Voor	Ago Group	Sav	Currently Drink Alcohol	Poolly Get Drunk	Overvieght	Lise Marijuana	Have Understanding Parents	s Missed_classes_without_permssion	Had sexual relation	Smoke sig current	Had fights	Pullind	Sot Seriously injured	No close friends	Attempted suicide
0	Argentina			Female			27.8	7.9	41.5	24.7	25.7	16.8	17.2	0	27.5	4.8	19.9
,	Argentina				44.9												
2				Female													
3					68.1												
4					49.3												
5				Female													
6				Male													
7	Barabados		13-15	Female													
8																	
9				Female													
10	Benin			Male													
11	Benin			Female													
12	Bhutan			Male													
13	Bhutan			Female													
14	Bhutan			Male													
15	Bhutan			Female													
16	Brunei Darussalam			Male													
17	Brunei Darussalam			Female													
18	Brunei Darussalam			Male													
19	Brunei Darussalam			Female													
20				Male													
21	Dominican Republic			Female													
22	Dominican Republic			Male													
23	Dominican Republic			Female													
24				Male													
25				Female													
26	E141			Male													
27	E141			Female													
28	Indonesia			Male													
29	Indonesia			Female													
30	Indonesia	2015		Male													
31	Indonesia	2015		Female													
32	Jamaica			Male	54.8												
33	Jamaica			Female													
34	Jamaica			Male													
35	Jamaica			Female													
36	Kiribati			Male													
37	Kiribati	2011		Female													
38		2015		Male													
39	Laos	2015		Female													
40		2015		Male													
41	Laos	2015		Female													
42	Malaysia	2012		Male													
43	Malaysia	2012		Female													
44	Malaysia	2012		Male													
45	Malaysia	2012		Female													
46	Mauritus			Male													
47	Mauritus			Female													

Liite 2. Python-koodi

.....

Data-analyysi ja tekoälyn perusteet Oppimistehtävä 2 Hanna Rantanen & Jenna Räty 201227

.....

import pandas as pd import seaborn as sns import matplotlib.pyplot as plt import numpy as np import scipy.stats as stats from scipy.stats.stats import pearsonr from scipy.stats import chi2_contingency

 $from \ sklearn. linear_model \ import \ Linear Regression$

from sklearn.model_selection import train_test_split

from sklearn.metrics import r2_score, mean_absolute_error, mean_squared_error

from sklearn.preprocessing import StandardScaler

from sklearn.compose import ColumnTransformer

from sklearn.preprocessing import OneHotEncoder

from sklearn.tree import DecisionTreeRegressor

from sklearn.ensemble import RandomForestRegressor

from tensorflow.keras.models import Sequential

from tensorflow.keras.layers import Dense, Dropout

#Luetaan aineisto dataframeen
df = pd.read_csv('GHSH_Pooled_Data1.csv')

DATAN ESIKÄSITTELY

#Tarkistetaan onko NaN-arvoja
print (df.isnull().values.any())

#Katsotaan NaN-arvojen lukumäärä print (df.isnull().sum().sum())

#Korvataan NaN-arvot nollalla df = df.fillna(0)

#Tarkistetaan onko NaN-arvoja
print (df.isnull().values.any())

#Katsotaan NaN-arvojen lukumäärä print (df.isnull().sum().sum())

```
#Korvataan virheellinen prosenttilukema
df.loc[df['Currently_Drink_Alcohol'] == 548, 'Currently_Drink_Alcohol'] = '54.8'
```

LINEAARINEN REGRESSIO

```
#Luodaan X- ja y-muuttujat
X = df.loc[:, ['Country', 'Year', 'Age Group', 'Sex', 'Currently_Drink_Alcohol', 'Overwieght',
       'Use Marijuana',
                            'Have Understanding Parents',
                                                                'Missed classes without permssion',
       'Had sexual relation', 'Smoke cig currently', 'Had fights', 'Bullied', 'Got Seriously injured']]
y = df.loc[:, ['Attempted suicide']]
# Poistetut sarakkeet: 'No close friends' 'Really Get Drunk' --> saatu R2-score optimoitua
#Dummy-muuttujien luominen Country-, Age Group- ja Sex-sarakkeista
X_{orig} = X
ct = ColumnTransformer(transformers=[('encoder', OneHotEncoder(drop='first'), ['Country', 'Age
       Group', 'Sex'])], remainder='passthrough')
X = ct.fit transform(X)
#Jaetaan aineisto opetus- ja testidataan (90% & 10%)
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.1, random_state=0)
#Skaalaus
scaler x = StandardScaler(with_mean=False)
X_train = scaler_x.fit_transform(X_train)
X test = scaler x.transform(X test)
#Opetetaan lineaarisen regression malli
model = LinearRegression()
model.fit(X_train, y_train)
#Ennustetaan testidatalla
y_pred = model.predict(X_test)
#Arvioidaan malli käyttäen metriikoita
R2 = r2 score(y test, y pred)
mae = mean_absolute_error(y_test, y_pred)
mse = mean_squared_error(y_test, y_pred)
rmse = np.sqrt(mse)
print (f'\nLinear Regression:\nR2: {R2}')
print (f'MAE: {mae}')
print (f'RMSE: {rmse}')
```

```
#Tehdään kolme testitapausta: keskiverto, suuri riski, pieni riski
new cases
             =
                  [{'Country':'Argentina',
                                            'Year':2017,
                                                           'Age
                                                                  Group':'16-17',
                                                                                    'Sex':'Female',
       'Currently Drink Alcohol':20, 'Overwieght':10, 'Use Marijuana':30, 'Have Understanding
       Parents':40, 'Missed classes without permssion':25, 'Had sexual relation':55, 'Smoke cig
       currently':25, 'Had fights':15, 'Bullied':30, 'Got Seriously injured':43},
{'Country':'Samoa', 'Year':2011, 'Age Group':'13-15', 'Sex':'Male', 'Currently_Drink Alcohol':90,
       'Overwieght':80, 'Use Marijuana':70, 'Have Understanding Parents':20, 'Missed classes
       without permssion':60, 'Had sexual relation':50, 'Smoke cig currently':80, 'Had fights':80,
       'Bullied':90, 'Got Seriously injured':80},
{'Country':'Indonesia', 'Year':2015, 'Age Group':'16-17', 'Sex':'Female', 'Currently Drink Alcohol': 0,
                       'Use Marijuana':0,
                                            'Have Understanding Parents':90,
       'Overwieght':0,
                                                                                 'Missed classes
       without permssion':0,
                               'Had sexual relation':0,
                                                          'Smoke cig currently':0,
                                                                                    'Had fights':0,
       'Bullied':0, 'Got Seriously injured':0}]
#Talletetaan testitapauksen data, tehdään dummyt sekä skaalaus
new_data = pd.DataFrame(new_cases)
new data = ct.transform(new data)
new_data = scaler_x.transform(new_data)
#Ennustetaan kolmen testitapauksen perusteella
new y = pd.DataFrame(model.predict(new data))
###### DECISIONTREE- & RANDOMFOREST-REGRESSIO (TOTEUTUS ALIOHJELMILLA) ##############
#Datan skaalaus
scaler_x = StandardScaler(with_mean=False)
X train = scaler x.fit transform(X train)
X_test = scaler_x.transform(X_test)
scaler y = StandardScaler()
y train = scaler y.fit transform(y train)
#DecisionTree regressio-mallin aliohjelma
def decisionTree(max depth, min samples split):
#Opetetaan DecisionTree regressio-malli opetusdatalla
model dt = DecisionTreeRegressor(max depth=max depth, min samples split=min samples split)
model_dt.fit(X_train, y_train)
#Ennustetaan testidatalla
y_pred_dt = scaler_y.inverse_transform(model_dt.predict(X_test).reshape(-1,1))
#Arvioidaan malli käyttäen metriikoita
R2 = r2 score(y test, y pred dt)
mae = mean_absolute_error(y_test, y_pred_dt)
mse = mean_squared_error(y_test, y_pred_dt)
rmse = np.sqrt(mse)
```

return R2, mae, rmse, y pred dt

```
#RandomForest regressio-mallin aliohjelma
def randomForest(n_estimators):
#Opetetaan RandomForest regressio-malli opetusdatalla
model_rf = RandomForestRegressor(n_estimators=n_estimators)
model_rf.fit(X_train, y_train.ravel())
#Ennustetaan testidatalla
y_pred_rf = scaler_y.inverse_transform(model_rf.predict(X_test).reshape(-1,1))
#Arvioidaan malli käyttäen metriikoita
R2 = r2_score(y_test, y_pred_rf)
mae = mean_absolute_error(y_test, y_pred_rf)
mse = mean_squared_error(y_test, y_pred_rf)
rmse = np.sqrt(mse)
return R2, mae, rmse, y_pred_rf
#Palautetaan metriikat DecisionTreelle
R2, mae, rmse, y_pred_dt = decisionTree(10,50)
print (f'\nDecision Tree:\nR2: {R2}')
print (f'MAE: {mae}')
print (f'RMSE: {rmse}')
#Palautetaan metriikat RandomForestille
R2, mae, rmse, y pred rf = randomForest(100)
print (f'\nRandom Forest:\nR2: {R2}')
print (f'MAE: {mae}')
print (f'RMSE: {rmse}')
```