

Data Mining und Maschinelles Lernen

Prof. Kristian Kersting
Zhongjie Yu
Johannes Czech



TECHNISCHE
UNIVERSITÄT
DARMSTADT

Sommersemester 2021
24. Juli 2021
Bonusübungsblatt 1

Die **Abgabefrist** dieser Bonusübung ist am **31.07.2021** um **23:59 Uhr**.

Aufgabe	1	2	3	4	5
Maximal Punktzahl	12	10	18	12	18
Erreichte Punktzahl					

Gruppe AE	Nachnahme	Vorname	Matrikelnummer
1	Liu	Lanmiao	2571210
2	Shao	Yuening	2504740
3	Li	Zongjian	2595681
4	Liu	Yue	2803140

Benötigte Dateien

Die benötigten Datensatz zu Aufgabe 1 und 2 sowie Skriptvorlagen finden Sie in unserem Moodle-Kurs:

<https://moodle.informatik.tu-darmstadt.de/course/view?id=1058>

Den Datensatz zu Aufgabe 3, 4 und 5 finden Sie auf [kaggle.com](https://www.kaggle.com):

<https://www.kaggle.com/slothkong/10-monkey-species>

Gruppeneinteilung

Bearbeiten Sie diese Übung in Dreier- oder Vierergruppen. Es steht Ihnen frei die Gruppen selbst zu bilden. Nutzen Sie hierfür die Gruppeneinteilung und das Forum in Moodle. Sehen Sie bitte aufgrund der aktuellen Coronakrise davon ab, sich lokal zu treffen und nutzen Sie stattdessen digitale Kommunikationskanäle. Zur gemeinsamen Bearbeitung der Abgabe können sie beispielsweise <https://overleaf.com> nutzen.

Theoretische Aufgaben

Bei theoretischen Übungsaufgaben, sind wir Ihnen sehr dankbar, wenn Sie diese in \LaTeX formatieren und als PDF einreichen. Nutzen Sie hierfür die \LaTeX -Vorlage und die vorgesehene Blöcke:

```
\begin{solution}  
% your solution goes here  
\end{solution}
```

Geben Sie dabei ihre Gruppenmitglieder und Gruppennummer in der Datei `group_members.tex` an.

Wenn Sie mit \LaTeX nicht ausreichend vertraut sind, können Sie auch einen hochauflösenden Scan einer handgeschriebenen Lösung einreichen. Bitte schreiben Sie ordentlich und leserlich.

Programmieraufgaben

Bei Aufgaben, die mit einem `</>` versehen sind, handelt es sich um Programmieraufgaben. Bearbeiten Sie in diesem Fall die vorgegebene Programmiervorlage. Verwenden Sie bevorzugt \geq **Python 3.7**, da wir diese Version zum Testen ihrer Lösung benutzen. Benennen Sie die Funktionsdateien nicht um und ändern Sie die angegebenen Funktionssignaturen nicht. Wenn Sie das Gefühl haben, dass es ein Fehler bei den Zuweisungen, fragen Sie uns auf Moodle.

Formalien zur Abgabe

Bitte laden Sie Ihre Lösungen in der entsprechenden Rubrik auf Moodle hoch. Sie müssen **nur eine Lösung pro Gruppe** einreichen. Wenn Sie keinen Zugang zu Moodle haben, setzen Sie sich bitte so schnell wie möglich mit uns in Verbindung. Laden Sie alle Ihre Lösungsdateien als eine einzige .zip-Datei hoch. Bitte beachten Sie, dass wir keine anderen als die angegebenen Dateiformate akzeptieren. Laden Sie **nicht den gegebenen Datensatz zur Programmieraufgabe** und **nicht die trainierten Modelle** in der Abgabe hoch.

Nutzen Sie folgende Namensgebung:

```
dmml_bonus_group<groupid>.zip
├ dmml_bonus_2021.pdf
├ 01_ptu_classification.py
├ 02_ptu_regression.py
├ 03_monkey_traditional_ml.py
├ 04_monkey_simple_cnn.ipynb
├ 05_monkey_transfer_cnn.ipynb
├ utils.py
└ requirements.txt
```

Geben Sie die Jupyter Notebooks **inklusive der Ausgabe aller Zellen** ab und geben Sie in `requirements.txt` die Liste der benötigten Bibliotheken an. In `utils.py` können Sie eigene Funktionen definieren, die Sie in mehreren Skripten nutzen.

Verspätete Abgaben

Verspätete Abgaben werden akzeptiert, aber für jeden Tag, an dem die Abgabefrist überschritten wird, werden 25 % der insgesamt erreichbaren Punkte abgezogen. Am fünften Tag nach Abgabefrist, können Sie die Aufgabe nicht mehr einreichen.

Bewertungsfaktoren

Die Bewertung dieser Übung hängt von den folgenden Faktoren ab:

- Richtigkeit der Antwort
- Klarheit der Präsentation der Ergebnisse
- Schreibstil

Wenn Sie bei einer Aufgabe nicht weiterkommen, versuchen Sie zu erklären warum und beschreiben Sie die Probleme, auf die Sie gestoßen sind, da Sie dafür Teilpunkte erhalten können.

Umgang mit Plagiaten

Sie dürfen gerne kursbezogene Themen in der Vorlesung oder in unseren Moodle Foren diskutieren. Sie sollten allerdings keine Lösungen mit anderen Gruppen teilen, und alles, was Sie einreichen, muss Ihre eigene Arbeit sein. Es ist Ihnen auch nicht gestattet, Material aus dem Internet zu kopieren. Sie sind verpflichtet, jede Informationsquelle, die Sie zur Lösung der Übungsaufgabe verwendet haben (d.h. andere Materialien als die Vorlesungsmaterialien), anzuerkennen. Zitierungen haben keinen Einfluss auf Ihre Note. Mindestens 60 % der abgegeben Lösung für individuelle Aufgaben muss selbst formuliert sein. Nicht anerkennen einer Quelle, die Sie verwendet haben, ist dagegen ein klarer Verstoß gegen die akademische Ethik. Beachten Sie, dass die Universität sehr ernsthaft mit Plagiaten umgeht.

Anrechnung der Bonuspunkte

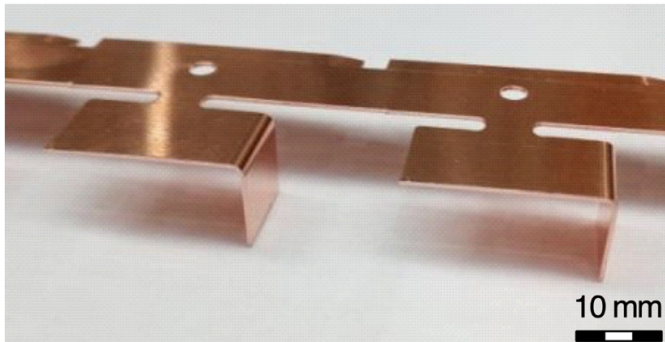
Die Klausur wird eine maximale Punktzahl von 100 Punkten haben. In dieser Übung können Sie bis zu **70 Punkte** erreichen. Die angerechneten Bonuspunkte für die Klausur ergeben sich über folgende Formel:

$$\text{Bonuspunkte für Klausur} = \frac{\text{erreichte Bonuspunkte}}{10.0}$$

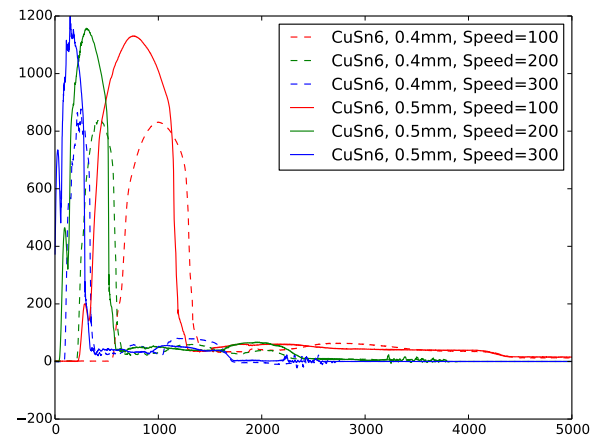
Das Intervall von fünf Punkten in der Klausur wird etwa einer Notenstufe entsprechen.

Aufgabe 1.1: </> Entscheidungsbäume - Klassifikation (12)

Im Institut für Produktionstechnik und Umformmaschinen, kurz PtU¹, der TU Darmstadt gibt es die Aufgabe eines Scherschneideverfahrens. Dabei wird mit Hilfe eines Stempels ein Loch in einen metallischen Werkstoff gestanzt, siehe Abbildung 1. Es gibt zwei Werkstoffe im Versuch: CuSn6 und 16MnCr5. Die Dicke des Materials beträgt entweder 0.4 mm oder 0.5 mm. Die Geschwindigkeit des Stempels liegt in den drei Stufen 100, 200 und 300 vor.



(a) Produziertes Werkstück.



(b) Beispielverlauf des Kraftsignals.

Abbildung 1: Veranschaulichung des PtU Prozesses.

Es gibt mehrere Sensoren, die den Status des Stanzen messen:

- Kraft Sensor / Force sensor (Fz)
- Verschiebungssensor / Displacement sensor (w)
- Beschleunigungssensor / Acceleration sensor (acc)

Im Experiment messen die Sensoren den Status des Stanzen von Beginn bis Ende jedes Prozesses. Die Zeitreihendaten von jedem Sensor werden als 1D-Vektor dargestellt. Das Interessante an dieser Aufgabe ist, dass der Status des Stanzen von der Art und Dicke des Materials abhängt. Andererseits ist es uns möglich, aus den Zeitreiheninformationen von den Sensoren auf die Art und Dicke des Materials und die Geschwindigkeit des Stanzen zu schließen.

- **Filename_Fz_raw.csv**: Enthält Daten aus dem Kraftsensor. Jeder Zeilenvektor repräsentiert das Kraftsignal. Die Anzahl der Zeilen in der Datei beträgt 2787, was der Anzahl der Experimente im Datensatz entspricht.
- **Filename_Speed.csv**: Enthält die Geschwindigkeit des Schlags der 2787 Proben. Mögliche Werte für Geschwindigkeit sind dabei {100, 200, 300}.
- **Filename_thickness.csv**: Enthält die Materialstärke in Millimeter der 2787 Proben.

In dieser Aufgabe verwenden wir die Daten des Kraftsensors, um die Proben nach der Geschwindigkeit des Stempels zu klassifizieren.

1.1a) </> 01_ptu_classification.py (9 Punkte)

Laden Sie die Daten des Kraftsensor aus `Filename_Fz_raw.csv` und die Geschwindigkeit des Stanzen aus `Filename_Speed.csv` und vervollständigen Sie den Code in `01_ptu_classification.py`:

¹<https://www.ptu.tu-darmstadt.de/>

- </>** Visualisieren Sie die ersten zehn Zeitreihen des Trainingsdatensatzes als Linien-Plot in der Methode `visualize_data()`. Stellen Sie dabei die unterschiedlichen Klassen farblich dar und geben Sie eine Legende an. Vergessen Sie nicht die Achsen zu beschriften und vermeiden Sie doppelte Einträge in der Legende. (2 Punkte)
- </>** Trainieren Sie einen Entscheidungsbaum zur Klassifikation in der Methode `fit_dt_classifier()` mithilfe von `sklearn` unter der Verwendung der Standardparameter. Der Entscheidungsbaum darf allerdings nur eine maximale Tiefe von 3 erreichen. (1 Punkt)
- </>** Ermitteln Sie die Testgenauigkeit in der Methode `get_test_accuracy()`. (1 Punkt)
- </>** Plotten Sie den resultierenden Entscheidungsbaum in `export_tree_plot()`. Sie können dabei die Funktion `tree.export_graphviz()` verwenden. (1 Punkt)
- </>** Führen Sie eine Grid-Suche für einen AdaBoost-Klassifier in der Methode `run_grid_search()` durch. Verwenden Sie dabei pro Konfiguration eine 5-fache Kreuzvalidierung. Der Suchraum beschränkt sich dabei auf die Hyperparameter: Maximale Baumtiefe der Basislerner (Entscheidungsbäume), Lernrate und Anzahl an Estimatoren. Sie können dabei die Funktion `sklearn.model_selection.GridSearchCV()` verwenden. Geben Sie zuletzt die Konfiguration der als best gefundenen Hyperparameter zurück. (3 Punkte)
- </>** Trainieren Sie in der Methode `fit_ada_boost()` einen AdaBoost-Klassifier auf den Trainingsdatensatz mit der als best gefundenen Hyperparameter-Konfiguration. (1 Punkt)

```

1 from sklearn import tree
2 from sklearn.model_selection import train_test_split
3 import numpy as np
4 import graphviz
5 from sklearn.metrics import accuracy_score
6 import matplotlib.pyplot as plt
7 from util import get_x_down_sampled
8 from sklearn.ensemble import AdaBoostClassifier, RandomForestClassifier
9 from sklearn.metrics import confusion_matrix
10 from sklearn.metrics import make_scorer
11 from sklearn.model_selection import GridSearchCV
12 import pydotplus
13 from collections import OrderedDict
14
15 def fit_dt_classifier(X_train: np.ndarray, y_train: np.ndarray, max_depth: int = None) -> tree.
    DecisionTreeClassifier:
16     """Creates and fits a decision tree classifier on the training data."""
17
18     # max_depth:The maximum depth of the tree
19     # Use sklearn.tree.DecisionTreeClassifier()
20     classifier=tree.DecisionTreeClassifier(max_depth=max_depth)
21
22     # fit():Build a decision tree classifier from the training set(X_train, y_train)
23     classifier=classifier.fit(X_train,y_train)
24
25     return classifier
26
27
28 def get_test_accuracy(clf, X_test: np.ndarray, y_test: np.ndarray) -> float:
29     """Evaluates the test accuracy for a given classifier and test dataset."""
30
31     acc=clf.score(X_test,y_test)
32     return acc
33
34
35 def export_tree_plot(clf, filename: str) -> None:
36     """Exports the tree plot for the given classifier as a pdf with given filename."""
37
38     # Export the decision tree classifier process as a graphviz.dot file.
39     dot_data=tree.export_graphviz(clf)
40
41     # A Dot class will be returned to represent the graph
42     graph=pydotplus.graph_from_dot_data(dot_data)
43
44     # Generate PDF file
45     graph.write_pdf(filename)

```

```

46
47
48 def visualize_data(X_train: np.ndarray, y_train: np.ndarray, nb_samples_to_plot: int) -> None:
49     """Visualizes the first nb_samples_to_plot in a plot with a legend of the class and exports the plot
50     as a pdf."""
51
52     # nb_samples_to_plot=10
53     X_train_plot=X_train[:nb_samples_to_plot]
54     y_train_plot=y_train[:nb_samples_to_plot]
55
56     fig,ax=plt.subplots()
57     t=np.linspace(0,5000,X_train.shape[1])
58     # Xaxis:t, Yaxis:=X_train[:10], true_label:y_train[:10]
59     for Y,true_label in zip(X_train_plot,y_train_plot):
60         if true_label==100:
61             ax.plot(t,Y,c='r',label='Speed=100')
62         elif true_label==200:
63             ax.plot(t,Y,c='g',label='Speed=200')
64         else:
65             ax.plot(t,Y,c='b',label='Speed=300')
66
67     # Return handles and labels for legend
68     handles,labels=ax.get_legend_handles_labels()
69
70     # Avoid duplicate entries in the legend
71     by_label=OrderedDict(zip(labels,handles))
72     ax.legend(by_label.values(),by_label.keys(),loc='best')
73
74     ax.set_title('first 10 time series of the training data')
75     plt.xlabel("T",fontsize=20)
76     plt.ylabel("Force",fontsize=20)
77     plt.savefig('./speed_data.pdf')
78     plt.show()
79
80 def run_grid_search(X_train, y_train, max_depth_range, learning_rate_range, n_estimators_range) -> (int,
81 float, int):
82     """Runs a grid search on the training data using the specified hyperparameter ranges using
83     a 5 fold cross-validation per configuration. At last, the best hyperparamter tuple is returned."""
84
85     df_clf=tree.DecisionTreeClassifier()
86
87     parameters={'base_estimator__max_depth':max_depth_range,'learning_rate':learning_rate_range,'
88     n_estimators':n_estimators_range}
89     Ada_clf=AdaBoostClassifier(df_clf)
90
91     # Exhaustive search over specified parameter values for an estimator.
92     clf=GridSearchCV(Ada_clf,parameters,n_jobs=-1)
93     clf.fit(X_train,y_train)
94
95     return clf.best_params_['base_estimator__max_depth'], clf.best_params_['learning_rate'], clf.
96     best_params_['n_estimators']
97
98 def fit_ada_boost(X_train: np.ndarray, y_train: np.ndarray, max_depth: int, learning_rate: float,
99 n_estimators: int) -> AdaBoostClassifier:
100     """Creates and fits an ada boost classifier on the training data."""
101     # Use the best hyperparameters(max_depth,learning_rate,n_estimators) obtained through grid search
102     df_clf=tree.DecisionTreeClassifier(max_depth=max_depth)
103
104     Ada_clf=AdaBoostClassifier(base_estimator=df_clf,learning_rate=learning_rate,n_estimators=n_estimators
105     )
106
107     Ada_clf.fit(X_train,y_train)
108     return Ada_clf
109
110 def main():
111     # for reproducibility
112     np.random.seed(42)

```

```

110
111 # load data
112 X_data = get_x_down_sampled('./PtU/FileName_Fz_raw.csv')
113 y_data = np.loadtxt(open('./PtU/FileName_Speed.csv', 'r'), delimiter=",", skiprows=0)
114
115 print("X_data.shape:", X_data.shape)
116 print("y_data.shape:", y_data.shape)
117 print("X_sample.shape:", X_data.shape)
118
119 # split training and test sets
120 X_train, X_test, y_train, y_test = train_test_split(X_data, y_data, test_size=0.2, random_state=42)
121
122 # visualize the task
123 visualize_data(X_train, y_train, 10)
124
125 # train
126 clf = fit_dt_classifier(X_train, y_train, max_depth=3)
127 # predict
128 acc = get_test_accuracy(clf, X_test, y_test)
129 print('Test Accuracy:', acc)
130
131 print("predict_proba:", clf.predict_proba(X_test))
132
133 # plot tree
134 export_tree_plot(clf, "classification_tree.pdf")
135
136 # run grid search
137 max_depth_range = list(range(1, 5))
138 learning_rate_range = [2 ** i for i in range(-2, 2)]
139 n_estimators_range = [2 ** i for i in range(5, 8)]
140 best_max_depth, best_lr, best_n_estimators = run_grid_search(X_train, y_train, max_depth_range,
141                                                             learning_rate_range, n_estimators_range)
142 clf = fit_ada_boost(X_train, y_train, best_max_depth, best_lr, best_n_estimators)
143 print(f'Best configuration: max_depth: {best_max_depth}, lr: {best_lr}, n_estimators: {
144     best_n_estimators}')
145
146 # predict
147 y_pred = clf.predict(X_test)
148 print("y_pred:", y_pred[:10], "...")
149 print("y_test:", y_test[:10], "...")
150
151 # show confusion matrix
152 print("Train Confusion Matrix for Ada Boost Classifier:\n", confusion_matrix(y_train, clf.predict(
153     X_train)))
154 print("Test Confusion Matrix for Ada Boost Classifier:\n", confusion_matrix(y_test, y_pred))
155
156 # evaluate
157 acc = accuracy_score(y_test, y_pred)
158 print('Ada Boost Test Accuracy:', acc)
159
160 if __name__ == '__main__':
161     main()

```

1.1b) Datenvisualisierung (1 Punkt)

Geben Sie die erstellte Visualisierung der Daten aus Unteraufgabe a) an.

Lösungsvorschlag:

s. Abb. 2.

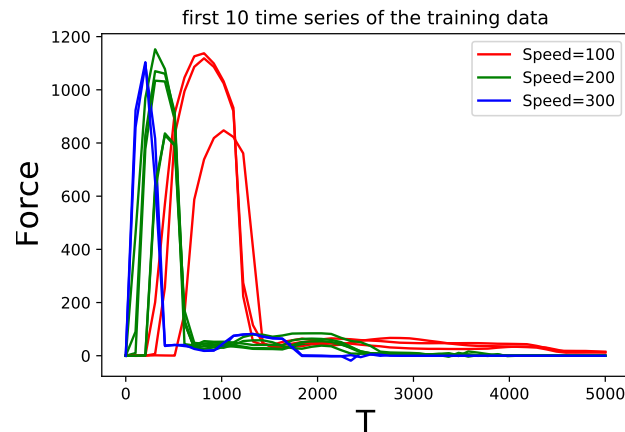


Abbildung 2: Visualisierung der Klassifikationsaufgabe.

1.1c) Visualisierung (1 Punkt)

Geben Sie die erstellte Visualisierung des Entscheidungsbaumes zur Klassifikation aus Unteraufgabe a) an.

Lösungsvorschlag:

s. Abb. 3.

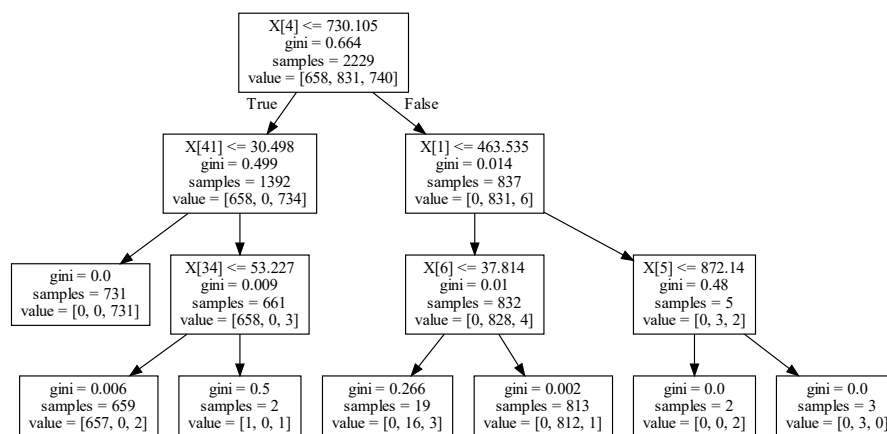


Abbildung 3: Struktur des Entscheidungsbaumes zur Klassifikation.

1.1d) Konfusionsmatrix (1 Punkt)

Geben Sie die Konfusionsmatrizen des Trainings- und Testdatensatzes mit absoluten Werten für den AdaBoost-Klassifizier an.

Lösungsvorschlag:

training	100	200	300	test	100	200	300
100	658	0	0	100	166	0	0
200	0	831	0	200	0	207	0
300	0	0	740	300	0	4	181

Aufgabe 1.2: Entscheidungsbäume - Regression (10)

Ein Entscheidungsbaum zur Regression kann auf den PtU Datensatz (s. Abb. 1) angewendet werden, um auf die Dicke des Materials zu schließen, welche ein kontinuierlicher Wert ist.

1.2a) 02_ptu_regression.py (7 Punkte)

Laden Sie die Daten des Kraftsensor aus `Filename_Fz_raw.csv` und die Materialstärken aus `Filename_thickness.csv` und vervollständigen Sie den Code in `02_ptu_regression.py`:

- `</>` Visualisieren Sie die ersten zehn Zeitreihen des Trainingsdatensatzes als Linien-Plot in der Methode `visualize_thickness_data()`. Bilden Sie dabei die Werte der Dicke farblich auf eine Farbtabelle/Colormap ab und geben Sie eine Beschriftung der Colormap an. Vergessen Sie nicht die Achsen zu beschriften. (3 Punkte)
- `</>` Visualisieren Sie die Verteilung der Werte für die Dicke mittels eines Histogramms mit 50 Bins. Färben Sie die Bins farblich mittels der gleichen Farbtabelle ein wie in der Methode `visualize_thickness_data()`. Vergessen Sie nicht die Achsen zu beschriften. (2 Punkte)
- `</>` Trainieren Sie einen Entscheidungsbaum zur Regression in `fit_dt_regressor()`. (1 Punkt)
- `</>` Berechnen Sie den mittleren quadratischen Fehler für den Testdatensatz in der Funktion `get_test_mse()`. (1 Punkt)

```

1 from sklearn import tree
2 from sklearn.model_selection import train_test_split
3 import numpy as np
4 import graphviz
5 from sklearn.metrics import mean_squared_error
6 import matplotlib.pyplot as plt
7 from util import get_x_down_sampled
8 import pydotplus
9
10 def fit_dt_regressor(X_train: np.ndarray, y_train: np.ndarray, max_depth: int = None) -> tree.
    DecisionTreeRegressor:
11     """Creates and fits a regression tree on the training data."""
12
13     # max_depth:The maximum depth of the tree
14     # Use sklearn.tree.DecisionTreeRegressor()
15     regressor = tree.DecisionTreeRegressor(max_depth=max_depth)
16
17     # fit():Build a decision tree regressor from the training set(X_train, y_train)
18     regressor.fit(X_train,y_train)
19
20     return regressor
21
22
23 def get_test_mse(clf, X_test: np.ndarray, y_test: np.ndarray) -> float:
24     """Evaluates the test mse for a given classifier and test dataset."""
25
26     y_test_pred = clf.predict(X_test)
27     # Return mean squared error regression loss
28     return mean_squared_error(y_test,y_test_pred)
29
30
31 def export_tree_plot(clf, filename: str) -> None:
32     """Exports the tree plot for the given classifier as a pdf with given filename."""
33
34     # Export the decision tree regressor process as a graphviz.dot file.
35     dot_data = tree.export_graphviz(clf)
36

```



```

37 # A Dot class will be returned to represent the graph
38 graph = pydotplus.graph_from_dot_data(dot_data)
39
40 # Generate PDF file
41 graph.write_pdf(filename)
42
43
44 def visualize_thickness_data(X_train: np.ndarray, y_train: np.ndarray, nb_samples_to_plot: int) -> None:
45     """Visualize the data by color encoding the thickness variable on a colormap and export the plot as
46     pdf."""
47
48     # print(X_train.shape)(2229,50)
49     # print(y_train.shape)(2229,)
50
51     # nb_samples_to_plot=10
52     X_train_plot = X_train[:nb_samples_to_plot]
53     y_train_plot = y_train[:nb_samples_to_plot]
54
55     plt.figure()
56     t = np.linspace(0, 5000, X_train.shape[1])
57     # Xaxis:t, Yaxis:=X_train[:10], true_label:y_train[:10]
58     for Y, true_label in zip(X_train_plot, y_train_plot):
59         plt.plot(t, Y, c=plt.cm.inferno((true_label-np.min(y_train_plot))/(np.max(y_train_plot)-np.min(
60             y_train_plot))))
61
62     sm=plt.cm.ScalarMappable(cmap="inferno", norm=plt.Normalize(np.min(y_train_plot), np.max(y_train_plot)))
63     cb=plt.colorbar(sm)
64     cb.set_label("Thickness", fontsize=20)
65     plt.xlabel("T", fontsize=20)
66     plt.ylabel("Force", fontsize=20)
67     plt.savefig('./thickness_data.pdf')
68     plt.show()
69
70 def plot_thickness_histogram(y_train: np.ndarray, bins: 50) -> None:
71     """Plots a histogram of 50 bins of the thickness distribution."""
72
73     plt.figure()
74     N,Bins,Patches=plt.hist(y_train,bins)
75     for i in range(bins):
76         Patches[i].set_facecolor(plt.cm.inferno(((Bins[i]+Bins[i+1])/2-np.min(y_train))/(np.max(y_train)-
77             np.min(y_train))))
78     plt.xlabel("Thickness", fontsize=20)
79     plt.ylabel("Frequency", fontsize=20)
80     sm=plt.cm.ScalarMappable(cmap="inferno", norm=plt.Normalize(np.min(y_train), np.max(y_train)))
81     cb=plt.colorbar(sm)
82     cb.set_label("Thickness", fontsize=20)
83     plt.savefig('./thickness_histogram.pdf')
84     plt.show()
85
86 def main():
87     # for reproducibility
88     np.random.seed(42)
89
90     # load data
91     X_data = get_x_down_sampled('./PtU/FileName_Fz_raw.csv')
92     y_data = np.loadtxt(open('./PtU/FileName_thickness.csv', 'r'), delimiter=",", skiprows=0)
93     print("X_data.shape:", X_data.shape)
94     print("y_data.shape:", y_data.shape)
95
96     # split training and test sets
97     X_train, X_test, y_train, y_test = train_test_split(X_data, y_data, test_size=0.2, random_state=42)
98
99     # visualize the task
100     visualize_thickness_data(X_train, y_train, 10)
101     plot_thickness_histogram(y_train, 50)
102
103     # train
104     clf = fit_dt_regressor(X_train, y_train)
105     # predict % evaluate

```

```

104 mse = get_test_mse(clf, X_test, y_test)
105 print('Test MSE:', mse)
106
107 # change max tree depth
108 # train
109 clf = fit_dt_regressor(X_train, y_train, max_depth=3)
110
111 # predict & evaluate
112 mse = get_test_mse(clf, X_test, y_test)
113 print('Test MSE:', mse)
114
115 # plot tree
116 export_tree_plot(clf, "regression_tree_d3.pdf")
117
118
119 if __name__ == '__main__':
120     main()

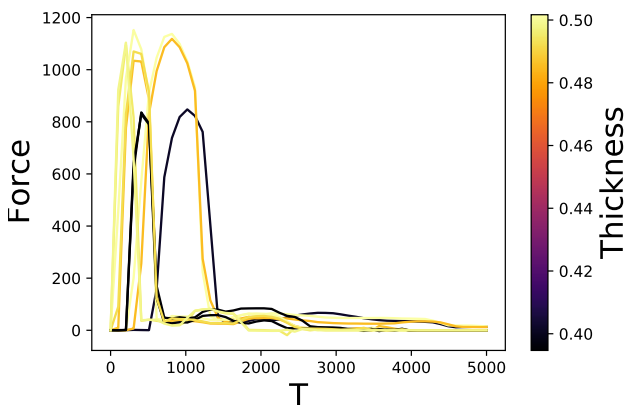
```

1.2b) Datenvisualisierung (2 Punkte)

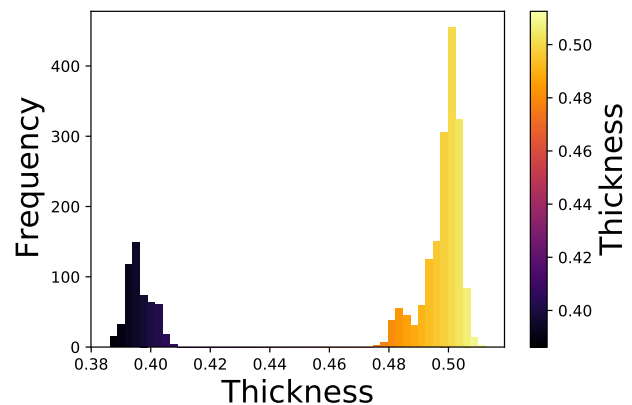
Geben Sie die erstellte Visualisierung der Daten und das Histogramm aus Unteraufgabe a) an.

Lösungsvorschlag:

s. Abb. 4.



(a) Materialdicke der ersten zehn Zeitreihen des Trainingsdatensatzes.



(b) Verteilung der Materialdicke des Trainingsdatensatzes.

Abbildung 4: Visualisierung der Regressionsaufgabe.

1.2c) Visualisierung (1 Punkt)

Geben Sie die erstellte Visualisierung des Entscheidungsbaumes zur Regression aus Unteraufgabe a).

Lösungsvorschlag:

s. Abb. 5.

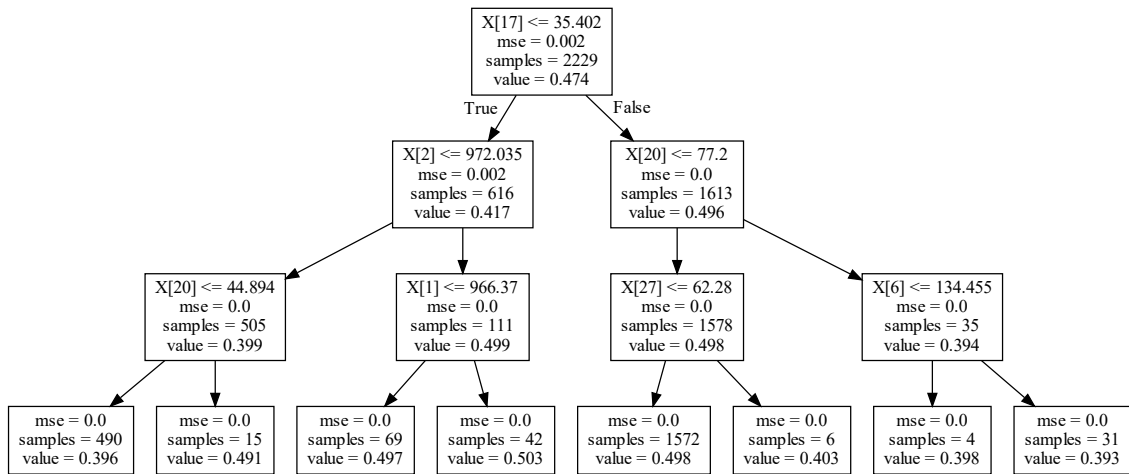


Abbildung 5: Struktur des Entscheidungsbaumes zur Regression.

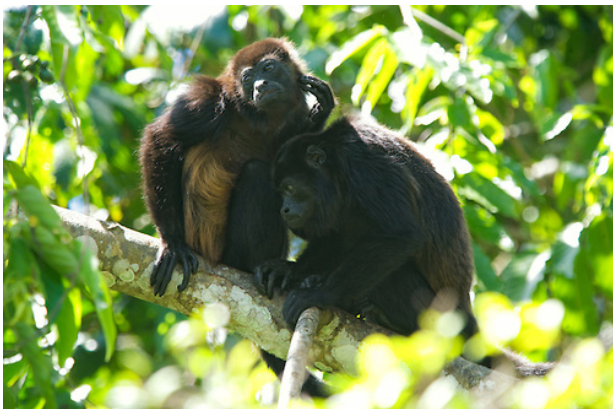
Aufgabe 1.3: Dimensionsreduktion und Klassifikation (18)

In dieser und den darauf folgenden Aufgaben werden wir den 10 Monkey Species Datensatz (s. Abb. 6) verwenden.

Laden Sie den Datensatz auf [kaggle.com](https://www.kaggle.com/slothkong/10-monkey-species) herunter.

<https://www.kaggle.com/slothkong/10-monkey-species>

Das Ziel des Datensatzes ist es die Bilder von zehn verschiedenen Affenarten korrekt mittels Verfahren des maschinellen Lernens richtig zu klassifizieren.



(a) Klasse n0: Mantelbrüllaffe (*alouatta palliata*)



(b) Klasse n7: Gewöhnlicher Totenkopffaffe (*saimiri sciureus*)

Abbildung 6: Beispielphoto der Klasse *saimiri sciureus* aus dem Datensatz 10 Monkey Species

1.3a) 03_monkey_traditional_ml.py (16 Punkte)

Vervollständigen Sie den Code in 03_monkey_traditional_ml.py:

Implementieren Sie die Methode `get_train_val_data()`, welche den Trainings- und Validierungsdatsatz aus dem Daten-Verzeichnis einliest und auf eine vorgegebene Bildgröße von 224×224 skaliert. Zum Laden und

Skalieren der Bilder können Sie die OpenCV Bibliothek verwenden. Wenn Sie nicht den vollständigen Datensatz in den Arbeitsspeicher laden können, können Sie eine kleinere Auflösung verwenden. Die Methode soll den Trainings- und Validierungsdatensatz als numpy-Array mit den drei Farbkanälen RGB zurückgeben. Die zehn Klassen sollen als Integer-Werte $\in [0, 9]$ beschrieben werden. (3 Punkte)

- `</>` Führen Sie in der Methode `apply_pca()` eine Dimensionalitätsreduktion auf zwei Dimensionen durch und geben Sie den transformierten Trainings- und Validierungsdatensatz zurück. Beachten Sie, dass Sie die Eigenwerte und Eigenvektoren nur anhand des Trainings- und nicht anhand des Validierungsdatensatzes bestimmen. Die Methode `apply_pca()` soll Tensoren mit 4 Dimensionen (z.B. Bilddaten), als auch Tensoren mit 2 Dimensionen (z.B. Matrizen), verarbeiten können. (3 Punkte)
- `</>` Visualisieren Sie den transformierten Trainingsdatensatz als 2D-Scatterplot in der Methode `plot_pca_data()`. Stellen Sie die unterschiedlichen Klassen farblich dar und beschriften Sie die Achsen. (1 Punkt)
- `</>` Evaluieren Sie in `evaluate_knn()` mittels des dimensionsreduzierten Datensatzes die Performance eines k-Nächste-Nachbarn Klassifiers für $k=1$. (2 Punkte)
- `</>` Extrahieren Sie ein Farbhistogramm mit 32 Bins pro Farbe für den Datensatz in der Methode `get_histogram_data()`. Pro Zeile sollen die Bins der Farbkanäle nacheinander angehängt werden. Also 96 Einträge pro Zeile, wobei die ersten 32 für Farbe Rot, die Einträge 33-64 für Farbe Grün und 65-96 für Farbe Blau reserviert sind. (3 Punkte)
- `</>` Trainieren Sie einen Klassifier Ihrer Wahl in der Methode `train_custom_classifier()` mittels des gegebenen Datensatzes. Sie können dabei einen beliebigen Modelltyp der `sklearn` Bibliothek unter folgenden Einschränkung verwenden: Die Trainingszeit sollte unter 5 Minuten auf einer herkömmlichen CPU betragen und Sie dürfen kein neuronales Netz / Multi-Layer-Perceptron verwenden. Versuchen Sie eine Validierungsgenauigkeit von $\geq 35\%$ zu erreichen. (4 Punkte)

Hinweis: Sie können Merkmale kombinieren und müssen nicht zwangsläufig alle Merkmale verwenden.

Lösungsvorschlag:

```

1 from sklearn.decomposition import PCA
2 from sklearn.metrics import accuracy_score
3 from sklearn.ensemble import RandomForestClassifier
4 from sklearn import preprocessing
5
6
7 def get_train_val_data(data_dir: str, shape: tuple) -> (np.ndarray, np.ndarray, np.ndarray, np.ndarray):
8     """Reads the training and validation image data from the given data directory and
9     scales all images to the given to the input shape.
10    The target vectors (y_train, y_val) should contain the class labels as integer values."""
11
12    import glob
13    train_folder = os.path.join(data_dir, "training/training")
14    x = []
15    y = []
16
17    for i in range(10):
18        # for filename in glob.glob(train_folder + "/n{}/*.jpg".format(i)):
19        for filename in glob.glob(train_folder + "/n{}/*".format(i)):
20            image = cv2.imread(filename)
21            image = cv2.resize(image, shape) # resize the image
22            image_RGB = cv2.cvtColor(image, cv2.COLOR_BGR2RGB) # change to RGB
23            x.append(image_RGB)
24            y.append(i)
25    X_train = np.array(x)
26    y_train = np.array(y)
27
28    # print("Training Dataset:")
29    # print(X_train.shape)
30    # print(y_train.shape)
31
32    test_folder = os.path.join(data_dir, "validation/validation")

```

```

33 x = []
34 y = []
35
36 for i in range(10):
37     for filename in glob.glob(test_folder + "/n{}/*".format(i)):
38         image = cv2.imread(filename)
39         image = cv2.resize(image, shape) # resize the image
40         image_RGB = cv2.cvtColor(image, cv2.COLOR_BGR2RGB) # change to RGB
41         x.append(image_RGB)
42         y.append(i)
43 X_val = np.array(x)
44 y_val = np.array(y)
45
46
47 # print("Validation Dataset:")
48 # print(X_val.shape)
49 # print(y_val.shape)
50 return X_train, y_train, X_val, y_val
51
52 def plot_pca_data(X_data_pca: np.ndarray, y_data: np.ndarray, filename: str):
53     """Plots the first two principal components in a 2d scatter plot and saves the figure."""
54     x1 = X_data_pca[:, 0]
55     x2 = X_data_pca[:, 1]
56     y = y_data
57
58     plt.scatter(x1, x2, c=y)
59     plt.title(filename)
60     plt.xlabel('the first principal component in X_data_pca')
61     plt.ylabel('the second principal component in X_data_pca')
62     plt.show()
63
64
65
66 def get_histogram_data(X_data, nb_bins, color_channels=3) -> np.ndarray:
67     """Extracts a histogram for each channel and returns the features """
68     x=[]
69
70     for i in range(X_data.shape[0]):
71         img = X_data[i, :, :, :]
72         hist_r = cv2.calcHist([img], [0], None, [nb_bins], [0, 255])
73         hist_g = cv2.calcHist([img], [1], None, [nb_bins], [0, 255])
74         hist_b = cv2.calcHist([img], [2], None, [nb_bins], [0, 255])
75         hist_end = (np.concatenate([hist_r, hist_g, hist_b])).T
76         x.append(hist_end)
77     X_hist_basis = np.array(x)
78     X_hist = X_hist_basis.squeeze()
79     return X_hist
80
81
82
83 def evaluate_knn(X_train_pca, y_train, X_val_pca, y_val, k: int) -> (float, float):
84     """Evaluates a k-nearest-Neighbour classifier with value k and returns the training and validation
85     accuracy."""
86     classifier = KNeighborsClassifier(n_neighbors=k, metric='minkowski', p=2)
87     classifier.fit(X_train_pca, y_train)
88     train_acc = classifier.score(X_train_pca, y_train)
89     val_acc = classifier.score(X_val_pca, y_val)
90
91     return train_acc, val_acc
92
93
94 def apply_pca(X_train: np.ndarray, X_val: np.ndarray, n_components: int) -> (np.ndarray, np.ndarray):
95     """Returns the transformed data for X_train and X_val with n_components as principal components."""
96     # if len(X_train.shape)>2:
97     #     X_train = (X_train.reshape(X_train.shape[0], X_train.shape[1] * X_train.shape[2] * X_train.shape
98     # [3]))
99     #
100     # if len(X_val.shape)>2:
101     #     X_val = (X_val.reshape(X_val.shape[0], X_val.shape[1] * X_val.shape[2] * X_val.shape[3]))

```

```

101 X_train_reshape = X_train.reshape(X_train.shape[0],-1)
102 X_train = X_train_reshape - np.mean(X_train_reshape, axis= 0).reshape(1,-1)
103 X_val_reshape = X_val.reshape(X_val.shape[0],-1)
104 X_val = X_val_reshape - np.mean(X_val_reshape,axis = 0).reshape(1,-1)
105
106
107
108
109 pca = PCA(n_components=n_components)
110 pca.fit(X_train)
111 X_train_pca = pca.transform(X_train)
112 X_val_pca = pca.transform(X_val)
113 return X_train_pca, X_val_pca
114
115
116 def train_custom_classifier(X_train, X_train_pca, X_train_pca_hist,
117                             X_val, X_val_pca, X_val_pca_hist,
118                             y_train, y_val) -> (float, float):
119     """Trains a custom non neural network classifier and returns the training and validation accuracy."""
120
121
122     X_train_forest = np.hstack((X_train_pca, X_train_pca_hist))
123     rfc = RandomForestClassifier(random_state=42)
124     rfc.fit(X_train_forest, y_train)
125     train_acc = rfc.score(X_train_forest, y_train)
126     X_val_forest = np.hstack((X_val_pca, X_val_pca_hist))
127     val_acc = rfc.score(X_val_forest, y_val)
128     return train_acc, val_acc
129
130 def main():
131     # for reproducibility
132     np.random.seed(42)
133
134     data_dir = './data/kaggle/10-monkey-species'
135     print('Reading image data...')
136     get_train_val_data(data_dir, (224, 224))
137     X_train, y_train, X_val, y_val = get_train_val_data(data_dir, (224, 224))
138
139     print('Apply PCA on original data...')
140     X_train_pca, X_val_pca = apply_pca(X_train, X_val, 2)
141     plot_pca_data(X_train_pca, y_train, 'pca_all_features.pdf')
142     train_acc, val_acc = evaluate_knn(X_train_pca, y_train, X_val_pca, y_val, 1)
143     print(f"knn (k=1) Validation-Acc: {val_acc}")
144
145     X_train_hist = get_histogram_data(X_train, 32)
146
147     X_val_hist = get_histogram_data(X_val, 32)
148     print('Apply PCA on color histogram data...')
149     X_train_pca_hist, X_val_pca_hist = apply_pca(X_train_hist, X_val_hist, 2)
150     plot_pca_data(X_train_pca_hist, y_train, 'pca_color_histogram.pdf')
151     train_acc, val_acc = evaluate_knn(X_train_pca_hist, y_train, X_val_pca_hist, y_val, 1)
152     print(f"knn (k=1) Validation-Acc: {val_acc}")
153
154     train_acc, val_acc = train_custom_classifier(X_train, X_train_pca, X_train_pca_hist,
155                                                  X_val, X_val_pca, X_val_pca_hist,
156                                                  y_train, y_val)
157
158     print(f"Validation Acc. for custom classifier: {val_acc}")
159
160
161 if __name__ == '__main__':
162     main()

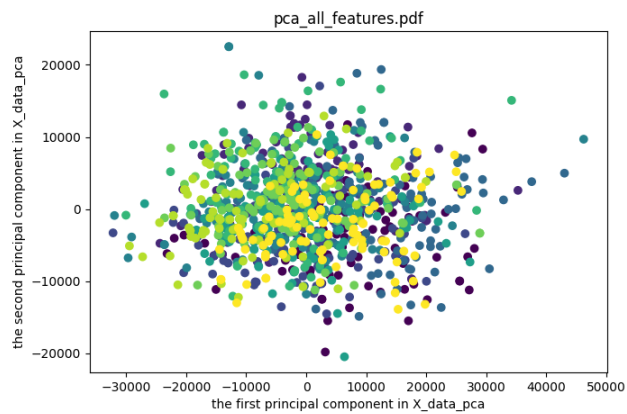
```

1.3b) Datenvisualisierung (2 Punkte)

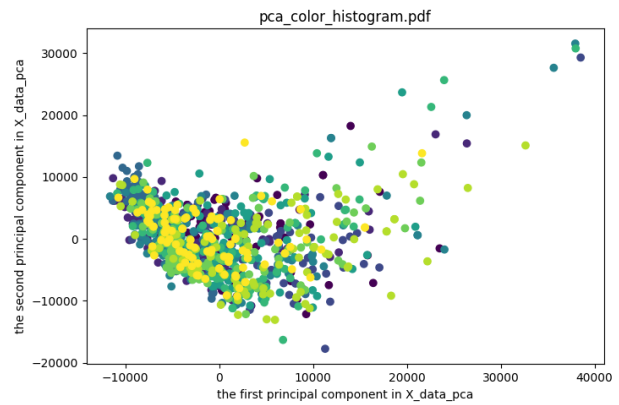
Geben Sie die beiden Scatterplots der Datenverteilung nach Transformation mittels PCA aus Unteraufgabe a) an.

Lösungsvorschlag:

s. Abb. 7.



(a) Datenverteilung nach Anwendung von PCA auf Rohdaten.



(b) Datenverteilung nach Anwendung von PCA auf Farbhistogramme.

Abbildung 7: Datenverteilung nach Transformation mittels PCA.

Aufgabe 1.4: Einfaches Convolutional Neural Network (12)

In dieser Aufgabe verwenden wir den 10 Monkey Species Datensatz (s. Abb. 6), um ein einfaches Convolutional Neural Networks (CNN) zu trainieren.

Verwenden Sie ein Deep Learning Framework Ihrer Wahl und eine passende Trainingskonfiguration. Wenn Sie mit noch keinen Deep Learning Framework vertraut sind, können wir Keras als Einstiegsframework empfehlen.

Wenn Sie Ihre Kenntnisse in Deep Learning weiter vertiefen möchten können Sie den Kurs *Deep Learning Methods & Architectures* (20-00-1034-iv) im nächsten Sommersemester besuchen. Ergänzende Literatur zum Thema Deep Learning und Transfer Learning finden Sie in *Modern mathematics of Deep Learning* [1] und *The Principles of Deep Learning Theory* [2].

1.4a) 04_monkey_simple_cnn.ipynb (10 Punkte)

Implementieren Sie Ihre Lösung im Jupyter-Notebook `04_monkey_simple_cnn.ipynb` und geben Sie das Jupyter-Notebook samt Trainingslog ab.

Erstellen und trainieren Sie ein einfaches Convolutional Neural Network mit bis zu vier Convolutional Layern. Sie können dabei Max-Pooling Layer oder Strided Convolution verwenden.

Beim Trainingsvorgang können Sie ebenfalls Datenaugmentation verwenden.

Versuchen Sie eine Validierungsgenauigkeit von $\geq 50\%$ zu erreichen. Die Trainingszeit sollte unter 15 Minuten auf einer NVIDIA GTX 1080ti vergleichbaren GPU betragen. (10 Punkte)

Hinweis: Sollte keiner ihrer Gruppenmitglieder eine GPU zur Verfügung haben, können Sie stattdessen mittels CPU trainieren oder Google Colab nutzen:

<https://colab.research.google.com/notebooks/intro.ipynb>

```

1 import pandas as pd
2 import tensorflow as tf
3 from tensorflow.keras import datasets, layers, models, initializers, preprocessing
4 import matplotlib.pyplot as plt
5 from PIL import Image
6 # todo
7 ### Load the image data or define the data loader
8
9 train_path = '/Users/llm/Desktop/dmml_bonus_2021_scripts_2/archive/training/training'
10 valid_path = '/Users/llm/Desktop/dmml_bonus_2021_scripts_2/archive/validation/validation'
11 label_path = '/Users/llm/Desktop/dmml_bonus_2021_scripts_2/archive/monkey_labels.txt'
12
13 # data augmentation
14 '''
15 reference : https://www.tensorflow.org/api_docs/python/tf/keras/preprocessing/image/
16             ImageDataGenerator
17 '''
18 train_datagen = preprocessing.image.ImageDataGenerator(
19     rescale = 1./255,
20     rotation_range = 20,
21     width_shift_range = 0.2,
22     height_shift_range = 0.2,
23     shear_range = 0.2,
24     zoom_range = 0.2,
25     horizontal_flip = True,
26     fill_mode = "nearest"
27 )
28 train_generator = train_datagen.flow_from_directory(
29     train_path,
30     target_size = (128, 128),
31     batch_size = 64,
32     seed = 7,
33     shuffle = True,
34     class_mode = 'categorical'
35 )
36 valid_datagen = preprocessing.image.ImageDataGenerator(

```



```

37     rescale = 1./255
38 )
39 valid_generator = valid_datagen.flow_from_directory(
40     valid_path,
41     target_size = (128, 128),
42     batch_size = 64,
43     seed = 7,
44     shuffle = False,
45     class_mode = 'categorical')
46 train_num = train_generator.samples
47 valid_num = valid_generator.samples
48
49
50 # ### Define your model architecture
51 '''
52 Reference : https://www.tensorflow.org/tutorials/images/cnn?hl=zh-eg
53 '''
54
55 model = models.Sequential()
56 model.add(layers.Conv2D(32, kernel_size=3, activation='relu', input_shape=[128, 128, 3]))
57 model.add(layers.MaxPooling2D(pool_size=2))
58 model.add(layers.Conv2D(64, kernel_size=3, activation='relu'))
59 model.add(layers.MaxPooling2D(pool_size=2))
60 model.add(layers.Conv2D(64, kernel_size=3, activation='relu'))
61 model.add(layers.MaxPooling2D(pool_size=2))
62 model.add(layers.Conv2D(64, kernel_size=3, activation='relu'))
63 model.add(layers.MaxPooling2D(pool_size=2))
64 model.add(layers.Flatten())
65 model.add(layers.Dense(64, activation='relu'))
66 model.add(layers.Dense(10, activation='softmax'))
67 model.summary()
68
69 # ### Define your training setup
70
71 # todo
72 model.compile(optimizer='adam',
73               loss="categorical_crossentropy",
74               metrics=['accuracy'])
75
76
77 # ### Initialize your model weights
78 # todo
79 '''
80 Reference: https://keras.io/api/layers/initializers/
81 '''
82
83 initializers.TruncatedNormal(mean=0.0, stddev=0.05, seed=None)
84
85 # ### Run the training process
86
87 '''
88 Reference: https://www.tensorflow.org/tutorials/images/cnn?hl=zh-eg
89 '''
90
91 history = model.fit_generator(train_generator,
92                               steps_per_epoch=train_num // 64,
93                               epochs=20,
94                               validation_data=valid_generator,
95                               validation_steps=valid_num // 64)
96
97 # ### Plot the learning curves
98
99 '''
100 Reference: https://www.tensorflow.org/tutorials/images/cnn?hl=zh-eg
101 '''
102
103
104
105
106

```

```

107 plt.plot(history.history['accuracy'], label='accuracy')
108 plt.plot(history.history['val_accuracy'], label = 'val_accuracy')
109 plt.xlabel('Epoch')
110 plt.ylabel('Accuracy')
111 plt.ylim([0.0, 1])
112 plt.legend()
113 plt.show()
114 plt.plot(history.history['loss'], label='loss')
115 plt.plot(history.history['val_loss'], label = 'val_loss')
116 plt.xlabel('Epoch')
117 plt.ylabel('Accuracy')
118 plt.ylim([1.0, 2.5])
119 plt.legend()
120 plt.show()

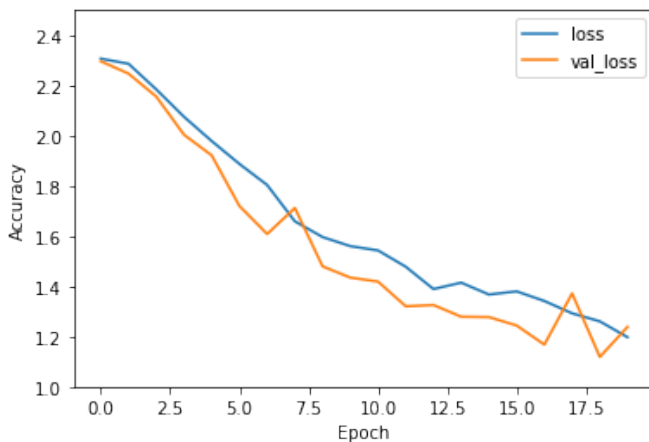
```

1.4b) Lernplots (2 Punkte)

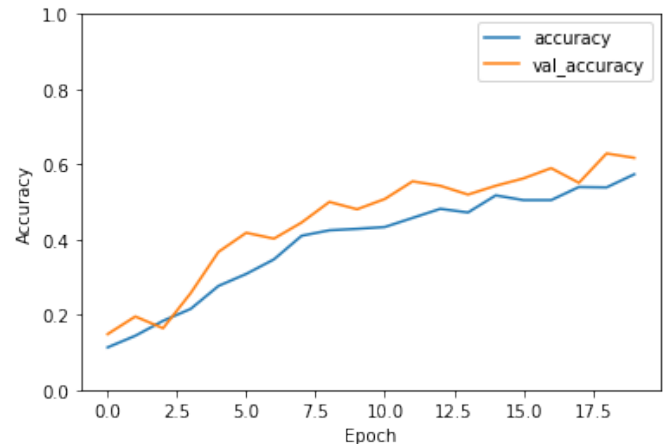
Geben Sie hier in zwei Plots den Verlauf des Trainings- und Validierungsfehlers sowie der Trainings- und Validierungsgenauigkeit an.

Lösungsvorschlag:

s. Abb. 8.



(a) Entwicklung des Trainings- und Validierungsfehlers.



(b) Entwicklung der Trainings- und Validierungsgenauigkeit.

Abbildung 8: Lernplots für ein einfaches Convolutional Neural Network.

Aufgabe 1.5: Convolutional Neural Network - Transfer Learning (18)

In dieser Aufgabe verwenden wir den 10 Monkey Species Datensatz (s. Abb. 6), um Transfer Learning anzuwenden.

1.5a) 05_monkey_transfer_cnn.ipynb (15 Punkte)

Implementieren Sie Ihre Lösung im Jupyter-Notebook 05_monkey_transfer_cnn.ipynb und geben Sie das Jupyter-Notebook samt Trainingslog ab.

Verwenden Sie ein vortrainiertes Deep-Learning-Modell aus dem Modell-Zoo des gewählten Deep-Learning Frameworks und transformieren Sie damit die Bilddaten nach der Ausgabe der letzten Convolutional Layer. (4 Punkte)

Hinweis: Als Empfehlung können Sie **MobileNetV1** verwenden. Dieses ist auf

<https://keras.io/api/applications/>

verfügbar und kann mittels der Keras API heruntergeladen werden.

Wenden Sie auf diesen durch das CNN transformierten Daten erneut eine Dimensionsreduktion nach PCA auf zwei Dimensionen an. (2 Punkte)

Visualisieren Sie analog zu Aufgabe 3a) diesen 2d-Datensatz in einem Scatterplot. (1 Punkt)

Evaluieren Sie analog zu Aufgabe 3a) diesen 2d-Datensatz mittels eines k-Nächste-Nachbarn Klassifiers, um die Affenarten zu klassifizieren. (2 Punkte)

Verwenden Sie das gleiche vortrainierte Modell aus dem Modell-Zoo des gewählten Deep-Learning Frameworks und fügen Sie ein passenden Klassifizierungskopf an. Optimieren Sie die Parameter der neu angefügten Schichten des neuronalen Netzwerkes. Die Parameter der vorangehenden Schichten bleiben fixiert. Beim Trainingsvorgang können Sie ebenfalls Datenaugmentation verwenden.

Geben Sie die erreichte Trainings- und Validierungsgenauigkeit an. Die Trainingszeit sollte unter 15 Minuten auf einer NVIDIA GTX 1080ti vergleichbaren GPU betragen. Finden Sie eine Konfiguration um eine Validierungsgenauigkeit $\geq 95\%$ zu erreichen. (6 Punkte)

```

1 # In[137]:
2
3
4 import tensorflow as tf
5 import pandas as pd
6 import numpy as np
7 from matplotlib import pyplot as plt
8 from keras.preprocessing.image import ImageDataGenerator
9 from sklearn.decomposition import PCA
10 from matplotlib.colors import ListedColormap
11 from sklearn.neighbors import KNeighborsClassifier
12 from keras.callbacks import ModelCheckpoint
13
14
15 # Load the image data or define the data loader
16
17 # In[138]:
18
19
20 cols = ['Label', 'Latin Name', 'Common Name', 'Train Image', 'Validation Images']
21 monkeys = pd.read_csv(f"monkey_labels.txt", names=cols, skiprows=1)
22 height = 224
23 width = 224
24 #make the data generator.
25 train_datagen = ImageDataGenerator(
26     rescale=1./255,
27     horizontal_flip=True,
28     fill_mode='nearest')
29 train_generator = train_datagen.flow_from_directory(
30     'training/training/',
31     target_size=(height, width),

```

```

32     batch_size=1,
33     shuffle=False)
34 test_datagen = ImageDataGenerator(rescale=1./255)
35 test_generator = train_datagen.flow_from_directory(
36     'validation/validation/',
37     target_size=(height, width),
38     batch_size=1,
39     shuffle=False)
40
41
42 # In[139]:
43
44
45 monkeys
46
47
48 # ### Load a pre-trained model
49
50 # In[140]:
51
52
53 mobileNet = tf.keras.applications.MobileNet(input_shape=(224, 224, 3), include_top=False, weights='
    imagenet')
54 mobileNet.trainable = False
55
56
57 # ### Transform the image data to the features of the last convolutional layer
58
59 # In[142]:
60
61
62 features = np.zeros([1098, 7, 7, 1024])
63 labels = np.zeros(1098)
64 labels_oh = np.zeros([1098, 10])
65 for i in range(1098):
66     feature = mobileNet.predict(train_generator[i][0])
67     features[i] = feature
68     labels_oh[i] = train_generator[i][1].reshape(-1)
69     labels[i] = train_generator[i][1].reshape(-1).argmax()
70     if i % 100 == 0:
71         print("finished {} procent!".format(i // 10))
72
73
74 # ### Apply PCA to the transformed CNN features
75
76 # In[143]:
77
78
79 features = features.reshape([1098, 7*7*1024])
80 pca_model = PCA(n_components=2)
81 features_d = pca_model.fit_transform(features)
82
83
84 # ### Plot the two dimensional data distribution
85
86 # In[145]:
87
88
89 cmap = ListedColormap(['b', 'r', 'g', 'y', 'c', 'm', 'k', 'olive', 'pink', 'purple', 'orange'])
90 plt.figure()
91 plt.scatter(features_d[:, 0], features_d[:, 1], c=labels, marker='.', cmap=cmap, edgecolors='none',
    label="train data")
92 plt.xlabel("$x_1$")
93 plt.ylabel("$x_2$")
94 plt.savefig("two_dim_data")
95 plt.show()
96
97
98 # ### Evaluate a kNN with k=1 on the two dimensional data distribution
99

```

```

100 # In[146]:
101
102
103 cls = KNeighborsClassifier(n_neighbors=1)
104 cls.fit(features_d, labels)
105 features_test = np.zeros([272, 7, 7, 1024])
106 y_label = np.zeros([272])
107 label_test_oh = np.zeros([272, 10])
108 for i in range(272):
109     features_test[i] = mobileNet(test_generator[i][0])
110     y_label[i] = test_generator[i][1].reshape(-1).argmax()
111     label_test_oh[i] = test_generator[i][1].reshape(-1)
112 features_test_d = pca_model.transform(features_test.reshape([272, 7*7*1024]))
113 y_pred = cls.predict(features_test_d)
114
115
116 # In[147]:
117
118
119 acc = (y_pred == y_label).sum() / 272
120 print("The accuracy of validation sets is: {}".format(acc))
121
122
123 # ### Add a custom classification head on the pre-trained model
124
125 # In[148]:
126
127
128 fc_model = tf.keras.Sequential([
129     tf.keras.layers.Dense(5000, activation='relu'),
130     tf.keras.layers.Dense(100, activation='relu'),
131     tf.keras.layers.Dense(10, activation='softmax')
132 ])
133 fc_model.build(input_shape=[None, 7*7*1024])
134 fc_model.summary()
135
136
137 # ### Define your training setup
138
139 # In[149]:
140
141
142 epochs = 50
143 batch_size = 64
144 checkpoint = ModelCheckpoint("monkey_model.h5f", monitor='val_acc', verbose=1, save_best_only=True,
145                             mode='max')
146 ls=tf.keras.losses.CategoricalCrossentropy()
147 features_test = features_test.reshape([272, 7*7*1024])
148 fc_model.compile(
149     optimizer = 'sgd',
150     loss = tf.keras.losses.CategoricalCrossentropy(from_logits = False),
151     metrics = ['acc'])
152
153 # ### Run the training process
154
155 # In[150]:
156
157
158 history = fc_model.fit(
159     features, labels_oh,
160     batch_size = batch_size,
161     validation_data = (features_test, label_test_oh),
162     steps_per_epoch = 5,
163     epochs = epochs,
164     callbacks=[checkpoint],
165     verbose=1)
166
167
168 # ### Plot the learning curves

```

```

169
170 # In[152]:
171
172
173 plt.figure()
174 steps = np.arange(1, epochs)
175 plt.plot(steps, history.history["loss"][1:], label="training sets")
176 plt.plot(steps, history.history["val_loss"][:1], label="validation sets")
177 plt.xlabel("step")
178 plt.ylabel("loss")
179 plt.title("loss of training sets and test sets")
180 plt.legend()
181 plt.savefig("loss_pic.png")
182 plt.show()
183 plt.figure()
184 plt.plot(steps, history.history["acc"][1:], label="training sets")
185 plt.plot(steps, history.history["val_acc"][:1], label="validation sets")
186 plt.xlabel("step")
187 plt.ylabel("accuracy")
188 plt.title("accuracy of training sets and test sets")
189 plt.legend()
190 plt.savefig("acc_pic.png")
191 plt.show()
192
193
194 # Here we use this model to predict two monkeys. From the training step, we can see that the accuracy
    of the most accurate model is 0.9596. This model is good to estimate the class of monkey.
195
196 # In[153]:
197
198
199 def predict_monkey(mobileNet, fc_model, testdata):
200     fc_model.trainable = False
201     real_name = monkeys["Common Name"][int(tf.argmax(testdata[1][0]))]
202     features = tf.reshape(mobileNet.predict(testdata[0]), [-1, 50176])
203     result = int(tf.argmax(tf.reshape(fc_model.predict(features), [10])))
204     pred_name = monkeys["Common Name"][result]
205     plt.figure()
206     plt.imshow(testdata[0][0])
207     plt.title("Predicted label is: " + real_name + "\nreal label is: " + pred_name)
208     predict_monkey(mobileNet, fc_model, test_generator[50])
209     predict_monkey(mobileNet, fc_model, test_generator[100])
    
```

1.5b) Datenvisualisierung: PCA (1 Punkt)

Geben Sie hier den Scatterplot der Datenverteilung nach Transformation mittels PCA aus Unteraufgabe a) an.

Lösungsvorschlag:

s. Abb. 9.

1.5c) Lernplots (2 Punkte)

Geben Sie hier in zwei Plots den Verlauf des Trainings- und Validierungsfehlers sowie der Trainings- und Validierungsgenauigkeit an.

Lösungsvorschlag:

s. Abb. 10.

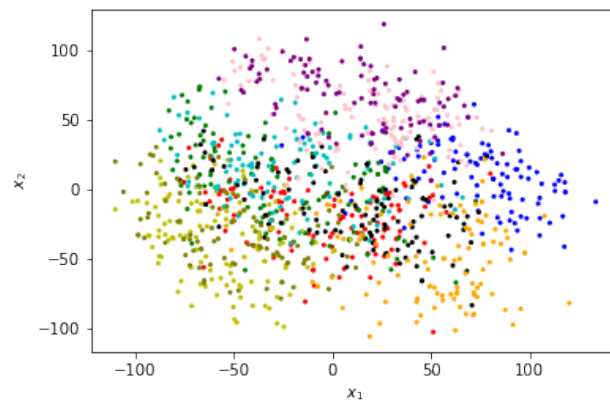
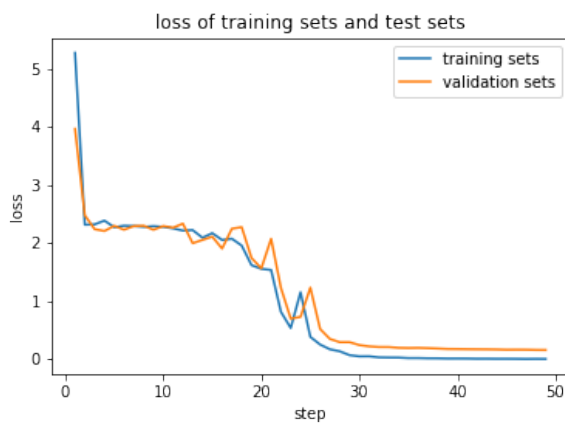
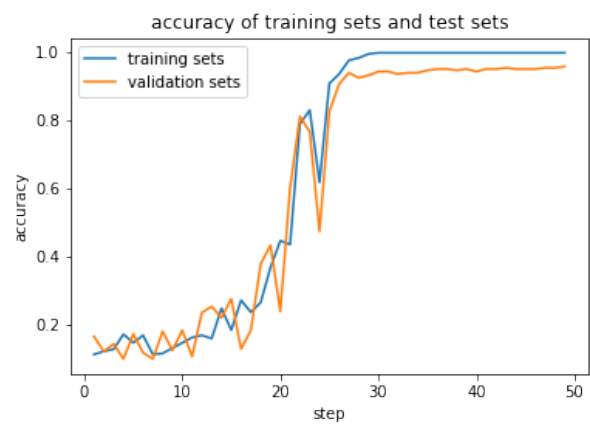


Abbildung 9: Datenverteilung der Merkmale des vortrainierten CNNs nach Transformation mittels PCA.



(a) Entwicklung des Trainings- und Validierungsfehlers



(b) Entwicklung der Trainings- und Validierungsgenauigkeit

Abbildung 10: Übersicht der Lernplots des Trainingsvorgangs

Literatur

- [1] Julius Berner, Philipp Grohs, Gitta Kutyniok, and Philipp Petersen. The modern mathematics of deep learning. *arXiv preprint arXiv:2105.04026*, 2021.
- [2] Daniel A Roberts, Sho Yaida, and Boris Hanin. The Principles of Deep Learning Theory. <https://deeplearningtheory.com/PDLT.pdf>, page 449, 2021.