



KI für Alle 2: Verstehen, Bewerten, Reflektieren

Themenblock Prognosemodelle (Klassifikation und Regression): 02_05Daten _Datenexploration

Datenexploration

Erarbeitet von

Dr. Katja Theune

ernziele	1
nhalt	2
Einstieg	
Datenvorverarbeitung	
Visualisierung des Outputs	
Visualisierung der Inputs	
Visualisierung der Zusammenhänge von Output und Inputs	5
Zusammenhänge und Auswahl von Inputs	7
Abschluss	8
Quellen	8
Disclaimer	8

Lernziele

- Du kannst aus Datenvisualisierungen Informationen über die Inputs und Outputs ableiten
- Du kannst aus Datenvisualisierungen Informationen über die Zusammenhänge zwischen Inputs und Outputs ableiten





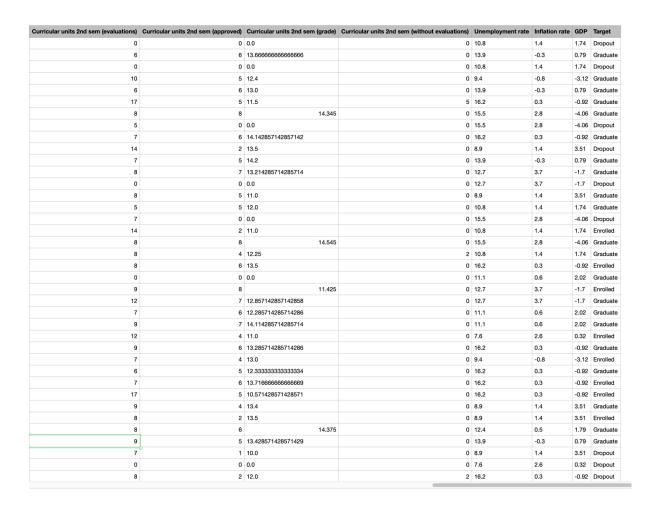


Inhalt

Einstieg

Bevor wir mit den uns vorliegenden Daten Prognosemodelle bilden, müssen wir uns die Daten einmal genauer ansehen. Nur dann bekommen wir wichtige erste Eindrücke davon, ob sich die Daten für unsere Analysen überhaupt eignen. Darüber hinaus erhalten wir durch eine detaillierte Datenexploration schon Hinweise auf Zusammenhänge der Inputs untereinander und zwischen den Inputs und dem Output. Das hilft uns, die Daten passend für unsere weiteren Analysen aufzubereiten.

Datenvorverarbeitung



Der Datensatz, den wir hier verwenden, hat schon einige Schritte der Vorverarbeitung durchlaufen. Aus dem Datensatz wurden z. B. schon alle Fehler, Ausreißer und auch fehlende Werte entfernt, bzw. entsprechend aufbereitet.

Quelle [1]



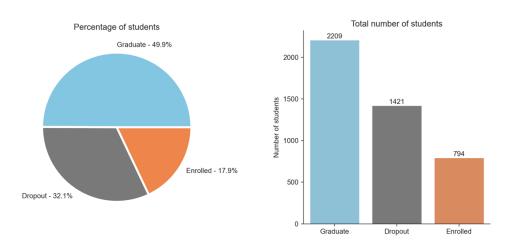




Zudem wissen wir ja schon, dass die meisten Verfahren am besten mit Zahlen umgehen können. Hier sehen wir, dass alle Inputs bereits numerisch sind, die Kategorien also durch Zahlen repräsentiert werden. Nur der Output, hier "Target" genannt, noch nicht. Daher werden wir die Kategorien dann ebenfalls in Zahlen umwandeln.

Visualisierung des Outputs

Jetzt wollen wir uns genauer den Output, verschiedene Inputs und schonmal einige mögliche Zusammenhänge untereinander ansehen. Am besten eignen sich dafür Visualisierungen. Beginnen wir mit unserem Output und stellen die absoluten Häufigkeiten der drei Klassen in einem Säulendiagramm und die relativen Häufigkeiten in % anhand eines Kuchendiagramms dar.



Wir haben hier die Klassen "Graduate", "Dropout" und "Enrolled", also "Studienerfolg", "Studienabbruch" und "noch studierend". Die größte Klasse ist hierbei "Graduate" mit 2209 Studierenden bzw. 49,9 %, gefolgt von "Dropout" mit 1421 Studierenden, bzw. 32,1 %. Die kleinste Gruppe stellen die noch eingeschrieben Studierenden mit 794 Studierenden, bzw. 17,9 %. Insgesamt haben wir 4424 Studierende.

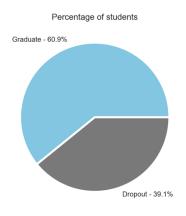
Da wir uns mit einem Klassifikationsproblem mit nur zwei Klassen beschäftigen wollen, entfernen wir alle Beobachtungen, die der Klasse "Enrolled" angehören. Damit verlieren wir natürlich Informationen, aber der Einfachheit halben belassen wir es hier so.

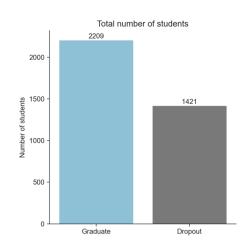
Es bleiben damit 3630 Studierende übrig, von denen 60,9 % der Klasse "Graduate" und 39,1 % der Klasse "Dropout" angehören. Wir haben hier also eine etwas ungleiche Klassenverteilung.





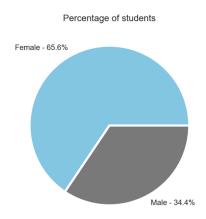






Visualisierung der Inputs

Jetzt schauen wir uns mal beispielhaft zwei der Inputs an. Beginnen wir mit dem "Geschlecht", bzw. "Gender".



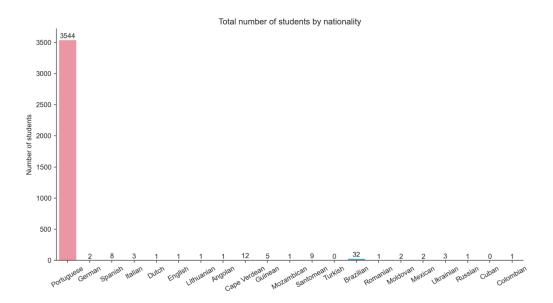
Wir sehen hier im Kuchendiagramm, dass 65,6 % der Befragten Frauen und 34,4 % Männer sind. Weitere Kategorien wurden hier leider nicht abgefragt. Männer sind also unterrepräsentiert, wenn wir grundsätzlich von einer Gleichverteilung der beiden Geschlechter in der Grundgesamtheit ausgehen.

Ein weiterer Input ist die Nationalität der Studierenden. Hier sehen wir, dass die allermeisten, nämlich fast 98 % der Studierenden Portugiesen sind und die anderen Nationalitäten nur sehr wenig besetzt sind.





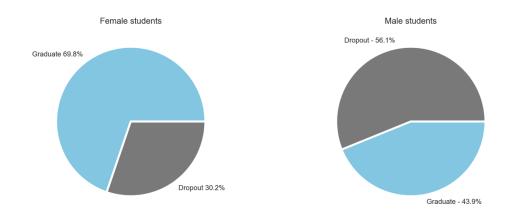




Diese sehr kleinen absoluten Fallzahlen einzelner Kategorien müssen auch bei der Interpretation der Ergebnisse immer mitberücksichtigen werden. Ggf. ist es sinnvoll, Kategorien zusammenzufassen oder den Input nicht mit in die Analysen einzubeziehen.

Visualisierung der Zusammenhänge von Output und Inputs

Als nächstes untersuchen wir, ob und wie Output und Inputs zusammenhängen. Um dafür einen ersten Eindruck zu bekommen, eignet sich zunächst auch wieder eine Visualisierung. Z. B. könnten wir uns ansehen, wie viel Prozent der Frauen im Vergleich zu den Männern "Graduates", oder "Dropouts" sind. Das sehen wir hier im Diagramm.



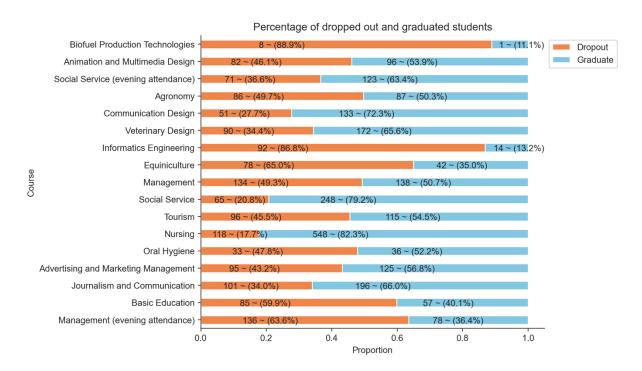
Es fällt auf, dass weibliche Studierende mit 30,2 % eine deutlich niedrigere Studienabbruchquote haben als männliche Studierende mit 56,1 %. Das ist natürlich nur ein erster Hinweis auf einen möglichen Zusammenhang zwischen Geschlecht und Studienabbruch und spiegelt auch keine Kausalität wider.



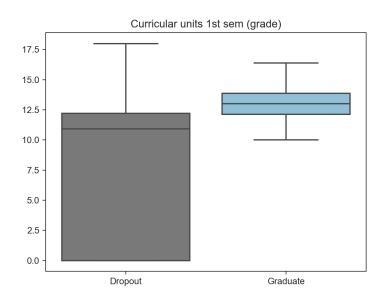




Ähnlich können wir uns einmal die Erfolgs- und Abbruchquoten in den einzelnen Kursen ansehen. Dieser Barplot zeigt, dass die Abbruchquoten zwischen den einzelnen Kursen stark variieren. So haben z. B. "Nursing" mit 17,7 % und "Social Service" mit 20,8 % die niedrigsten Abbruchquoten. Kurse wie z. B. "Biofuel Production Technologies" mit 88,9 % und "Informatics Engineering" mit 86,8 % haben sehr hohe Abbruchquoten. Auch das ist nur ein erster Hinweis auf einen möglichen Zusammenhang zwischen der Fachrichtung und einem Studienabbruch.



Um uns auch mal einen metrischen Input anzuschauen, sehen wir hier einen Boxplot der Noten am Ende des ersten Semesters für Dropouts und Graduates. Man erkennt deutlich, dass die Noten bei Dropouts niedriger, also hier schlechter sind und stark streuen.



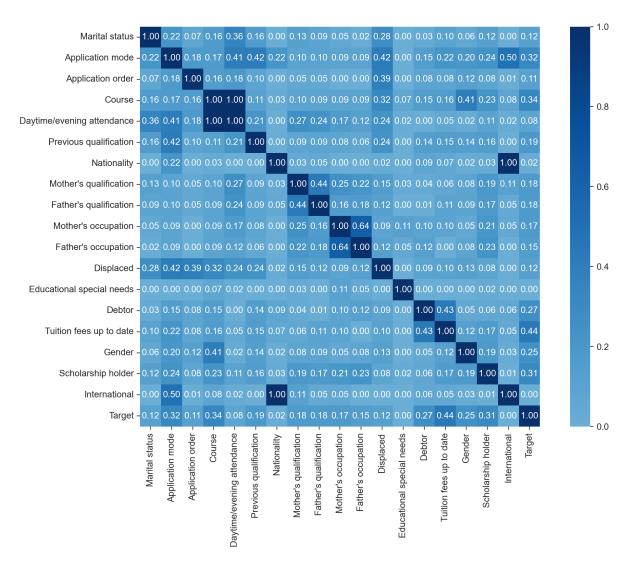




Zusammenhänge und Auswahl von Inputs

Um zu sehen, ob und welche Zusammenhänge es zwischen den einzelnen Inputs und zwischen Inputs und Output gibt, können wir uns eine sogenannte Korrelationsmatrix ansehen, in welcher die einzelnen Korrelationskoeffizienten ausgegeben werden. Wichtig ist hier erstmal nur: Die Farbintensität gibt die Stärke des Zusammenhangs wieder – je dunkler, desto stärker also der Zusammenhang.

Wir sehen hier zur Veranschaulichung nur eine Auswahl der Inputs. Wir können erkennen, dass einige Inputs stark miteinander korrelieren, z. B. der Berufsstatus bzw. die Qualifikation der Mutter und des Vaters. Das bedeutet, dass diese Inputs ähnliche Informationen liefern und davon einige mit dem Ziel entfernt werden können, die Anzahl an Inputs und damit die Dimension zu reduzieren. Man behält meistens den Input mit dem stärkeren Zusammenhang zum Output. Das wäre hier z. B. die Qualifikation der Mutter. Ähnlich gehen wir auch bei den anderen Inputs vor. Zudem entfernen wir Inputs, die nur sehr wenig mit dem Output korrelieren.







Allerdings gibt es hier kein fest vorgeschriebenes Vorgehen. Wir behalten z. B. häufig auch Inputs, die bei unseren Analysen im Fokus stehen, auch wenn sie auf den ersten Blick nicht relevant erscheinen.

Abschluss

Wir haben uns nun einige Details zu unserem Output und verschiedene Inputs angesehen und einen ersten Eindruck der Zusammenhänge in den Daten bekommen. Insbesondere Visualisierungen und Kennzahlen wie Korrelationen liefern hier wichtige Informationen. Diese helfen uns, die Daten passend für weitere Analysen aufzubereiten.

Quellen

Quelle [1] Realinho, V., Machado, J., Baptista, L., Martins, M.V. (2022). Predicting Student Dropout and Academic Success. Data, 7(11), 146. https://doi.org/10.3390/data7110146

Grafiken: Dr. Ludmila Himmelspach

Disclaimer

Transkript zu dem Video "Prognosemodelle (Klassifikation und Regression): Datenexploration", Dr. Katja Theune.

Dieses Transkript wurde im Rahmen des Projekts ai4all des Heine Center for Artificial Intelligence and Data Science (HeiCAD) an der Heinrich-Heine-Universität Düsseldorf unter der Creative Commons Lizenz CC-BY 4.0 veröffentlicht. Ausgenommen von der Lizenz sind die verwendeten Logos, alle in den Quellen ausgewiesenen Fremdmaterialien sowie alle als Quellen gekennzeichneten Elemente.

