

## Metadaten Datensatz 1: Titanic - Überleben einer Katastrophe



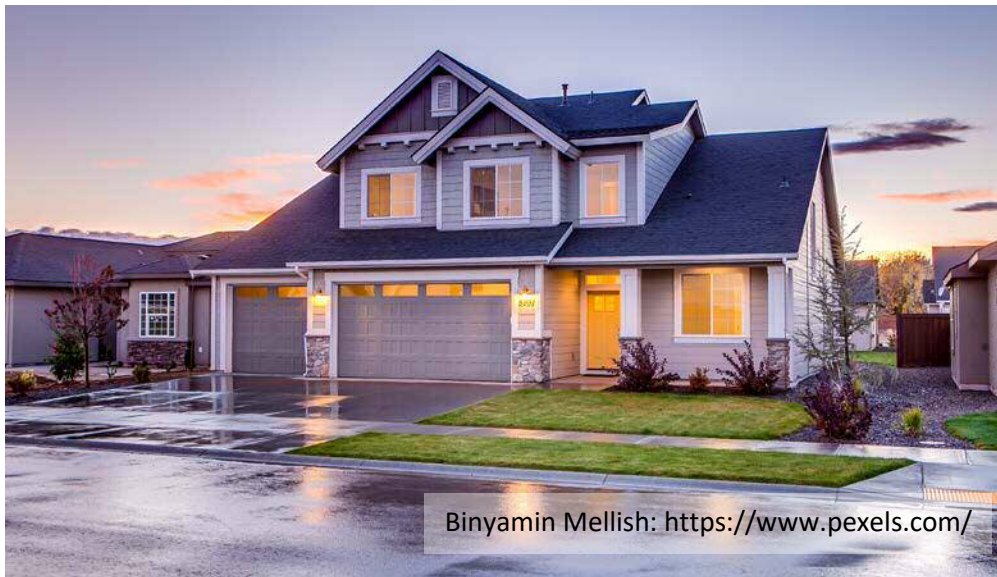
<b>Daten</b>	<b>Titanic.zip</b>
<b>Inhalte</b>	Passagierliste der Titanic (train_with_survival_info.csv, test_no_survival_info.csv)
<b>Informationen</b>	<p>Der Train-Datensatz beinhaltet die Informationen bzgl. Überlebensstatus. Mit diesem Datensatz können Sie das Modell entwickeln.</p> <p>Der Test-Datensatz soll genutzt werden, um die Performance des Modells zu testen. Hier gibt es keine Informationen zum Überlebensstatus der Passagiere und das Modell soll diesen Vorhersagen.</p>
<b>Fragestellung</b>	Welche Passagiere im Testdatensatz können die Katastrophe überleben und welche nicht? Welche Passagiermerkmale haben den größten Einfluss auf die Überlebenschancen?
<b>Vokabular</b>	<p>Survival - Überlebensstatus    0 = Nein, 1 = Ja</p> <p>pclass   Ticket Klassen    1 = Erste Klasse (Upper) , 2 = Zweite Klasse, 3 = Dritte Klasse (zur Einschätzung des sozio-ökonomischen Status der Passagiere)</p> <p>Name</p> <p>sex      Geschlecht      male /female</p> <p>Age      Alter (Jahre)</p> <p>sibsp    Anzahl der Geschwister / Ehepartner an Bord der Titanic</p> <p>parch    Anzahl der Eltern / Kinder an Bord der Titanic</p> <p>ticket   Ticket-Nummer</p> <p>fare     Ticketpreis</p> <p>cabin    Kabinen-Nummer</p> <p>embarked    Zustiegshafen    C = Cherbourg, Q = Queenstown, S = Southampton</p>
<b>Quelle</b>	<a href="https://www.kaggle.com/competitions/titanic">https://www.kaggle.com/competitions/titanic</a>

## Metadaten Datensatz 2: Diabetes – Diagnostische Vorhersage einer Erkrankung



<b>Daten</b>	<b>diabetes.csv</b>
<b>Inhalte</b>	Patientenliste vom US National Institute of Diabetes and Digestive and Kidney Diseases
<b>Informationen</b>	Dieser Datensatz enthält die medizinischen Messdaten von Patientinnen, die von nordamerikanischen Indigenen abstammen (Pima-Indianer). Diese Gruppe zeigt ein besonders erhöhtes Risiko, an Diabetes zu erkranken.
<b>Fragestellung</b>	Ziel des Datensatzes ist es, auf der Grundlage bestimmter diagnostischer Merkmale vorherzusagen, ob eine Patientin am Beginn einer Diabeteserkrankung steht oder nicht. Welche Patientinnen-Merkmale haben den größten Einfluss auf das Erkrankungsrisiko?
<b>Vokabular</b>	Medizinische Variable (Prädiktoren) im Datensatz sind: <ul style="list-style-type: none"><li>• Pregnancies = Anzahl der Schwangerschaften</li><li>• Glucose = Glukosekonzentration nach Glukosetoleranztest</li><li>• BloodPressure = Blutdruck (mg Hg)</li><li>• SkinThickness = Dicke der Trizephshautfalte (mm)</li><li>• Insulin = Insulin (<math>\mu\text{U}/\text{ml}</math>)</li><li>• BMI = Body Mass Index (in <math>\text{kg}/\text{Körperhöhe in m}^2</math>)</li><li>• DiabetesPedigreeFunction = Diabetes Stammbaumfunktion (bewertet die Wahrscheinlichkeit von Diabetes aufgrund der Familienanamnese)</li><li>• Age = Alter (Jahre)</li><li>• Outcome = nicht erkrankt (=0), erkrankt (=1)</li></ul>
<b>Quelle</b>	<ul style="list-style-type: none"><li>• Smith, J.W., Everhart, J.E., Dickson, W.C., Knowler, W.C., &amp; Johannes, R.S. (1988). <a href="#">Using the ADAP learning algorithm to forecast the onset of diabetes mellitus</a>. In <i>Proceedings of the Symposium on Computer Applications and Medical Care</i> (pp. 261--265). IEEE Computer Society Press.</li><li>• <a href="https://www.kaggle.com/datasets/uciml/pima-indians-diabetes-database">https://www.kaggle.com/datasets/uciml/pima-indians-diabetes-database</a></li></ul>

### Metadaten Datensatz 3: Immobilienpreise – welche Eigenschaften des Umfeldes sind entscheidend?



<b>Daten</b>	<b>Hauspreise_Boston.csv</b>
<b>Inhalte</b>	In diesem Datensatz sind Informationen über die Preise von über 500 Häusern in Boston (USA) enthalten. Zusätzlich finden sich im Datensatz charakteristische Merkmale zum Umfeld der Häuser, wie z.B. Kriminalitätsrate, Luftqualität oder Bevölkerung.
<b>Informationen</b>	Der Datensatz dient der Analyse von Einflussfaktoren auf amerikanische Immobilienpreise.
<b>Fragestellung</b>	Mit welcher Genauigkeit lassen sich aus charakteristischen Umfeldmerkmalen die Hauspreise vorhersagen? Welche Merkmale sind die stärksten Einflussfaktoren?
<b>Vokabular</b>	<ul style="list-style-type: none"> <li>• CRIME: pro Kopf Kriminalitätsrate</li> <li>• ZN: Anteil der Wohnbauflächen für Grundstücke mit einer Fläche von mehr als 25.000 m<sup>2</sup></li> <li>• INDUS: Anteil der Nicht-Einzelhandelsflächen</li> <li>• RIVER: Dummy Variable für Fluss (1 – grenzt an Fluss, 0 - sonst)</li> <li>• AIRQ: Luftqualitätsindikator NOx Konzentration in parts per 10 million</li> <li>• RM: durchschnittliche Anzahl der Zimmer pro Wohnung</li> <li>• AGE: Anteil der vor 1940 gebauten selbstgenutzten Wohneinheiten</li> <li>• DIS: gewichtete Entfernungen zu fünf Beschäftigungszentren</li> <li>• RAD: Index der Erreichbarkeit von Autobahnen und Schnellstraßen</li> <li>• TAX: Vollwertiger Grundsteuersatz pro 10.000 \$</li> <li>• PTRATIO: Index für Schüler-Lehrer Verhältnis</li> <li>• CPRATIO: Index zur ethnischen Zusammensetzung der Bevölkerung</li> <li>• LSTAT: Prozentanteil Bevölkerung mit niedrigem Lebensstatus</li> <li>• PRICE: gemittelter der Eigenheime in US\$</li> </ul>
<b>Quelle</b>	<ul style="list-style-type: none"> <li>• The Boston house-price data of Harrison, D. and Rubinfeld, D.L. 'Hedonic prices and the demand for clean air', J. Environ. Economics &amp; Management, vol.5, 81-102, 1978.</li> <li>• <a href="https://www.kaggle.com/code/alisonmachadolui/boston-house-prices-multiple-regressions/input">https://www.kaggle.com/code/alisonmachadolui/boston-house-prices-multiple-regressions/input</a></li> </ul>