ETH zürich



Übungslektion 12 – Machine Learning

Informatik II

6. / 7. Mittwoch 2025

Heutiges Programm

Wiederholung der Vorlesung

Theoretische Übungen

Praktische Übungen

Homework

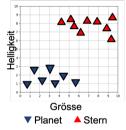
1. Wiederholung der Vorlesung

Machinelles Lernen

Trainieren von Modellen mithilfe von Beispielen. ML hat Anwendungen in:

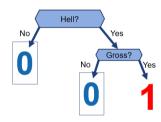
- Astronomie
- Physik
- Medizin
- Etc.

Sterne und Planeten Datensatz



Sterne sind gross und hell. Planeten sind kleiner und weniger hell.

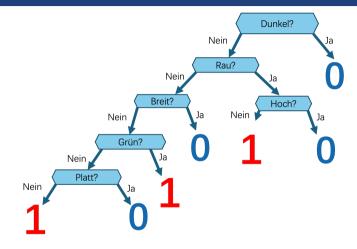
Entscheidungsbaum

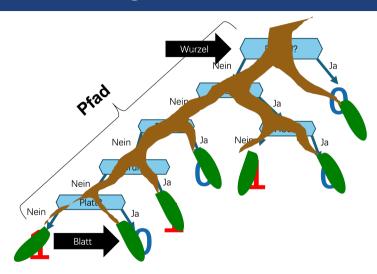


Planeten: 0, Sterne: 1

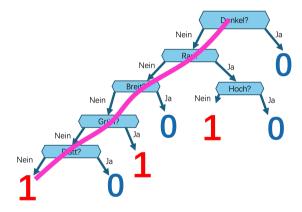
Tiefe des Baumes: Anzahl von Knoten im längsten Pfad (zwischen Wurzel und Blatt).

4





Tiefe = 5



■ **Datensatz** Erstellen Sie einen Datensatz und teilen Sie ihn in den Trainingssatz *D* und den Validierungssatz *D'* auf.

- **Datensatz** Erstellen Sie einen Datensatz und teilen Sie ihn in den Trainingssatz *D* und den Validierungssatz *D'* auf.
- Modell auswählen Modell \mathcal{H} : eine Menge von Entscheidungsfunktionen (z.B. Entscheidungsbäume der Tiefe \leq 3, logistische Regression, oder spezifische NN Architektur)

- **Datensatz** Erstellen Sie einen Datensatz und teilen Sie ihn in den Trainingssatz *D* und den Validierungssatz *D'* auf.
- Modell auswählen Modell \mathcal{H} : eine Menge von Entscheidungsfunktionen (z.B. Entscheidungsbäume der Tiefe \leq 3, logistische Regression, oder spezifische NN Architektur)
- **Verlustfunktion (loss)** Funktion L(D, f), wobei $f \in \mathcal{H}$, die misst, wie gut f in D abschneidet.

- **Datensatz** Erstellen Sie einen Datensatz und teilen Sie ihn in den Trainingssatz *D* und den Validierungssatz *D'* auf.
- Modell auswählen Modell \mathcal{H} : eine Menge von Entscheidungsfunktionen (z.B. Entscheidungsbäume der Tiefe \leq 3, logistische Regression, oder spezifische NN Architektur)
- **Verlustfunktion (loss)** Funktion L(D, f), wobei $f \in \mathcal{H}$, die misst, wie gut f in D abschneidet.
- **Training** Finden einer Entscheidungsfunktion $f^* \in \mathcal{H}$, für die $L(D, f^*)$ klein ist.

- **Datensatz** Erstellen Sie einen Datensatz und teilen Sie ihn in den Trainingssatz *D* und den Validierungssatz *D'* auf.
- Modell auswählen Modell \mathcal{H} : eine Menge von Entscheidungsfunktionen (z.B. Entscheidungsbäume der Tiefe \leq 3, logistische Regression, oder spezifische NN Architektur)
- **Verlustfunktion (loss)** Funktion L(D, f), wobei $f \in \mathcal{H}$, die misst, wie gut f in D abschneidet.
- **Training** Finden einer Entscheidungsfunktion $f^* \in \mathcal{H}$, für die $L(D, f^*)$ klein ist.
- Validierung Auswerten der Entscheidungsfunktion f^* an einem Dataset $D' \neq D$.

Initialisierung des Pseudo-Zufallszahlen-Generators

Zufälligkeit spielt beim maschinellen Lernen eine entscheidende Rolle.

■ Wir haben Zufälligkeit bei der Datenerfassung, bei der Datenaufteilung (in Trainings-/Testsätze), bei den Trainingsmodellen usw.

Als Quelle der Zufälligkeit wird normalerweise ein **Pseudozufallszahlengenerator** verwendet.

- Es handelt sich um eine mathematische Funktion, die eine Folge nahezu zufälliger Zahlen generiert.
- Die Sequenz ist deterministisch und wird mit einer (oder mehreren) Anfangszahl(en) initialisiert, die seed genannt wird.

In Code Expert legen wir den Ausgangspunkt fest und machen alles deterministisch. Sie erhalten reproduzierbare Ergebnisse.

Datensatz erstellen

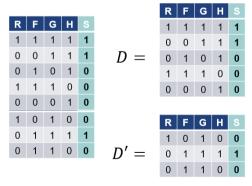
Features oder Merkmale

Class label

	Rot	Flackern	Gross	Hell	Stern
D =	1	1	1	1	1
	0	0	1	1	1
	0	1	0	1	0
	1	1	1	0	0
	0	0	0	1	0

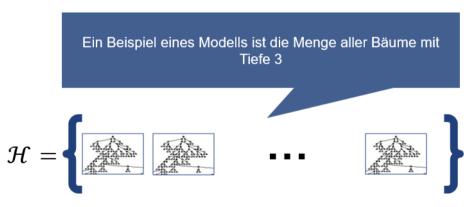
Aufteilen des Datensatzes

lacktriangle Den Datensatz in einen Trainingsdatensatz D und einen Validierungsdatensatz D' aufteilen.



Modell auswählen

Modell \mathcal{H} : eine Menge von Entscheidungsfunktionen



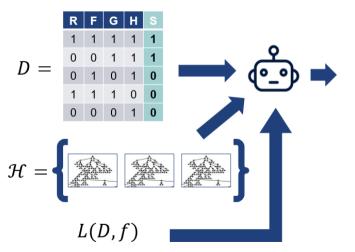
Verlustfunktion

- Funktion L(D, f), wobei $f \in \mathcal{H}$, die misst, wie gut f in D abschneidet.
- Normalerweise benutzt man die 0/1 Verlustfunktion für das Trainieren von Bäumen.

$$L(D,f) = \frac{\textit{\#Beispiele in } \textit{D}, \textit{die von } \textit{f} \textit{ falsch klassifiziert wurden}}{\textit{\#Beispiele in } \textit{D}}$$

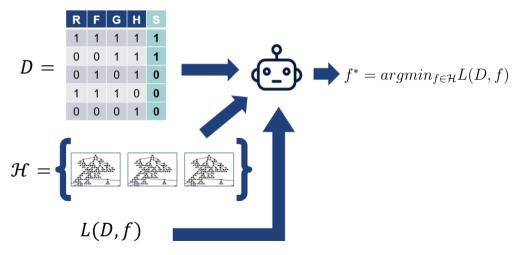
Trainieren

Finden Sie einen Schätzer $f^* \in \mathcal{H}$, für die der Wert von $L(D, f^*)$ niedrig ist.



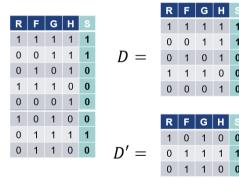
Trainieren

Finden Sie einen Schätzer $f^* \in \mathcal{H}$, für die der Wert von $L(D, f^*)$ niedrig ist.



Validierung

■ Der Schätzer f^* auf dem Validierungsdatensatz $D' \neq D$ auswerten.



Verwechslungsmatrix

- Eine Verwirrungsmatrix ist eine Tabelle mit der Verteilung der Klassifikatorleistung auf die Daten.
- Es handelt sich um eine $N \times N$ -Matrix, die zur Bewertung der Leistung eines Klassifizierungsmodells verwendet wird.

		Predicted		
		Ungiftig	Giftig	
True Label	Ungiftig	True Positive (TP)	False Positive (FP)	
	Giftig	False Negative (FN)	True Negative (TN)	

Balanced Accuracy

- Die "Balanced Accuracy" ist die durchschnittliche Genauigkeit aller Klassen.
- Sie kann als Durchschnitt der Diagonalen der normalisierten Verwirrungsmatrix berechnet werden.
- Dies ist hilfreich beim Umgang mit unausgeglichenen Daten.

		Predicted		
		Ungiftig	Giftig	
True Label	Ungiftig	20	70	
True Labet	Giftig	30	5000	

Ausgewogene Genauigkeit

		Predicted		
		Ungiftig	Giftig	
True Label	Ungiftig	20	70	
True Labet	Giftig	30	5000	

- Berechnen wir die Genauigkeit dieser Vorhersage: $\frac{TP+TN}{TP+FN+FP+TN} \approx 98,05\%$.
- Dieses Ergebnis ist beeindruckend, allerdings wird die Ungiftig-Spalte in der Vorhersage selbst nicht richtig behandelt.
- Verwenden Sie die "Balanced Accuracy": $\frac{1}{2}(\frac{TP}{TP+FN}+\frac{TN}{TN+FP})\approx 60,80\%$
- Dadurch ist die Punktzahl niedriger als von der Genauigkeit vorhergesagt, da beide Klassen das gleiche Gewicht erhalten.

Regression: Finden einer Beziehung (Funktion) zwischen einer abhängigen und unabhängigen Variable.

Lineare Regression: Finden einer Linearen Abbildung, die den Vektor $\mathbf x$ von unabhängigen Variablen auf eine abhängige Variable y abbildet.

Regression: Finden einer Beziehung (Funktion) zwischen einer abhängigen und unabhängigen Variable.

Lineare Regression: Finden einer Linearen Abbildung, die den Vektor $\mathbf x$ von unabhängigen Variablen auf eine abhängige Variable y abbildet.

■ **Datensatz**: Jeder Datenpunkt besteht aus einem Vektor $\mathbf x$ und einem Skalar y.

Regression: Finden einer Beziehung (Funktion) zwischen einer abhängigen und unabhängigen Variable.

Lineare Regression: Finden einer Linearen Abbildung, die den Vektor \mathbf{x} von unabhängigen Variablen auf eine abhängige Variable y abbildet.

- **Datensatz**: Jeder Datenpunkt besteht aus einem Vektor $\mathbf x$ und einem Skalar y.
- **Modell**: Die Menge der linearen Abbildungen $\mathbb{R}^n \to \mathbb{R}$. Wir finden den Vektor $\mathbf{w} : \mathbf{w} \cdot \mathbf{x} = y$.

Regression: Finden einer Beziehung (Funktion) zwischen einer abhängigen und unabhängigen Variable.

Lineare Regression: Finden einer Linearen Abbildung, die den Vektor $\mathbf x$ von unabhängigen Variablen auf eine abhängige Variable y abbildet.

- **Datensatz**: Jeder Datenpunkt besteht aus einem Vektor $\mathbf x$ und einem Skalar y.
- **Modell**: Die Menge der linearen Abbildungen $\mathbb{R}^n \to \mathbb{R}$. Wir finden den Vektor $\mathbf{w} : \mathbf{w} \cdot \mathbf{x} = y$.
- **Verlustfunktion**: Die Summe der Quadrate: $\sum_{\mathbf{x} \in D} (\mathbf{w} \cdot \mathbf{x} y)^2$

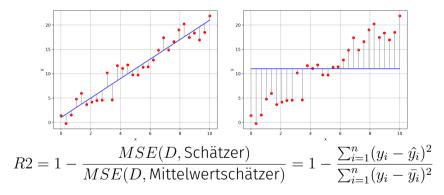
Regression: Finden einer Beziehung (Funktion) zwischen einer abhängigen und unabhängigen Variable.

Lineare Regression: Finden einer Linearen Abbildung, die den Vektor $\mathbf x$ von unabhängigen Variablen auf eine abhängige Variable y abbildet.

- **Datensatz**: Jeder Datenpunkt besteht aus einem Vektor $\mathbf x$ und einem Skalar y.
- **Modell**: Die Menge der linearen Abbildungen $\mathbb{R}^n \to \mathbb{R}$. Wir finden den Vektor $\mathbf{w} : \mathbf{w} \cdot \mathbf{x} = y$.
- **Verlustfunktion**: Die Summe der Quadrate: $\sum_{\mathbf{x} \in D} (\mathbf{w} \cdot \mathbf{x} y)^2$
- **Trainieren**: Mit der Methode der kleinsten Quadrate (nicht klausurrelevant).
- **Validierung**: Messen der Summe der Quadrate auf dem Validierungsdatensatz *D'*.

R2-Score

■ Wie gut ist ein Schätzer verglichen mit dem Mittelwertschätzer?



- lacksquare 0 < R2 < 1: Besser als Mittelwertschätzer
- $-\infty < R2 < 0$: Schlechter als Mittelwertschätzer

Gini Index

Wie unrein sind die Partitionen eines Entscheidungsbaums?

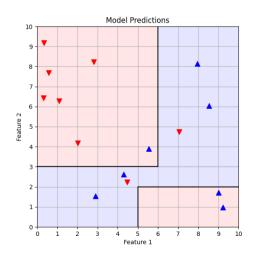
■ Einzelne Partition:

$$G(D, \mathsf{Part}_m) = \sum_k \hat{p}_{m,k} \cdot (1 - \hat{p}_{m,k})$$

 $\hat{p}_{m,k}$: Anteil der Klasse k in Part $_m$.

■ Insgesamt:

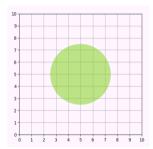
$$G(D, f) = \sum_{m} \frac{|\mathsf{Part}_{m}|}{|\mathsf{Part}|} \cdot G(D, \mathsf{Part}_{m})$$



2. Theoretische Übungen

Machen Sie diese Übung nur mit Bleistift und Papier. Gegeben sei der Kreis auf dem Bild.

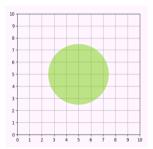
Führen Sie die folgenden Schritte aus:



Machen Sie diese Übung nur mit Bleistift und Papier. Gegeben sei der Kreis auf dem Bild.

Führen Sie die folgenden Schritte aus:

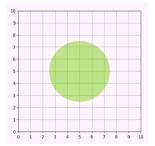
1. Erstellen Sie einen Trainingsdatensatz *D* mit ca. 10 Datenpunkten. (Geben Sie die x- und y-Koordinaten und die Klasse an.)



Machen Sie diese Übung nur mit Bleistift und Papier. Gegeben sei der Kreis auf dem Bild.

Führen Sie die folgenden Schritte aus:

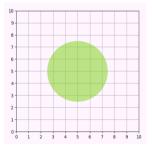
- 1. Erstellen Sie einen Trainingsdatensatz *D* mit ca. 10 Datenpunkten. (Geben Sie die x- und y-Koordinaten und die Klasse an.)
- 2. Wählen Sie ein Modell H. aus.



Machen Sie diese Übung nur mit Bleistift und Papier. Gegeben sei der Kreis auf dem Bild.

Führen Sie die folgenden Schritte aus:

- 1. Erstellen Sie einen Trainingsdatensatz *D* mit ca. 10 Datenpunkten. (Geben Sie die x- und y-Koordinaten und die Klasse an.)
- 2. Wählen Sie ein Modell H. aus.
- 3. Wählen Sie eine Verlustfunktion L(D, f) aus.

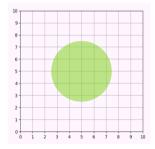


Übung 1: Training Reproduzieren

Machen Sie diese Übung nur mit Bleistift und Papier. Gegeben sei der Kreis auf dem Bild.

Führen Sie die folgenden Schritte aus:

- 1. Erstellen Sie einen Trainingsdatensatz *D* mit ca. 10 Datenpunkten. (Geben Sie die x- und y-Koordinaten und die Klasse an.)
- 2. Wählen Sie ein Modell ${\cal H}$ aus.
- 3. Wählen Sie eine Verlustfunktion L(D, f) aus.
- 4. Erstellen Sie einen Validierungsdatensatz D' mit ca. 5 Datenpunkten.



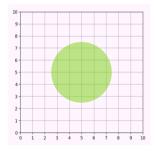
Das Innere sollte als 1 klassifiziert werden, das Äussere als 0.

Übung 1: Training Reproduzieren

Machen Sie diese Übung nur mit Bleistift und Papier. Gegeben sei der Kreis auf dem Bild.

Führen Sie die folgenden Schritte aus:

- 1. Erstellen Sie einen Trainingsdatensatz *D* mit ca. 10 Datenpunkten. (Geben Sie die x- und y-Koordinaten und die Klasse an.)
- 2. Wählen Sie ein Modell ${\cal H}$ aus.
- 3. Wählen Sie eine Verlustfunktion $\mathcal{L}(D,f)$ aus.
- 4. Erstellen Sie einen Validierungsdatensatz D' mit ca. 5 Datenpunkten.
- 5. Finden Sie einen Schätzer f^* , der L(D, f) minimiert.



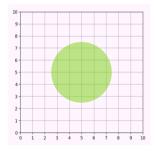
Das Innere sollte als 1 klassifiziert werden, das Äussere als 0.

Übung 1: Training Reproduzieren

Machen Sie diese Übung nur mit Bleistift und Papier. Gegeben sei der Kreis auf dem Bild.

Führen Sie die folgenden Schritte aus:

- 1. Erstellen Sie einen Trainingsdatensatz *D* mit ca. 10 Datenpunkten. (Geben Sie die x- und y-Koordinaten und die Klasse an.)
- 2. Wählen Sie ein Modell \mathcal{H} aus.
- 3. Wählen Sie eine Verlustfunktion $\mathcal{L}(D,f)$ aus.
- 4. Erstellen Sie einen Validierungsdatensatz D' mit ca. 5 Datenpunkten.
- 5. Finden Sie einen Schätzer f^* , der L(D, f) minimiert.
- 6. Validieren Sie den Schätzer f^* an D'.



Das Innere sollte als 1 klassifiziert werden, das Äussere als 0.

Beispiel Antwort 1

Beispiel Antwort 1

Übung 2: Training

Machen Sie diese Übung nur mit Bleistift und Papier. Gegeben sei der Trainingsdatensatz D für ein dreieckiges Muster im Bild.

	x-Koord	y-Koord	Klasse
	1	1	1
	3	5	1
	5	9	1
•	6	7	1
\mathcal{D} =	8	3	1
	1	3	0
	2	6	0
	4	9	0
	7	8	0
	9	4	0

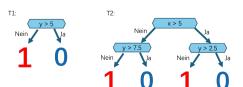
Finden Sie einen Baum mit der Tiefe 3, der alle Punkte im obigen Datensatz richtig klassifiziert.

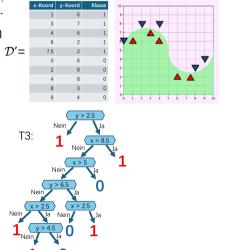
Beispiel Antwort 2

Übung 3: Validierung

Machen Sie diese Übung nur mit Bleistift und Papier. Gegeben seien drei Entscheidungsbäume T1, T2, T3 mit Tiefe 1, 2, und 6, und ein Validierungsdatensatz D', der dem Muster im Bild entspricht.

- Berechnen Sie für jeden Baum, wie viele Punkte in D' durch den Baum inkorrekt klassifiziert werden.
- Entscheiden Sie, welcher Baum der Beste ist.





Antwort 3

3. Praktische Übungen

Übersicht

Code Expert:

https://expert.ethz.ch/enrolled/SS25/mavt2/codeExamples Exercise 12 - In-class

Zwei Programmierübungen:

- 1. **Iris**: Klassifikation mit Entscheidungsbäumen
- 2. **House prices**: Regression mit der Linearen Regression

Iris Datensatz: klassischer Datensatz aus dem Jahr 1936. Er enthält Messungen von Blütenblattern (petal) und Kelchblättern (sepal) von 150 Blüten, die zu einer von 3 Arten von Iris (*Iris setosa, Iris versicolor, Iris virginica*) gehören.

Iris Datensatz: klassischer Datensatz aus dem Jahr 1936. Er enthält Messungen von Blütenblattern (petal) und Kelchblättern (sepal) von 150 Blüten, die zu einer von 3 Arten von Iris (*Iris setosa, Iris versicolor, Iris virginica*) gehören.

Machen Sie die folgenden Übungen in Code Expert:

1. Lesen Sie den Datensatz data.csv mit pandas ein.

Iris Datensatz: klassischer Datensatz aus dem Jahr 1936. Er enthält Messungen von Blütenblattern (petal) und Kelchblättern (sepal) von 150 Blüten, die zu einer von 3 Arten von Iris (*Iris setosa, Iris versicolor, Iris virginica*) gehören.

- 1. Lesen Sie den Datensatz data.csv mit pandas ein.
- 2. Teilen Sie den Datensatz in einen Trainingsdatensatz und einen Testdatensatz auf.

Iris Datensatz: klassischer Datensatz aus dem Jahr 1936. Er enthält Messungen von Blütenblattern (petal) und Kelchblättern (sepal) von 150 Blüten, die zu einer von 3 Arten von Iris (*Iris setosa, Iris versicolor, Iris virginica*) gehören.

- 1. Lesen Sie den Datensatz data.csv mit pandas ein.
- 2. Teilen Sie den Datensatz in einen Trainingsdatensatz und einen Testdatensatz auf.
- 3. Trainieren Sie einen Entscheidungsbaum für die Klassifikation der Arten von Iris aus den Messungen.

Iris Datensatz: klassischer Datensatz aus dem Jahr 1936. Er enthält Messungen von Blütenblattern (petal) und Kelchblättern (sepal) von 150 Blüten, die zu einer von 3 Arten von Iris (*Iris setosa, Iris versicolor, Iris virginica*) gehören.

- 1. Lesen Sie den Datensatz data.csv mit pandas ein.
- 2. Teilen Sie den Datensatz in einen Trainingsdatensatz und einen Testdatensatz auf.
- 3. Trainieren Sie einen Entscheidungsbaum für die Klassifikation der Arten von Iris aus den Messungen.
- 4. Verwenden Sie den Baum, um die Blumen im Testdatensatz zu klassifizieren.

Iris Datensatz: klassischer Datensatz aus dem Jahr 1936. Er enthält Messungen von Blütenblattern (petal) und Kelchblättern (sepal) von 150 Blüten, die zu einer von 3 Arten von Iris (*Iris setosa, Iris versicolor, Iris virginica*) gehören.

- 1. Lesen Sie den Datensatz data.csv mit pandas ein.
- 2. Teilen Sie den Datensatz in einen Trainingsdatensatz und einen Testdatensatz auf.
- 3. Trainieren Sie einen Entscheidungsbaum für die Klassifikation der Arten von Iris aus den Messungen.
- 4. Verwenden Sie den Baum, um die Blumen im Testdatensatz zu klassifizieren.
- 5. Lesen Sie X_final.csv mit pandas ein, verwenden Sie den Baum um die Blumen zu klassifizieren, und geben Sie die Vorhersagen zurück.

Iris Datensatz: klassischer Datensatz aus dem Jahr 1936. Er enthält Messungen von Blütenblattern (petal) und Kelchblättern (sepal) von 150 Blüten, die zu einer von 3 Arten von Iris (*Iris setosa, Iris versicolor, Iris virginica*) gehören.

- 1. Lesen Sie den Datensatz data.csv mit pandas ein.
- 2. Teilen Sie den Datensatz in einen Trainingsdatensatz und einen Testdatensatz auf.
- 3. Trainieren Sie einen Entscheidungsbaum für die Klassifikation der Arten von Iris aus den Messungen.
- 4. Verwenden Sie den Baum, um die Blumen im Testdatensatz zu klassifizieren.
- 5. Lesen Sie X_final.csv mit pandas ein, verwenden Sie den Baum um die Blumen zu klassifizieren, und geben Sie die Vorhersagen zurück.
- 6. Die Vorhersagen werden von Code Expert automatisch bewertet.

Iris Datensatz: klassischer Datensatz aus dem Jahr 1936. Er enthält Messungen von Blütenblattern (petal) und Kelchblättern (sepal) von 150 Blüten, die zu einer von 3 Arten von Iris (*Iris setosa, Iris versicolor, Iris virginica*) gehören.

Machen Sie die folgenden Übungen in Code Expert:

- 1. Lesen Sie den Datensatz data.csv mit pandas ein.
- 2. Teilen Sie den Datensatz in einen Trainingsdatensatz und einen Testdatensatz auf.
- 3. Trainieren Sie einen Entscheidungsbaum für die Klassifikation der Arten von Iris aus den Messungen.
- 4. Verwenden Sie den Baum, um die Blumen im Testdatensatz zu klassifizieren.
- 5. Lesen Sie X_final.csv mit pandas ein, verwenden Sie den Baum um die Blumen zu klassifizieren, und geben Sie die Vorhersagen zurück.
- 6. Die Vorhersagen werden von Code Expert automatisch bewertet.

32

In dieser Übung trainieren wir eine Funktion mithilfe eines fiktiven Datensatzes, die den Preis eines Hauses anhand von seiner Fläche schätzen kann.

In dieser Übung trainieren wir eine Funktion mithilfe eines fiktiven Datensatzes, die den Preis eines Hauses anhand von seiner Fläche schätzen kann.

Machen Sie die folgenden Übungen in Code Expert:

1. Lesen Sie den Datensatz data.csv mit pandas ein.

In dieser Übung trainieren wir eine Funktion mithilfe eines fiktiven Datensatzes, die den Preis eines Hauses anhand von seiner Fläche schätzen kann.

- 1. Lesen Sie den Datensatz data.csv mit pandas ein.
- 2. Teilen Sie den Datensatz in einen Trainingsdatensatz und einen Testdatensatz auf.

In dieser Übung trainieren wir eine Funktion mithilfe eines fiktiven Datensatzes, die den Preis eines Hauses anhand von seiner Fläche schätzen kann.

- 1. Lesen Sie den Datensatz data.csv mit pandas ein.
- 2. Teilen Sie den Datensatz in einen Trainingsdatensatz und einen Testdatensatz auf.
- 3. Trainieren Sie ein lineares Regressionsmodel mit dem Trainingsdatensatz.

In dieser Übung trainieren wir eine Funktion mithilfe eines fiktiven Datensatzes, die den Preis eines Hauses anhand von seiner Fläche schätzen kann.

- 1. Lesen Sie den Datensatz data.csv mit pandas ein.
- 2. Teilen Sie den Datensatz in einen Trainingsdatensatz und einen Testdatensatz auf.
- 3. Trainieren Sie ein lineares Regressionsmodel mit dem Trainingsdatensatz.
- 4. Verwenden Sie das Modell, um die Hauspreise vorherzusagen. Bewerten Sie die Qualität der Vorhersage mit dem R2-Wert.

In dieser Übung trainieren wir eine Funktion mithilfe eines fiktiven Datensatzes, die den Preis eines Hauses anhand von seiner Fläche schätzen kann.

- 1. Lesen Sie den Datensatz data.csv mit pandas ein.
- 2. Teilen Sie den Datensatz in einen Trainingsdatensatz und einen Testdatensatz auf.
- 3. Trainieren Sie ein lineares Regressionsmodel mit dem Trainingsdatensatz.
- 4. Verwenden Sie das Modell, um die Hauspreise vorherzusagen. Bewerten Sie die Qualität der Vorhersage mit dem R2-Wert.
- 5. Lesen Sie X_final.csv mit pandas ein, verwenden Sie das Modell, um die Preise vorherzusagen und geben Sie die Vorhersagen zurück.

In dieser Übung trainieren wir eine Funktion mithilfe eines fiktiven Datensatzes, die den Preis eines Hauses anhand von seiner Fläche schätzen kann.

- 1. Lesen Sie den Datensatz data.csv mit pandas ein.
- 2. Teilen Sie den Datensatz in einen Trainingsdatensatz und einen Testdatensatz auf.
- 3. Trainieren Sie ein lineares Regressionsmodel mit dem Trainingsdatensatz.
- 4. Verwenden Sie das Modell, um die Hauspreise vorherzusagen. Bewerten Sie die Qualität der Vorhersage mit dem R2-Wert.
- 5. Lesen Sie X_final.csv mit pandas ein, verwenden Sie das Modell, um die Preise vorherzusagen und geben Sie die Vorhersagen zurück.
- 6. Die Vorhersagen werden von Code Expert automatisch bewertet.

In dieser Übung trainieren wir eine Funktion mithilfe eines fiktiven Datensatzes, die den Preis eines Hauses anhand von seiner Fläche schätzen kann.

Machen Sie die folgenden Übungen in Code Expert:

- 1. Lesen Sie den Datensatz data.csv mit pandas ein.
- 2. Teilen Sie den Datensatz in einen Trainingsdatensatz und einen Testdatensatz auf.
- 3. Trainieren Sie ein lineares Regressionsmodel mit dem Trainingsdatensatz.
- 4. Verwenden Sie das Modell, um die Hauspreise vorherzusagen. Bewerten Sie die Qualität der Vorhersage mit dem R2-Wert.
- 5. Lesen Sie X_final.csv mit pandas ein, verwenden Sie das Modell, um die Preise vorherzusagen und geben Sie die Vorhersagen zurück.
- 6. Die Vorhersagen werden von Code Expert automatisch bewertet.

Detaillierte Anweisungen befinden sich auf Code Expert.

4. Homework

Übungslektion 10: Intro ML I

On https://expert.ethz.ch/enrolled/SS25/mavt2/exercises

- Cancer Detection
- Diabetes Prediction
- Circles

Abgabedatum: Montag 19.05.2025, 20:00 MEZ Das Grading funktioniert wie bei den praktischen Übungen.

NO HARDCODING

Fragen oder Anregungen?