

## Exercise 2

Christof Pegrisch, Simon König, Philipp Eberl

December 18, 2021

# Table of Contents

Datasets

Preprocessing

Approach

Implementation

Evaluation metrics

Results using own implementation

Results using AutoML (TPOT)

Results using AutoML (Hyperopt-sklearn)

Conclusion

# Dataset 1: Solar flares

- ▶ Divided into two sections, one with more error correction
- ▶ 10 Attributes
- ▶ 1389 Instances
- ▶ No missing values
- ▶ Instances describe state of certain region of the sun
- ▶ <https://archive-beta.ics.uci.edu/ml/datasets/solar+flare>

## Dataset 2: Wine quality

- ▶ Two separate files for red and white wine
- ▶ Information about wine like citric acid or residual sugar
- ▶ 12 attributes each
- ▶ 4898 instances for white wine
- ▶ 1599 instances for red wine
- ▶ No missing values
- ▶ <https://www.kaggle.com/brendan45774/wine-quality?select=winequality-red.csv>

## Dataset 3: Coronavirus

- ▶ Dataset with population and vaccination data of countries
- ▶ 10 attributes
- ▶ About 24.000 instances with different timestamps
- ▶ No missing values
- ▶ <https://www.kaggle.com/sinakaraji/covid-vaccination-vs-death/activity>

# Table of Contents

Datasets

Preprocessing

Approach

Implementation

Evaluation metrics

Results using own implementation

Results using AutoML (TPOT)

Results using AutoML (Hyperopt-sklearn)

Conclusion

# What has been done

- ▶ Solar flares
  - ▶ Converting to numeric (requirement of used implementations)
  - ▶ One hot encoding (for region class, largest spot and spot distribution)
- ▶ Wine quality
  - ▶ Converting to numeric
- ▶ Coronavirus
  - ▶ Converting to numeric
  - ▶ Normalization by dividing through population (eg ratio of people vaccinated)
  - ▶ One hot encoding (for country names)

# Table of Contents

Datasets

Preprocessing

Approach

Implementation

Evaluation metrics

Results using own implementation

Results using AutoML (TPOT)

Results using AutoML (Hyperopt-sklearn)

Conclusion



- ▶ Implementation of both hillclimbing and simulated annealing due to their similarities.
- ▶ Simulated annealing for:
  - ▶ Solar flares
  - ▶ Wine white
  - ▶ Wine red
- ▶ Hillclimbing for:
  - ▶ Covid  
(due to the size of the dataset and shorter runtime of hillclimbing)

# Used Algorithms

- ▶ Hillclimbing (finds local maximum by searching the immediate neighborhood)
  - ▶ Starts with a random set of hyperparameters
  - ▶ Searches close neighborhood for better solution and keeps it
  - ▶ Repeats until no solution in the neighborhood is better than the current
- ▶ Simulated annealing (tries to break out of local minima using probabilities)
  - ▶ Starts with a random set of hyperparameters
  - ▶ Takes random solution from the close neighborhood
  - ▶ Keeps the new solution according to a certain probability dependent on the quality and a decreasing temperature value
  - ▶ Repeats until either no changes occur for a certain amount of epochs or until a maximum number of epochs is reached

# Chosen Regressors

- ▶ Linear SVR

https:

`//scikit-learn.org/stable/modules/generated/  
sklearn.svm.LinearSVR.html#sklearn.svm.LinearSVR`

- ▶ K-Neighbours Regressor

https:

`//scikit-learn.org/stable/modules/generated/  
sklearn.neighbors.KNeighborsRegressor.html`

- ▶ DecisionTree Regressor

`https://scikit-learn.org/stable/modules/  
generated/sklearn.tree.DecisionTreeRegressor.html`

# Table of Contents

Datasets

Preprocessing

Approach

**Implementation**

Evaluation metrics

Results using own implementation

Results using AutoML (TPOT)

Results using AutoML (Hyperopt-sklearn)

Conclusion

Super vector regression using a linear kernel. By allowing an error in the model a support vector regressor finds the optimal hyperplane to fit the data.

- ▶ Strengths: Effective for high dimensional feature spaces, memory efficient
- ▶ Weaknesses: Not suitable for large datasets, hard to understand parameter selection (Black Box)

# K-Neighbours Regressor

Uses K-nearest-neighbours algorithm to perform classification. To make a prediction the regressor takes the mean of k nearest neighbours.

- ▶ Strengths: Effective for large datasets, robust to noisy data, fast model building (Lazy learner)
- ▶ Weaknesses: Prediction is computationally costly, hard to determine which attributes contribute to regression

# DecisionTree Regressor

Builds a decision tree with real values as leave nodes. In the model building process mean squared error is used to split node in sub-nodes.

- ▶ Strengths: Easy to interpret and visualise, little influence by outliers, useful in data exploration
- ▶ Weaknesses: Cannot extrapolate, tends to overfit

# Table of Contents

Datasets

Preprocessing

Approach

Implementation

**Evaluation metrics**

Results using own implementation

Results using AutoML (TPOT)

Results using AutoML (Hyperopt-sklearn)

Conclusion



- ▶ Negative Mean squared error
  - ▶ Describes how close a regression line is to a specific set of points. Lower scores indicate a better fit
- ▶  $R^2$  score
  - ▶ R-squared explains to what extent the variance of one variable explains the variance of the second variable. Best result is 1.0, which means all of the variation can be explained by the models inputs.

# Table of Contents

Datasets

Preprocessing

Approach

Implementation

Evaluation metrics

Results using own implementation

Results using AutoML (TPOT)

Results using AutoML (Hyperopt-sklearn)

Conclusion

- ▶ Covid dataset
  - ▶ Result negative MSE:  $-1822.8864578877833$
  - ▶ Hyperparameters: (loss="epsilon\_insensitive", tol=1e-5, C=10, epsilon=1e-3)
  - ▶ Result  $R^2$ :  $-1.6148099403159493$
  - ▶ Hyperparameters: (loss="epsilon\_insensitive", tol=1e-5, C=1e-4, epsilon=1)
- ▶ Red wine dataset
  - ▶ Result negative MSE:  $-0.651999565611231$
  - ▶ Hyperparameters: (loss="epsilon\_insensitive", tol=1e-5, C=1e-4, epsilon=1e-3)
  - ▶ Result  $R^2$ :  $0.19967708634304884$
  - ▶ Hyperparameters: (loss="epsilon\_insensitive", tol=1e-3, C=1e-2, epsilon=1)

- ▶ White wine dataset
  - ▶ Result negative MSE:  $-0.9578209077406313$
  - ▶ Hyperparameters: (loss="epsilon\_insensitive", tol=1e-5, C=1, epsilon=1)
  - ▶ Result  $R^2$ : 0.19560956565824933
  - ▶ Hyperparameters: (loss="epsilon\_insensitive", tol=1e-5, C=1e-2, epsilon=1e-3)
- ▶ Solar flares dataset
  - ▶ Result negative MSE:  $-0.5772637203438242$
  - ▶ Hyperparameters: (loss="epsilon\_insensitive", tol=1e-1, C=10, epsilon=1e-1)
  - ▶ Result  $R^2$ : 0.010585959265640965
  - ▶ Hyperparameters: (loss="epsilon\_insensitive", tol=1e-5, C=1e-2, epsilon=1e-1)

# K-Neighbours Regressor

- ▶ Covid dataset
  - ▶ Result negative MSE:  $-528.5757039291524$
  - ▶ Hyperparameters: (n\_neighbors=4, weights="uniform", p=2)
  - ▶ Result  $R^2$ : 0.23799133614878842
  - ▶ Hyperparameters: (n\_neighbors=10, weights="distance", p=1)
- ▶ Red wine dataset
  - ▶ Result negative MSE:  $-0.5745512186731103$
  - ▶ Hyperparameters: (n\_neighbors=10, weights="distance", p=1)
  - ▶ Result  $R^2$ : 0.029412732130088725
  - ▶ Hyperparameters: (n\_neighbors=6, weights="distance", p=1)

# K-Neighbours Regressor

- ▶ White wine dataset
  - ▶ Result negative MSE:  $-0.7485020711918213$
  - ▶ Hyperparameters: (n\_neighbors=7, weights="uniform", p=1)
  - ▶ Result  $R^2$ : 0.01870206497013156
  - ▶ Hyperparameters: (n\_neighbors=6, weights="uniform", p=1)
- ▶ Solar flares dataset
  - ▶ Result negative MSE:  $-0.614328428953588$
  - ▶ Hyperparameters: (n\_neighbors=4, weights="uniform", p=2)
  - ▶ Result  $R^2$ :  $-0.06532347561422938$
  - ▶ Hyperparameters: (n\_neighbors=9, weights="uniform", p=1)

# Decision Tree Regressor

- ▶ Covid dataset
  - ▶ Result negative MSE:  $-47.644739644785794$
  - ▶ Hyperparameters: (max\_depth=2, min\_samples\_split=5, min\_samples\_leaf=6)
  - ▶ Result  $R^2$ : 0.9914236829965463
  - ▶ Hyperparameters: (max\_depth=4, min\_samples\_split=10, min\_samples\_leaf=2)
- ▶ Red wine dataset
  - ▶ Result negative MSE:  $-0.4948440103954347$
  - ▶ Hyperparameters: (max\_depth=6, min\_samples\_split=10, min\_samples\_leaf=10)
  - ▶ Result  $R^2$ : 0.2015905613389884
  - ▶ Hyperparameters: (max\_depth=6, min\_samples\_split=10, min\_samples\_leaf=10)

# Decision Tree Regressor

- ▶ White wine dataset
  - ▶ Result negative MSE:  $-0.5860529460556473$
  - ▶ Hyperparameters: (max\_depth=5, min\_samples\_split=7, min\_samples\_leaf=4)
  - ▶ Result  $R^2$ : 0.2402569677761802
  - ▶ Hyperparameters: (max\_depth=4, min\_samples\_split=2, min\_samples\_leaf=10)
- ▶ Solar flares dataset
  - ▶ Result negative MSE:  $-0.576831949429341$
  - ▶ Hyperparameters: (max\_depth=9, min\_samples\_split=6, min\_samples\_leaf=10)
  - ▶ Result  $R^2$ :  $-0.2005731727387947$
  - ▶ Hyperparameters: (max\_depth=4, min\_samples\_split=3, min\_samples\_leaf=9)



# Table of Contents

Datasets

Preprocessing

Approach

Implementation

Evaluation metrics

Results using own implementation

Results using AutoML (TPOT)

Results using AutoML (Hyperopt-sklearn)

Conclusion

- ▶ Covid dataset
  - ▶ Result negative MSE:  $-676.798302552437$
  - ▶ Hyperparameters: (loss="squared\_epsilon\_insensitive", tol=1e-05, C=1, epsilon=1)
  - ▶ Result  $R^2$ : 0.03303043368725276
  - ▶ Hyperparameters: (loss="squared\_epsilon\_insensitive", tol=1e-05, C=20, epsilon=1)
- ▶ Red wine dataset
  - ▶ Result negative MSE:  $-0.4336334344973999$
  - ▶ Hyperparameters: (loss="squared\_epsilon\_insensitive", tol=1e-05, C=20, epsilon=0.001)
  - ▶ Result  $R^2$ : 0.3260756862051558
  - ▶ Hyperparameters: (loss="squared\_epsilon\_insensitive", tol=1e-05, C=1, epsilon=0.001)

- ▶ White wine dataset
  - ▶ Result negative MSE:  $-0.5781509600594222$
  - ▶ Hyperparameters: (loss="squared\_epsilon\_insensitive", tol=1e-05, C=1, epsilon=0.1)
  - ▶ Result  $R^2$ : 0.27084918178993517
  - ▶ Hyperparameters: (loss="squared\_epsilon\_insensitive", tol=1e-05, C=20, epsilon=0.001)
- ▶ Solar flares dataset
  - ▶ Result negative MSE:  $-0.5342683423739354$
  - ▶ Hyperparameters: (loss="squared\_epsilon\_insensitive", tol=0.1, C=0.01, epsilon=0.001)
  - ▶ Result  $R^2$ : 0.15284999357211257
  - ▶ Hyperparameters: (loss="squared\_epsilon\_insensitive", tol=0.001, C=0.01, epsilon=0.001)

# K-Neighbours Regressor

- ▶ Covid dataset
  - ▶ Result negative MSE:  $-183.02381848495688$
  - ▶ Hyperparameters: (n\_neighbors=2, weights="distance", p=2)
  - ▶ Result  $R^2$ : 0.7384318978131217
  - ▶ Hyperparameters: (n\_neighbors=2, weights="distance", p=2)
- ▶ Red wine dataset
  - ▶ Result negative MSE:  $-0.4736539716577819$
  - ▶ Hyperparameters: (n\_neighbors=9, weights="distance", p=1)
  - ▶ Result  $R^2$ : 0.26291900754678466
  - ▶ Hyperparameters: (n\_neighbors=9, weights="distance", p=1)

# K-Neighbours Regressor

- ▶ White wine dataset
  - ▶ Result negative MSE:  $-0.5383166244166239$
  - ▶ Hyperparameters: (n\_neighbors=9, weights="distance", p=1)
  - ▶ Result  $R^2$ : 0.321785064120657
  - ▶ Hyperparameters: (n\_neighbors=9, weights="distance", p=1)
- ▶ Solar flares dataset
  - ▶ Result negative MSE:  $-0.5606531262445241$
  - ▶ Hyperparameters: (n\_neighbors=9, weights="uniform", p=1)
  - ▶ Result  $R^2$ : 0.11782061039157306
  - ▶ Hyperparameters: (n\_neighbors=9, weights="uniform", p=1)

# Decision Tree Regressor

- ▶ Covid dataset
  - ▶ Result negative MSE:  $-15.57304409288794$
  - ▶ Hyperparameters: (max\_depth=9, min\_samples\_split=8, min\_samples\_leaf=4)
  - ▶ Result  $R^2$ :  $0.973973891378282$
  - ▶ Hyperparameters: (max\_depth=9, min\_samples\_split=9, min\_samples\_leaf=8)
- ▶ Red wine dataset
  - ▶ Result negative MSE:  $-0.45011627543766675$
  - ▶ Hyperparameters: (max\_depth=4, min\_samples\_split=4, min\_samples\_leaf=8)
  - ▶ Result  $R^2$ :  $0.29964107282339586$
  - ▶ Hyperparameters: (max\_depth=4, min\_samples\_split=8, min\_samples\_leaf=9)

# Decision Tree Regressor

- ▶ White wine dataset
  - ▶ Result negative MSE:  $-0.5801390660131089$
  - ▶ Hyperparameters: (max\_depth=5, min\_samples\_split=2, min\_samples\_leaf=8)
  - ▶ Result  $R^2$ :  $0.2673631830667058$
  - ▶ Hyperparameters: (max\_depth=4, min\_samples\_split=4, min\_samples\_leaf=8)
- ▶ Solar flares dataset
  - ▶ Result negative MSE:  $-0.5529494224260217$
  - ▶ Hyperparameters: (max\_depth=2, min\_samples\_split=8, min\_samples\_leaf=7)
  - ▶ Result  $R^2$ :  $0.12554988981936674$
  - ▶ Hyperparameters: (max\_depth=1, min\_samples\_split=9, min\_samples\_leaf=7)

# Table of Contents

Datasets

Preprocessing

Approach

Implementation

Evaluation metrics

Results using own implementation

Results using AutoML (TPOT)

Results using AutoML (Hyperopt-sklearn)

Conclusion



## ► Covid dataset

- Result negative MSE:  $-1343.9272236143343$
- Result  $R^2$ :  $-0.975865899654315$
- Hyperparameters: ('learner': SVR( $C=6.14430240999864e-05$ , cache\_size=512, degree=1, epsilon=471.85790078465106, gamma='auto', kernel='linear', max\_iter=84056886.0, tol=0.00019776043393699547), 'preprocs': (), 'ex\_preprocs': ())

## ► Red wine dataset

- Result negative MSE:  $-0.47382854090840987$
- Result  $R^2$ : 0.23298333807017352
- Hyperparameters: ('learner': SVR( $C=0.021658668215710258$ , cache\_size=512, degree=1, epsilon=0.5767107013240668, gamma='auto', kernel='linear', max\_iter=814470115.0, tol=4.057465654030927e-05), 'preprocs': (), 'ex\_preprocs': ())

- ▶ White wine dataset
  - ▶ Result negative MSE:  $-0.7742476360259494$
  - ▶ Result  $R^2$ :  $0.0030735412517071347$
  - ▶ Hyperparameters: ('learner': SVR( $C=9.231350836902164e-05$ , cache\_size=512, degree=1, epsilon=0.003526261250305432, gamma='auto', kernel='linear', max\_iter=450281484.0, tol=0.007903935921111719), 'preprocs': (), 'ex\_preprocs': ())
- ▶ Solar flares dataset
  - ▶ Result negative MSE:  $-0.6209445377362287$
  - ▶ Result  $R^2$ :  $-0.21250010342979167$
  - ▶ Hyperparameters: ('learner': SVR( $C=0.011018759854312025$ , cache\_size=512, degree=1, epsilon=0.4722286042118413, gamma='auto', kernel='linear', max\_iter=422105859.0, shrinking=False, tol=4.828809035580296e-05), 'preprocs': (), 'ex\_preprocs': ())

# K-Neighbours Regressor

- ▶ Covid dataset
  - ▶ Result negative MSE:  $-447.24793185054614$
  - ▶ Result  $R^2$ :  $0.35903633021606923$
  - ▶ Hyperparameters: ('learner':  
KNeighborsRegressor(metric='euclidean', n\_jobs=1,  
n\_neighbors=23), 'preprocs': (), 'ex\_preprocs': ())
- ▶ Red wine dataset
  - ▶ Result negative MSE:  $-0.5461105011420537$
  - ▶ Result  $R^2$ :  $0.12185830259391994$
  - ▶ Hyperparameters: ('learner':  
KNeighborsRegressor(metric='manhattan', n\_jobs=1,  
n\_neighbors=22, p=1, weights='distance'), 'preprocs': (),  
'ex\_preprocs': ())

# K-Neighbours Regressor

- ▶ White wine dataset
  - ▶ Result negative MSE:  $-0.6911960646052259$
  - ▶ Result  $R^2$ :  $0.1087286219856329$
  - ▶ Hyperparameters: ('learner':  
KNeighborsRegressor(metric='euclidean', n\_jobs=1,  
n\_neighbors=34, weights='distance'), 'preprocs': (),  
'ex\_preprocs': ())
- ▶ Solar flares dataset
  - ▶ Result negative MSE:  $-0.5189269288987464$
  - ▶ Result  $R^2$ :  $0.02340265519013567$
  - ▶ Hyperparameters: ('learner':  
KNeighborsRegressor(metric='euclidean', n\_jobs=1,  
n\_neighbors=44), 'preprocs': (), 'ex\_preprocs': ())

# Decision Tree Regressor

- ▶ Decision Tree Regressor could not be tested for Hyperopt-sklearn since it is not yet implemented according to its github page.

<https://github.com/hyperopt/hyperopt-sklearn>

# Table of Contents

Datasets

Preprocessing

Approach

Implementation

Evaluation metrics

Results using own implementation

Results using AutoML (TPOT)

Results using AutoML (Hyperopt-sklearn)

Conclusion

# Interpretation and comparison of results

- ▶ The experiments have shown that both implemented AutoML algorithms performed very well, finding better sets of hyperparameters than they started with. They could however not reach the quality of the hyperparameter-sets of existing implementations. Especially TPOT performed extremely well, outperforming every other algorithm in almost every test.