

**Projektbericht zum Modul Information Retrieval und  
Visualisierung Sommersemester 2021**

**Visualisierung von verschiedenen Weindaten**

Richard Brennecke

Matrikelnummer:

# Inhaltsverzeichnis

Inhaltsverzeichnis .....	II
1. Einleitung .....	1
1.1 Anwendungshintergrund .....	1
1.2 Zielgruppen .....	3
1.3 Überblick und Beiträge .....	4
2. Daten .....	4
2.1 Technische Breitstellung der Daten .....	5
2.2 Datenvorverarbeitung .....	6
3. Visualisierung .....	7
3.1 Analyse der Anwendungsaufgaben .....	7
3.2 Anforderungen an die Visualisierungen .....	8
3.3 Präsentation der Visualisierung .....	9
3.3.1 Visualisierung Eins .....	9
3.3.2 Visualisierung Zwei .....	10
3.3.3 Visualisierung Drei .....	11
3.4 Interaktion .....	12
4. Implementierung .....	12
5. Anwendungsfälle .....	13
5.1 Anwendung Visualisierung Eins .....	13
5.2 Anwendung Visualisierung Zwei .....	14
5.3 Anwendung Visualisierung Drei .....	16
6. Verwandte Arbeiten .....	17
7. Zusammenfassung und Ausblick .....	18
Literaturverzeichnis .....	i
Anhang .....	ii

# 1. Einleitung

Weltweit wurden 2020 über mehr als 260 Millionen Hektoliter (mhl) Wein produziert [1]. Umgerechnet in Liter sind dies 26.000.000.000 Liter (26 Milliarden Liter). Davon wurden von Deutschland allein 8,4 mhl produziert [1]. Was jedoch von der Führenden Weinproduktionsnation Italien mit 49,1 mhl rund um das neunfache übertroffen wird [1]. Aus diesen Zahlen lässt sich bereits erahnen, dass es Weine in den verschiedensten Arten und Geschmäckern gibt. Aus dieser Vielfalt an Weinen wurde ein Datensatz zusammengestellt mit rund 21 Tausend Datensätzen rund um verschiedene Weine. Diese wurden dabei in den verschiedensten Kategorien bewertet und werden nun mithilfe von verschiedenen Visualisierungstechniken in diesem Projektbericht genauer analysiert.

Dabei ist das Hauptziel aus diesen verschiedenen Daten möglichst viel Erkenntnisse zu generieren, umso vielleicht neue Einsichten in das Thema rund um die Weine zu erhalten. Insbesondere ob es verschiedene Zusammenhänge zwischen den verschiedenen Eigenschaften, welche in dem Datensatz analysiert worden sind, gibt. Darüber hinaus wäre es wichtig zu wissen, ob es auch einen Zusammenhang zwischen den Produktionsmengen von Weinen gibt und den verschiedenen Weinarten innerhalb eines Landes. Zusätzlich sollte überprüft werden, ob es besondere Datensätze gibt, welche in einigen Kategorien besonders hervorstechen. Oder auch ob generelle Trends oder Zusammenhänge zwischen den Kategorien vorhanden sind. Diese Fragen werden anschließend in diesem Bericht mithilfe von den drei Visualisierungstechniken des Scatterplots, Baumhierarchie und Parallele Koordinaten beantwortet.

## 1.1 Anwendungshintergrund

Die Visualisierungstechniken, die in diesem Projektbericht verwendet werden, sind die des Scatterplots, der Parallelen Koordinaten und der Baumhierarchie. Diese werden nachfolgend kurz vorgestellt, um ein besseres Verständnis für diese zu schaffen.

Der Scatterplot stellt dabei zwei verschiedene numerische Variablen gegenüber, diese werden mithilfe von Punkten dargestellt. Die Position dieser Punkte gibt dabei den Wert auf der Horizontalen und vertikalen Achse an. Diese Darstellung ist besonders wichtig, um verschiedene Zusammenhänge zwischen zwei Variablen herauszufinden. So können diese Punkte als Ganzes betrachtet werden, verschiedene Muster anzeigen. So ist es möglich, einfach verschiedene Korrelationsbeziehungen festzustellen. Bei solchen Fällen ist es dann einfach vorherzusagen, wo ein gewisser Horizontaler Wert liegen würde, wenn wir einen vertikalen Wert haben. Darüber hinaus ist es mithilfe eines Scatterplots auch möglich, verschiedene Daten in verschiedene Gruppen zu unterteilen, wenn diese nahe beieinander liegen. So lassen sich Ausreißer oder Lücken in den Daten erkennen, was nützlich ist, wenn die Daten eingeteilt werden sollen. [2]

Bei den Parallelen Koordinaten handelt es sich um einen Ansatz, Mehrdimensionale Daten zu analysieren. Dabei werden die Daten auf verschiedenen Achsen eingezeichnet, wobei jede Eigenschaft als eine Dimension dargestellt wird. Diese Daten werden über die verschiedenen Achsen miteinander verbunden. Die Stärke der parallelen Koordinaten liegt vor allem darin, mehrdimensionale Muster und Vergleiche zu tätigen. Um dies zu erreichen, ist es wichtig, nicht zu denken, dass Linien eine Kodierung von Zeitreihen darstellen (und somit eine Veränderung des Wertes von Zeitpunkt a nach Zeitpunkt b darstellen). Stattdessen stellen eine Linie im Parallelen Koordinatensystem eine Verbindung von einer Reihe von Werten dar. So lässt sich beispielsweise einfacher erkennen, ob sich Werte innerhalb oder außerhalb des Durchschnittes stehen oder aber besondere Ausreißer darstellen. Weiterhin ist es auch so möglich, generelle Aussagen über die

verschiedenen Werte zu treffen, ob beispielsweise einige Werte insgesamt größer sind als andere die restlichen Werte. [3]

Die Baumhierarchie ist eine Darstellung von Daten welche hierarchisch aufgebaut sind. Dabei beginnt die Baumhierarchie bei einem Element und verzweigt sich dabei mindestens zweimal. Diese, durch die Verzweigung entstanden Elemente, können sich wiederum auch verzweigen, um so immer weiter eine Hierarchische Beziehung darzustellen. Das fertige Diagramm ähnelt dabei einen Baum mit seinem Stamm und den Ästen. Dieses Diagramm kann dabei helfen von einer sehr generellen Kategorie mit feinen Detailstufen zu unterteilen. [4]

Diese Visualisierungstechniken visualisieren verschiedene Daten und somit verschiedene Eigenschaften von Wein. Auf diese Eigenschaften wird nun kurz eingegangen, um so ein grundlegendes Verständnis für diese Eigenschaften zu schaffen. Die Eigenschaften, welche in den Visualisierungstechniken verwendet werden, sind Alkohol, Trinktemperatur, Süße, Säuregehalt, Körper, Gerbstoff und Jahr.

Der Alkoholgehalt des Weines entsteht bei der Gehrung des Weins. Dabei handelt es sich um einen natürlichen Prozess, welcher abläuft, nachdem die Trauben zerquetscht worden sind. Je nachdem wie süß und reif die Trauben während der Ernte gewesen sind was einen höheren Zuckergehalt zur Folge hat steigt der Alkoholgehalt des Weines. Der Alkoholgehalt eines Typischen Weines liegt dabei zwischen 9 und 14 Volumenprozent. Falls ein Wein eines höheren Alkoholgehaltes haben sollte, wurde bei der Produktion künstlich Zucker oder Alkohol hinzugefügt. Die Zugabe von Zucker geht dabei auch in der Regel mit einer Minderung der Güteklasse des Weines einher, da der Qualitätsanspruch bei diesen niedriger ist. [5]

Auch die Trinktemperatur hat einen Einfluss darauf wie Wein bei Verkostung schmeckt. Dabei besitzt jeder Wein in der Regel seine eigene ideale Trinktemperatur (welcher meist auf dem Etikett des Weines vermerkt worden ist). Eine Faustregel ist dabei für junge Rotweine eine Trinktemperatur von 14-15 Grad, Kräftige Rotweine hingegeben bei 15-17 Grad und schwere Rotweine bei 17-19 Grad. Falls Rotwein zu warm getrunken werden, kann der Geschmack nicht mehr richtig zur Geltung kommen und der Alkohol überwiegt in dem Wein. Darüber hinaus könnte es sein das dieser Wein brennt im Hals. Bei Weißweinen hingegen liegt die Temperatur niedriger als bei den Rotweinen. So besitzt ein junger Weißwein zwischen neun und elf Grad, ein Reifer Weißwein sollte bei elf bis 13 Grad getrunken werden. Falls Weißwein zu kalt getrunken werden, sollte überwiegt dabei die Säure und die restlichen Aromen des Weines verloren gehen. Da jeder Geschmack individuell ist, ist es auch möglich mit der verschiedenen Trinktemperaturen des Weines zu experimentieren. Als Faustregel gilt dabei auch sobald ein Wein kühler getrunken wird, verstärken sich die Säuren und Tannine wobei sich die Aromatik abschwächt. Wenn jedoch die Trinktemperatur erhöht wird kommt die Aromatik, Körper, Süße und der Alkohol des Weines mehr zum Vorschein. Wenn der Wein serviert wird, ist darüber hinaus in tendenziell etwas kühler als die optimale Trinktemperatur zu servieren, da dieser sich nach dem Ausschanken schnell um zwei bis drei grad erhöht. [6]

Bei der Süße des Weines handelt es sich um eine der wichtigsten Bestandteile des Weines und verleiht mit anderen Komponenten dem Wein erst seinen Geschmack. Die Süße des Weines beschreibt dabei den Restzucker innerhalb des Weines. Dieser kann je nach Anbaugebiet und Rebsorte auch unterschiedlich ausfallen. Dabei gilt aber vor allem je reifer die Trauben sind desto mehr süßer ist der Wein. Innerhalb des EU – Weingesetzes ist festgelegt das ein süßer (oder auch lieblicher) Wein als mindestens eines Restzuckergehalt von 45 Gramm pro Liter aufweisen muss. Darüber hinaus ist es möglich dieses zu erhöhen in dem man im Nachhinein noch Zucker hinzufügt, um die natürliche Süße zu erhöhen. Dieser Vorgang ist als „Aufspriten“ bekannt. [7]

Die Säure innerhalb eines Weines soll ihm die Finesse und Eleganz eines Weines geben. Diese entsteht gemeinsam an der Rebe beim Reifungsprozess dieser. Dabei steigt der Zuckergehalt der Traube an, während der Säuregehalt sinkt. Das Verhältnis von dieser zu und Abnahme ist dabei nicht immer gleich. So können kühle Nächte den Säurerückgang verzögern. Im Wein sind besonders zwei verschiedene Säurearten anzutreffen, welche die Weinsäure und die Apfelsäure sind. Dabei ist die Weinsäure eine meist erwünschte Säure da diese den Wein weich und angenehm schmecken macht. Die Apfelsäure hingegen macht den Wein kantig und hart. Dies kann bei Weißweinen zum Teil erwünscht sein, bei Rotweinen jedoch wird diese meistens in Milchsäure umgewandelt. Darüber hinaus ist der Apfelsäuregehalt auch ein guter Indikator für die Qualität des Jahrganges, da dieser von der Witterung abhängt. So wird weniger Apfelsäure bei sonnigen Jahrgängen produziert und in kühlen Jahrgängen mehr. [8]

Bei dem Körper des Weines handelt es sich um eine Möglichkeit den Eindruck des Weines zu beschreiben welcher dieser im Mund hinterlässt. Der Körper beschreibt somit die Gerbstoffe, Restsüße und die Säure des Weins. Dabei wird dieser meistens in einen leichten, mittleren und voluminösen Körper unterschieden was zu der Bezeichnung leichtgewichtigem, mittelgewichtigem und schwergewichtigem Wein führt. So sind vor allem die Gerbstoffe für den Körper verantwortlich. So führt meistens ein höherer Gerbstoffgehalt zu einem voluminöseren Körper. Aber auch wenn der Alkohol kein Bestandteil des Körpers ist, spielt dieser beim Körper eine Rolle, denn je höher der Alkoholgehalt des Weines ist desto voluminöser erscheint dieser auch im Mund. [9]

Die Gerbstoffe (auch Tannine genannt) kommen primär im Wein vor bei dem die ganze Traube verarbeitet wird. Dies liegt daran, dass die meisten Gerbstoffe in dem Kern der Traube und in der Schale vorzufinden sind. Dementsprechend besitzt ein Rotwein mehr als ein Weißwein, da hierbei die gesamte Traube verarbeitet wird. Jedoch können auch bei der Gröhrung in Holzfässern weitere Gerbstoffe mit hinzukommen. Beim Trinken des Weines rufen die Gerbstoffe den herben bis bitteren Geschmack und das Gefühl sich die Mundschleimhäute zusammenziehen oder der Mund trocken ist, hervor. Je nach Stärke des Weines sind die Gerbstoffe im hinteren Bereich des Mundes zu finden, bei besonders Kräftigen Gerbstoffen auch im gesamten Mund. Zu viele Gerbstoffe verleihen dem Wein einen eher starken herben oder bitteren Geschmack, zu wenig Gerbstoffe führt jedoch zu einem eher flachen Geschmack. Bei der Lagerung kann es darüber hinaus sein dass sich die Gerbstoffe Moleküle zu größeren zusammenfügen und somit dem Wein seinem besonderen Geschmack verleihen. [10]

Ein großer Unterschied von Wein zu vielen anderen Getränken ist dass dieser weiterhin in der Flasche reifen kann. Dies bedeutet dass Weine in der Lage sind mit der Zeit an Qualität zu gewinnen. Jedoch haben die meisten Weine ihren Trinkhöhepunkt schon nach einigen Jahren ihrer Abfüllung erreicht und sollten deswegen nicht übermäßig lange gelagert werden. Da ansonsten mit der Zeit Sauerstoff an dem Verschluss des Weines vorbeikommt und im schlimmsten Falle Essig bildet. Dabei hängt es vom Gerbstoffgehalt des Weines ab, wie viel Sauerstoff er verträgt. Dieser macht die Gerbstoffe nämlich weicher und somit kann sich die eine gute Frucht innerhalb des Weines entwickeln. Somit eignen sich tendenziell Rotweine mehr zum Lagern als Weißweine, da diese weniger Gerbstoffe beinhalten. Beispielsweise ist einer der ältesten Weine ein 1870er Chateau Lafite ein Rotwein. [11]

## 1.2 Zielgruppen

Für die verschiedenen Visualisierungstechniken gibt es drei verschiedene Zielgruppen, welche potenziell an diesen einen Mehrwert finden könnten. Diese Zielgruppen sind: Weininteressierte, Weineinkäufer (-verkäufer) und der Weinexperte. Auf diese wird nachfolgend nun eingegangen.

Die Weininteressierten sind eine Gruppe welche gerne Wein trinken, welche jedoch aber kaum bis gar kein Vorwissen zu dem Thema Wein besitzen. Deswegen könnten diese mithilfe von

verschiedenen Visualisierungen neue Erkenntnisse gewinnen rund um das Thema der Weine. So könnten Sie die Zusammenhänge zwischen der Säure und der Süße des Weines erkennen. Darüber hinaus könnten sie mithilfe dieser Visualisierung die Weine, die diese bisher getrunken haben, viel besser einordnen und somit vielleicht auch neuen Weine entdecken, welche ihnen potenziell schmecken könnte.

Bei der Gruppe der Weineinkäufer bzw. Weinverkäufer handelt es sich um die Gruppe, welche das meiste Wissen über Weine besitzen sollte. Da diese ihre Kunden betrauten sollten und entsprechende gute Weine zu guten Preisen finden und einkaufen müssen. Somit ist ihr Vorwissen sehr gut bis exzellent. Mithilfe der Visualisierungen wäre es möglich verschiedene neue Sorten zu entdecken welche in das Sortiment des einzelnen Händlers passen könnten. Darüber hinaus könnten neue oder außergewöhnliche Sorten entdeckt werden, welche vielleicht nur spezielle Kunden infrage kommen würde. Weiterhin könnte diese zur Kundenberatung verwendet werden, um die speziellen Geschmäcker aller Kunden effektiv abdecken zu können.

Die Weinexperten sind eine Gruppe von Leuten welche gerne Wein trinken und sich bereits mit diesem Thema auseinandergesetzt haben. Dementsprechend sollte ihr Wissenstand durchschnittlich bis gut sein. Diese Gruppe könnte mithilfe der Visualisierungen neue Sorten für sich entdecken, welche ihren Geschmack entspricht und so vielleicht noch tiefer in die verschiedenen Weinrichtungen einzutauchen. Weiterhin wäre es auch möglich mithilfe der verschiedenen Visualisierungen ihren Weingeschmack weiter zu erforschen, da diese mithilfe der Darstellungen diesen besser erforschen könnten.

### 1.3 Überblick und Beiträge

In diesem Projekt wurden die Daten von der Webseite Kaggle verwendet. Dabei handelt es sich um Daten rund um das Thema Wein. Dabei sind in den Daten die verschiedenen Namen der Weine zu finden, deren Produzenten und woher diese sind (inkl. der einzelnen Standorte). Weiterhin sind zu dem Wein an sich weitere Informationen vorhanden. So ist der Typ und die Verwendung des Weines zu finden. Weiterhin sind in diesen Daten Eigenschaften wie Alkoholgehalt, Trinktemperatur, Süße, Säure, Körper, Gerbstoffe, Preis, Jahr und Größe einer Flasche vermerkt. Genauere Informationen was diese Eigenschaften bedeuten ist im Kapitel Anwendungshintergrund zu finden. Diese Daten werden anschließend anhand ihrer verschiedenen Eigenschaften an einem Scatterplot, Parallelen Koordinaten und einer Baumhierarchie dargestellt. Wie diese Diagramme aufgebaut sind, ist im Kapitel Anwendungshintergrund zu finden. Dabei ist es möglich mithilfe des Scatterplots zwei verschiedene Eigenschaften des Weines zu vergleichen, um entsprechende Muster zwischen diesen beiden Eigenschaften zu erkennen und identifizieren. Dies könnte zu neuen Erkenntnissen rund um diese Eigenschaften führen. Mithilfe der Parallelen Koordinaten können die verschiedenen Paare an Eigenschaften verglichen werden, um so für ein einzelnes Datenpaar herauszufinden, wie es sich gegenüber den anderen Datenpaaren verhält. Somit könnten schnell besondere Datenpaare herausgefiltert werden und entsprechende neue Informationen gewonnen werden. Mithilfe der Baumhierarchie wird es möglich sein die verschiedenen Weine ihren Regionen zuzuordnen und zu herausfinden wie viele Weine pro Region vorhanden sind. Diese Informationen kann dann wiederum mit anderen Daten abgeglichen werden, um einen weiteren Erkenntnisgewinn zu ermöglichen.

## 2. Daten

Die Originaldaten, welche auf der Plattform Kaggle zu finden sind, wurden von einem Nutzer von einer Koreanischen Webseite (welche nicht genauer angegeben worden ist) gesammelt und bereitgestellt. Diese Originaldatei enthält 32 Spalten mit insgesamt 21605 Datensätzen. Da diese

Datensätze teilweise Koreanische Symbole enthielten wurde eine zweite Datei angelegt, um die entsprechenden Zeichen herauszufiltern und den Datensatz somit zu bereinigen. Diese Datei heißt „cleasingWine.csv“ und besitzt 31 Spalten und 21600 Datensätze. Auf dieser Grundlage wurden die nachfolgend in diesem Kapitel beschriebenen Datenbearbeitungsschritte vorgenommen. [12]

Der Datensatz von CleasingWine beginnt mit der „wineID“ welche eine einfache Index Nummer darstellt für die einzelnen Daten. Nach dieser Spalte folgt die Spalte mit dem Namen „name“, in dieser Spalte sind entsprechend der ganzen Namen der einzelnen Weine mit eingetragen. Die Produzenten dieser Weine sind in der darauffolgenden Spalte der „producer“ zu finden. Aus welchem Land dieser Wein stammt wird anschließend in der Spalte der „nation“ beantwortet. Anschließend folgen fünf Spalten mit dem Namen „local“ und der entsprechenden Nummer, welche die entsprechende Region wo der Wein herkommt, darstellen. Anschließend sind finden sich die Spalten „varieties“ und die entsprechende Nummer, welche bis zwölf geht, um die Weine ihre entsprechende Sorte zuordnen zu können. Nachdem diesen Spalten folgt die Spalte „typ“ welche den entsprechenden Typ des Weines enthält. Anschließend wird in der Spalte „use“ die Verwendung des Weines genauer festgelegt. Auf diese Spalte folgt, die der „abv“ welche den Alkoholgehalt pro Volumen darstellt. Danach wird in der Spalte der „degree“ die optimale Trinktemperatur des Weines festgehalten. Darauffolgend wird mit den Spalten „sweet“, „acidity“, „body“, „tannin“, „price“, „year“, „ml“ die jeweilige Süße, Säure, Körper, Gerbstoff, Preis, Herstellungsjahr und Größe der Flasche des Weines definiert. Dabei sind ab der Spalte der „nation“ leer Werte in der Datei enthalten. Zahlenwerte sind aber der Spalte „abv“ in diesem Datensatz zu finden. [12]

Dabei eignen sich diese Daten für die Zielgruppe der Weininteressierten gut. Da diese ausreichend tiefe bieten um neue Weine kennen zu lernen und auch Zusammenhänge zwischen den Weinen Eigenschaften zu erkennen. Für die Gruppen der Weineinkäufer oder Weinexperten jedoch ist dieser Datensatz ausreichend. Da hierbei vor allem durch die vielen leeren Werte keine vollständige Datenbasis vorhanden ist. Somit ist nicht garantiert das für jeden Wein der ggf. entdeckt wird die entsprechenden Daten vorhanden sind, die ggf. für eine Entscheidung benötigt werden ausreichen. Weiterhin ist die Tiefe der Daten für diese Gruppen nicht ausreichend, da die Eigenschaften der Weine (Süße, Säure, Körper und Gerbstoffe) nur auf einer Skala von ein bis fünf angegeben werden. Dies führt dazu, dass diese Daten ungenau sind, wobei diese hätten, genauer sein können und somit weniger geeignet für Gruppen welche genauen Daten zu den Weinen wünschen würden.

Dabei werden reichen diese Daten jedoch aus um neue Erkenntnisse rund um das Thema des Weins und ihren unterschiedlichen Eigenschaften zu generieren. Dementsprechend reichen diese Daten auch aus, um die Fragestellungen der Zielgruppen zu beantworten, wenn auch nicht so ausführlich wie diese es sich wünschen würden. Dabei könnten die Weineinkäufer und Weinexperten diese Daten mehr als Richtwerte nehmen um sich dann entsprechend weiter Recherchieren.

Darüber hinaus wurde noch für das Baumdiagramm eine weitere Datei erstellt, um die Länder welche nur in der Datei vorhanden, warum um die Kontinente und Einteilungen dieser zu ergänzen. Um so eine bessere Hierarchische Darstellung der Weine und ihrer Länder zu erhalten.

## 2.1 Technische Breitstellung der Daten

Die Technische Bereitstellung der Daten erfolgt mithilfe des GitHubs Repositorie, in welchem das Visualisierungsprojekt umgesetzt worden ist. Dabei sind die Daten innerhalb des Ordners „Daten“ zu finden. In diesem Ordner sind unter dem Unterordner „Quelldaten“ sind die Ursprünglichen Daten vorhanden welche von der Plattform „Kaggle“ heruntergeladen werden konnten. Eine genauere Beschreibung dieser Daten ist im Kapitel Daten zu finden. Im anderen Unterordner namens „Aufbereitete Daten“ sind alle Daten zu finden welche weiterverarbeitet wurden und entsprechend

selektiert wurden. Wie die Weiterverarbeitung erfolgte, ist im Nachfolgendem Kapitel Datenvorverarbeitung zu lesen.

Dabei wurden zwei wesentliche Dateiformate für die Bereitstellung der Datei verwendet. Diese sind CSV- und JSON-Dateien. Die CSV-Datei, auf welche der Scatterplot und die Parallelen Koordinaten zugreifen heißt dabei „WineInformationExcelAufbereitetKlein“. Der Name der JSON Datei auf, welche das Baumdiagramm zugreift, lautet „WineInformationGeoKleinKlein“. Dabei sind alle Daten innerhalb der CSV Datei mithilfe eines Kommas getrennt. Darüber hinaus werden alle Zahlen mit einem Dezimaltrennzeichen als einen Punkt angegeben. Falls eine Null in den Feldern stehen sollte, ist dies gleichbedeutend mit einem Feld welches Leer ist. In den JSON Datei hingegen wurden nur die Beziehungen zwischen den Daten abgebildet. Somit enthalten die „data“-Felder nur eine „id“ welches nur den Namen enthält. Die Beziehung wurde mithilfe der „children“-Felder realisiert. Darüber hinaus wurden in der finalen JSON Datei alle Länder entfernt, welche nach der CSV-Datei keine Weine produzieren.

## 2.2 Datenvorverarbeitung

Bei der Datenverarbeitung wurden im Wesentlichen drei Schritten durchgeführt, um die Daten weiter zu bearbeiten. Diese Schritte sind das Sichten-, Bearbeiten- und anschließende Überführung der Daten. Was in diesen Schritten genauer geschehen ist, wird nachfolgend erklärt.

Bei der Sichtung der Daten war es das Ziel die Daten in ein Lesbares und leicht zu verarbeitendes Format zu bringen. Aufgrund dessen wurde die Datei „cleasingWine“ in eine Exceldatei konvertiert, da eine Exceldatei genau die Zielstellung dieses Schrittes erfüllt. Diese Exceldatei ist unter dem Unterordner „Aufbereitete Daten“ mit dem Namen „WineInformationExcel“ zu finden. Durch diesen Schritt konnten diese Dateien nun einfacher eingelesen und bearbeitet werden. Bei der JSON Datei wurde hierbei eine entsprechende zu ergänzende Datei gefunden, welche im nächsten Schritt dann ergänzt werden konnte. Diese Datei wurde durch den GitHub Nutzer „Curran Kelleher“ zur Verfügung gestellt [13]. Die Datei ist dabei unter „Aufbereitete Daten“ mit dem Namen „WineInformationGeo“ zu finden.

In der Bearbeitung der Daten sollten diese so bereitgestellt werden, dass sie von Code, welcher eingesetzt wurde, einfach zu verarbeiten sind. Dabei wurden sechs wesentliche Schritte getätigt um diese Daten entsprechend bereitzustellen. Zuerst wurden die Namen aus den Spalten „sweet“, „acidity“, „body“ und „tannin“ vor den Zahlen entfernt, da diese ansonsten verhindert hätten die Zahlen in ein entsprechendes Datenformat zu bringen. Anschließend wurden die Zahlenwerte überarbeitet. So standen in der Spalte „abv“ und „degree“ eine Tilde, welche die minimale und maximalen Werte miteinander verbunden hat. Da jedoch kein Zahlenformat eine Tilde zulässt wurden diese beiden Werte getrennt und anschließend aus den beiden Werten der Durchschnitt gebildet. Falls nur ein Wert bereits in der Spalte stand, wurde dieser einfach übernommen. Darüber hinaus wurde der Preis für die Weine in der Spalte „price“ in südkoreanischen Won angegeben. Aufgrund dessen wurde dieser Preis mithilfe eines Währungskurses von 1 Euro zu 1355,382 Won umgerechnet. [] Nachdem diese Zahlen erfolgreich überarbeitet worden sind, wurden die Namen in der Spalte „name“ angepasst. So ist es bei der Konvertierung vorgekommen, dass Apostrophe und Umlaute nicht richtig übersetzt worden sind. Diese Fehler wurden korrigiert. Um die Daten anschließend besser in die JSON Datei einfügen zu können wurde außerdem eine neue Spalte erstellt, in welcher die Namen der Weine mit der entsprechenden JSON Schreibweise kombiniert wurden, um so eine zusammenfügen der Weinamen und der Länder einfacher zu gestalten. Diese Änderungen wurden alle in der Exceldatei „WineInformationExcel“ durchgeführt. Anschließend wurden in einer neuen auf „WineInformationExcel“ basierenden Datei, die Namen der Spalten von Englischen ins Deutsche übersetzt. Zusätzlich dazu wurden die Spalte in der Excel „Column1“ und die



Spalte mit den Schreibweisen für die JSON Datei herausgelöst. Dieser Stand präsentiert findet sich in der Exceldatei mit dem Namen „WineInformationExcelAufbereitet“ wieder. Zum Abschluss dieser Bearbeitung der Daten wurden alle Datensätze herausgelöscht welche in den Spalten „alc“ bis „ml“ ein leeres Feld oder eine Null als Werte enthielten entfernt. Diese Änderung ist in der Datei mit dem Namen „WineInformationExcelAufbereitetKlein“ wiederzufinden. Innerhalb der JSON Datei gab es nur zwei Schritte welche zur Bearbeitung der JSON Dateien. Innerhalb des ersten Schritts wurden alle Weinnamen dort eingefügt, wo diese bereits auch in der CSV eine Zuordnung zu einem Land erhalten haben. Dieses Ergebnis dieses Schritts ist unter der Datei „WineInformationGeoKlein“ zu finden. Anschließend wurden im letzten Schritt alle Länder entfernt, welche keine Weine nach der CSV-Datei hergestellt haben. Dies ist in der Datei „WineInformationGeoKleinKlein“ zu erkennen.

Innerhalb der Überführung der Daten sollten die Daten wieder für den Programmcode lesbar gemacht werden. Somit wurden alle Daten, welche in einer Exceldatei gespeichert waren, wieder in eine CSV-Datei zurücküberführt. Dabei wurde darauf geachtet, dass die Daten mithilfe eines Kommas getrennt worden sind und das Dezimaltrennzeichen ein Punkt war. Da die JSON Datei nicht umgewandelt wurde, war hierbei auch keine Überführung der Daten in ein anderes Datenformat nicht nötig.

Im Schritt „Bearbeitung der Daten“ wurden die verschiedenen Durchschnitte gebildet, um die Daten besser lesbar sowohl für die Menschen als auch die Maschine zu machen. So hätten ansonsten noch mehr Kategorien zur Auswahl gestanden, welche durch die Differenz der Zahlen keinen großen Mehrwert für die Visualisierungen gebracht hätte. Weiterhin wurden verschiedene Datensätze in diesem Schritt entfernt. Dies wurde getan, um eine gemeinsame konsistente Datenbasis für alle Visualisierungen zu schaffen, so jeder Wein in allen Visualisierungen zu finden. Darüber hinaus führte die Reduktion der Datensätze dazu, dass die Visualisierungen nicht überfüllt wirken. Zusätzlich wurden alle Datensätze entfernt, welche nicht ins Deutsche übersetzt werden konnten, da diese noch in Koreanischer Sprache geschrieben waren und somit die Weine nicht mehr eindeutig identifiziert werden konnten.

### 3. Visualisierung

In dem nachfolgenden Kapitel wird genauer auf die einzelnen Visualisierungen und ihren Zweck eingegangen.

#### 3.1 Analyse der Anwendungsaufgaben

Wie bereits im Kapitel Einleitung erwähnt worden ist, gibt es ein großes Hauptziel, welches mit diesen Visualisierungen erreicht werden soll. Dieses ist möglichst viele Erkenntnisse zu gewinnen, rund um das Thema der Weine. Dabei kann das Thema vor allem aus den Eigenschaften des Weines beleuchtet werden, welche innerhalb der ursprünglichen CSV Datei vorhanden waren. Also an den Eigenschaften des Alkoholgehaltes, Trinktemperatur, Süße, Säure, Körpers, Gerbstoffe, Preis, Jahres und Größe der Flasche. Diese Eigenschaften können nun mithilfe von den verschiedenen Visualisierungen gegenübergestellt werden und so sollten neue Erkenntnisse rund um das Thema Wein gewonnen werden können. Als Unterstützung für dieses Hauptziel gibt es noch drei etwas konkretere Fragen welche bereits in der Einleitung genannt worden sind. Auf diese wird nachfolgend eingegangen.

Um möglicherweise verschiedene Zusammenhänge zwischen verschiedenen Eigenschaften herauszufinden, wird eine Darstellung benötigt, in welcher die verschiedenen Eigenschaften des Weines gegenübergestellt werden kann. Darüber hinaus sollte ein gewisses Vorwissen über verschiedene Weine bereits vorhanden sein, um die angezeigten Daten interpretieren zu können. Mit

dieser Kombination ist es nun möglich verschiedene Erkenntnisse, welche bereits bekannt sind zu bestätigen oder auch neue Erkenntnisse aus diesen Daten zu ziehen.

Bei der Frage, ob es einen Zusammenhang zwischen den Produktionsmengen und der Anzahl der von Weinen in einem Land, ist es nötig diese zwei Informationen bereitzustellen und anschließend miteinander zu vergleichen. Dafür ist die Information der einzelne Produktionsmengen der Länder nötig. Diese Information ist beispielsweise im „State of the Vitivinicultural World in 2020“ nachzulesen [1]. So finden sich unter dem Top 3 der am meisten Produzierenden Ländern Italien mit 49,1 mhl, Frankreich mit 46,6 mhl und Spanien mit 40,7 mhl. Diese Angaben sollten nun mit einer Darstellung verglichen werden aus der die Weine entsprechend ihren Ländern zugeordnet werden können, um so einen Abgleich mit diesen Zahlen durchführen zu können.

In der letzten unterstützenden Frage geht um besondere Datensätze und verschiedene Trends, die sich in diesem gesamten Datensatz erkennen lassen. Um diese Frage beantworten wird eine Darstellung benötigt, in welcher die verschiedenen Eigenschaften dargestellt werden können, aber gleichzeitig die verschiedenen Datensätze noch unterscheidbar sind. Weiterhin sollte die Darstellung so aufgebaut das gewisse Trends in den Daten erkennbar sind. Gemeinsam mit einem gewissen Vorwissen, was besonders für die einzelnen Eigenschaften ist, könnte diese Frage beantwortet werden.

Die eingesetzten Darstellungen des Scatterplots, Parallelen Koordinaten und Baumdiagramms können dabei helfen dem Betrachteten das Verstehen der Daten vereinfachen. So ist es für diesen möglich sich mithilfe der Parallelen Koordinationen einen Überblick über die gesamten Datensatz zu erhalten und bereits gewisse Trends zu erkennen. Falls dieser anschließend tiefere Analyse und Darstellung dieser Daten haben möchte kann, dieser zu dem Scatterplot greifen. Hierbei ist es möglich zwei Eigenschaften der Daten gegenüberzustellen und diese Gegenüberstellung weitere Analysen oder Trends zu erkennen. Weiterhin ist es möglich mithilfe des Baumdiagramms herauszufinden was die Herkunft dieser Weine ist, wenn es dort gewisse Präferenzen vom Betrachter geben sollte.

Diese Struktur der Darstellungen ist dabei nicht zwingen nötig um die verschiedenen oben genauer beschrieben Fragen beantworten zu können. Da diese verschiedenen Fragen sich mit einer Darstellung beantworten lassen. Somit muss es keine Vernetzung dieser Darstellungen geben. Es könnte jedoch sein, dass eine solche Struktur gebraucht werden kann, um die Hauptfrage zu beantworten, da Erkenntnisse gegeben falls auch zwischen den einzelnen Darstellungen rekombiniert werden können. Dann wäre eine solche Strukturierung sinnvoll und würde bei der Erkenntnisgewinn hilfreich sein.

### 3.2 Anforderungen an die Visualisierungen

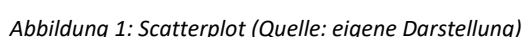
Durch die Analyse der verschiedenen Ziele ergeben sich verschiedene Anforderungen an die Darstellungen. Für das Hauptziel ist es wichtig das die Visualisierungen so dargestellt sind, dass die Darstellungen entsprechend einfach verständlich und es sich trotzdem Erkenntnisse zu dem Thema Weine ziehen lassen. Darüber hinaus sollten die verschiedenen Eigenschaften innerhalb der Anwendungen integriert sein.

Die Anforderungen aus der Frage nach den Zusammenhängen zwischen Eigenschaften sind das sich verschiedene Eigenschaften miteinander vergleichen lassen. Dementsprechend sollten es eine Auswahl geben, in welcher der Anwender nach seinen individuellen Wünschen Eigenschaften auswählen kann, welche Anschließend automatisch angezeigt werden. Darüber hinaus sollten sich in der Darstellung verschieden Trends und Besonderheiten erkennen lassen.

In der letzten Frage, welche sich um die Trends und besondere Datensätze herausfinden möchte, sollte die Darstellung es ermöglichen verschiedene Eigenschaften mit dem gesamten Datensatz darzustellen. Dabei sollte jeder Datensatz trotzdem noch nachverfolgbar sein, um diesen gegeben falls als besonderen Datensatz zu identifizieren. Darüber hinaus sollte die Darstellung des Datensatz so gut erfolgen, dass sich noch Trends aus der Darstellung ablesen lassen, und dies nicht untergehen aufgrund einer beispielsweise gedrängten Darstellungsweise.

Nachfolgend werden die verschiedenen Visualisierungen vorgestellt, welche bei der Realisierung des Projektes verwenden worden sind. Dies sind der Scatterplot, Parallele Koordinaten und die Buamdiagramm.

Die erste Visualisierung innerhalb dieses Projektes ist ein Scatterplot, in welchem immer zwei verschiedene Eigenschaften des Weines gegenübergestellt werden. Dabei nimmt immer eine Eigenschaft eine Achse des Scatterplots ein und anschließend werden die Werte wie XY Koordinaten in das entstandene Koordinatensystem eingetragen. Die dabei entstandenen Punkte werden in dieser Visualisierung als Kreise dargestellt. Wenn mit der Maus über diese Kreise gefahren wird, werden farblich und der Text, welcher über diesen Kreis auftaucht, zeigt den Namen des Weins, die X und Y Eigenschaft dieses Punktes an. Darüber hinaus können die Eigenschaften angepasst werden. Dafür sind über dem Diagramm verschiedene Buttons zu finden welche die Eigenschaften beinhalten, welche es insgesamt in diesem Scatterplot dargestellt werden können. Wenn diese Buttons angeklickt werden, ändert sich je nachdem in welcher Reihe der Button angeklickt wurde die jeweilige Achse des Scatterplots. Der Scatterplot ist in Abbildung 1 zu erkennen.



Die Anforderungen an den Scatterplot, welche es gibt konnten erfüllt werden. So kann der Scatterplot zwei Eigenschaften des Weines gegenüberstellen. Darüber hinaus können die Angezeigten Eigenschaften mithilfe der Buttons geändert werden und diese Veränderung passiert auch komplett ohne neu laden der Seite. Weiterhin lassen sich mithilfe der Darstellung der Koordinaten als Kreise entsprechende Trends oder Besonderheiten in den Daten gut erkennen.

In den Anforderungen dieses Projektes, sollte die erste Darstellung ein Scatterplot oder ein Zeitreihendiagramm werden. Da jedoch die Zeitreihendiagramme, immer eine Zeitliche Dimension benötigen, um Veränderungen über die Zeit darzustellen, was jedoch nicht mit diesen Daten gegeben war. So sind die einzigen Zeitdaten die der Herstellung des Weines. Dabei handelt es sich jedoch um nicht um eine zeitliche Veränderung. Dementsprechend eignen sich diese Daten nicht um über ein Zeitreihendiagramm dargestellt zu werden. Darüber hinaus lassen sich in einem Zeitreihendiagramm, zu dem beispielsweise ein Liniendiagramme gehören können, nur schwer zwei Eigenschaften gegenüberstellen. Da beide Eigenschaften dafür eine Zeitlichen verlauf benötigen würden. Da dies jedoch nicht der Fall war wurde, bei dieser Darstellung auf den Scatterplot zurückgegriffen. Mithilfe von diesem Diagramm lassen sich zwei Eigenschaften gegenüberstellen und dies unabhängig von einer Zeitlichen Achse.

### 3.3.2 Visualisierung Zwei

Bei der zweiten Visualisierung handelt es sich um einen Parallele Koordinaten Diagramm. In einem solchen Diagramm werden die Eigenschaften der Daten als Achse dargestellt. Diese Achsen werden durch Linien miteinander verbunden. Eine Linie stellt dabei einen Datensatz dar und dessen werten auf den verschiedenen Achsen. So gibt es in dieser Darstellung vier verschiedene Achsen, welche alle die verschiedenen Eigenschaften der Weine darstellen. Die dargestellten Linien sind dabei die unterschiedlichen Weine. Die Achsen können dabei individuell den Eigenschaften zugewiesen werden. Dies erfolgt mithilfe von verschiedenen Buttons, welche über der Darstellung zu finden sind. Wenn dabei ein Button gedrückt wird für die entsprechende Achse, nimmt diese die gewünschte Eigenschaft des Nutzers an. Somit sind alle Achsen individuell einstellbar. Die Parallelen Koordinaten sind in Abbildung 2 zu erkennen.

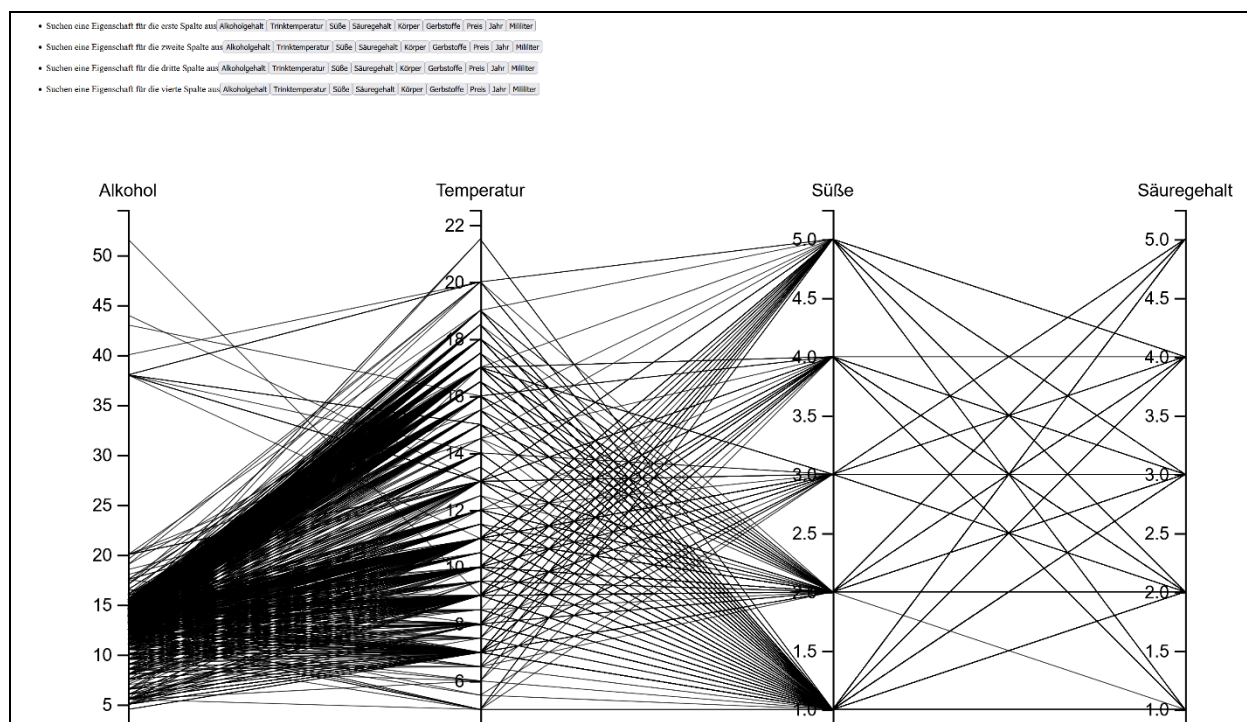


Abbildung 2: Parallele Koordinaten (Quelle: eigene Darstellung)

Die Anforderungen an die diese Darstellung, welche sich aus der Hilfefrage nach den Trends und besonderen Daten gibt, wurden teilweise erfüllt. So ist es mithilfe dieses Diagramms möglich die verschiedenen Eigenschaften gegenüberzustellen und daraus verschiedenen Trends abzuleiten, da diese so Übersicht dargestellt werden können. Jedoch ist es in diesem Diagramm nicht möglich über eine einzelne Linie drüber zu fahren und den entsprechenden Wert angezeigt bekommen. Somit ist die nach Verfolgbarkeit der Daten nur eingeschränkt möglich. Durch das Erfüllen der Übersichtlichkeit und erkennbar machen von Trends und das erschwerte Nachverfolgenden der Daten sind die Anforderungen an dieses Diagramm nur Teilweise erfüllt.

Innerhalb der Anforderungen für das Projekt sollte die zweite Anforderung eine Darstellung von Mehrdimensionalen Darstellungen sein. Für diese Darstellungen würden die Scatterplots, Projektion und Selektion, Parallelen Koordinaten, K-Means und Datentinte in Frage kommen. Bei der Datentinte wird versucht sich auf die Hauptaussage der Visualisierung zu konzentrieren und versucht somit den Teil der Visualisierung zu löschen welcher Verlustfrei gelöscht werden kann. Da jedoch aufgrund der verschiedenen Zielgruppen (vor allem der Weineinkäufer und Weinexperten) keine Daten verloren gehen sollen fällt diese für die Darstellung von Mehrdimensionalen Daten heraus. Dasselbe gilt für die K-Means welche versuchen den Abstand (quadriert) der Punkte zu den Prototypen zu minimieren. Dies würde mit einem Datenverlust einhergehen, was durch die Zielgruppen nicht erwünscht ist. Bei der Projektion und Selektion wird versucht ein mehrdimensionalen raum zu Visualisieren. Dafür wird dieser mithilfe von Projektionen auf verschiedenen 2D Darstellungen abgebildet. Da diese Darstellung jedoch viele Darstellungen nach sich ziehen würde, wäre dies nicht im Interesse und einfachen Erkennbarkeit von Besonderen Daten und Trends und widerspricht somit einer Anforderung an dieses Diagramm. Auch bei den Scatterplots müssen die Mehrdimensionalen Daten auf verschiedene Scatterplots aufgeteilt werden und widersprechen somit auch wie die Projektion und Selektion der Übersicht. Die Parallelen Koordinaten können mehrere Dimensionen in einer Darstellung darstellen und es trotzdem erreichen das gewisse Trends und besondere Daten einfach erkennbar bleiben. Deswegen wurde diese Darstellung mithilfe dieser umgesetzt.

### 3.3.3 Visualisierung Drei

Bei der letzten und somit dritten Visualisierung handelt es sich um Baumdiagramm. Mithilfe eines solchen Diagramms können verschiedene Hierarchische Beziehungen dargestellt werden. So sind diese meistens untereinander angeordnet und mithilfe von Linien verbunden. In der aktuellen Darstellung werden die Länder und Weinnamen mithilfe von Kreisen dargestellt. Diese werden anschließend je nach Beziehung mithilfe von einer Linie verbunden. Der Ausgangskreis ist „World“. Bei dem Baumhierarchie gibt es keine weiteren Interaktionsmöglichkeiten. Das Baumdiagramm ist in Abbildung 3 zu finden.

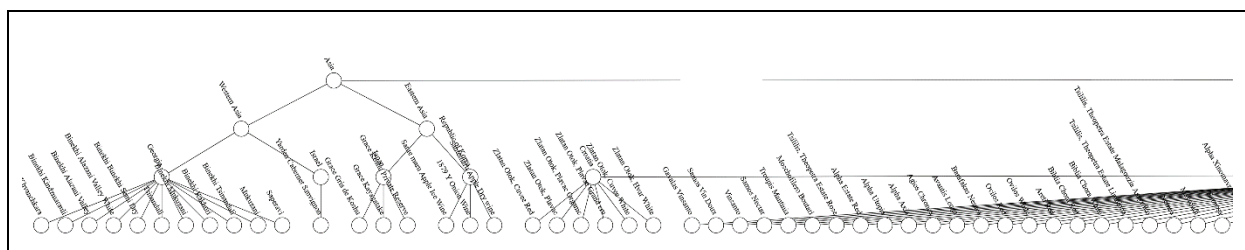


Abbildung 3: Auszug aus der Baumhierarchie

Die Anforderungen an das Baumdiagramm wurden nur teilweise erfüllt. Dabei sollte mithilfe dieser Darstellung die Frage nach den Produktionsmengen und Weinarten beantwortet werden. So stellt

diese Hierarchie zwar den einzelnen Ländern und deren Weine dar, jedoch gibt es keinen Aufschluss darüber wie viele diese sind. Weiterhin ist auch kein Abgleich der Produktionsmengen mit diesen Zahlen möglich, da diese nicht aufgeführt worden sind. Dementsprechend gibt es auch keine Gegenüberstellung dieser beiden Größen. Somit wurde nur die Anforderung teilweise erfüllt die Länder und ihre Weine darzustellen, die restlichen Anforderungen wurde hingegen nicht erfüllt.

Für diese Darstellung sollte nach den Anforderungen des Projektes hierbei ein Baum oder Graphen Technik verwendet werden. Beide Techniken eignen sich Zusammenhänge darstellen und somit für die Darstellung von Hierarchien. Da jedoch in dieser Darstellung um eine weitere einfache Hierarchie dargestellt werden sollte (die von Ländern und Weinamen) mussten keine größeren Zusammenhänge dargestellt werden, wodurch eine mehrfache Zuweisung auf einen Punkt in der Darstellung hätte erfolgen können. Dementsprechend wurde sich aufgrund der Übersichtlichkeit für ein Baumdiagramm entschieden, bei welchem es verschiedene Techniken gibt solche einfachen Hierarchischen Beziehungen übersichtlich darzustellen.

### 3.4 Interaktion

Bei den unterschiedlichen Darstellungen gibt es verschiedene Interaktionsmöglichkeiten, welche betätigt werden können. Die Hauptinteraktionsart ist dabei der Button. Mithilfe dieser können die bei den einzelnen Darstellungen die Eigenschaften welche aktuell angezeigt werden geändert werden. So ist es dem Nutzenden möglich die Darstellungen so anzupassen, wie diese es benötigen. Weiterhin besitzen diese die Möglichkeit mithilfe des Schwebens der Maus über die Kreise des Scatterplots den genauen Datensatz angezeigt bekommen. Damit können die Nutzenden die verschiedenen Datensätze identifizieren und so die besonderen Datensätze der entsprechenden Darstellung herausfinden. Dabei gibt es keine Interaktionsmöglichkeiten zwischen den Darstellungen, da diese ansonsten sich gegeben falls untereinander beeinflusst hätten, was nicht von den Nutzenden gewünscht sein könnte, da diese verschiedenen Eigenschaften ggf. genauer betrachten möchte wären er einen generellen über die anderen Eigenschaften erhalten möchte. Die Buttons sollten dabei jedem Nutzenden bekannt sein, weswegen auf diese Darstellungsform zurückgegriffen worden ist. Zusätzlich können die Nutzer die Daten des Scatterplots genauer untersuchen indem diese mithilfe der Maus über die Kreise dieser schweben. Auf diese Darstellungsform wurde zurückgegriffen, da diese am übersichtlichen ist und somit auch nicht die Erkennung von verschiedenen Trends oder Besonderheiten beeinflusst. Auf weitere Interaktionsmöglichkeiten aufgrund der Übersichtlichkeit nicht zurückgegriffen.

## 4. Implementierung

Das Projekt wurde auf der Grundlage der verschiedenen Übungen angefertigt, welche bereits vor dem Projekt fertiggestellt werden sollten. Dabei dienten insbesondere die Übung eins und Übung drei mit dem Scatterplot und der Interaktionsmöglichkeit als Grundlage für den Scatterplot. Bei den parallelen Koordinaten wurde die Übung sieben als Grundlage genutzt. Innerhalb der Baumhierarchie wurde die Übung zehn als Grundlage verwendet. Zusätzlich wurde für den Scatterplot und die parallelen Koordinaten die Übung acht mit ihrem CSV-Decoder als Grundlage verwendet. Dabei wurde bei beiden Darstellungen jeweils die Übung als Grundlage verwendet und die nun versucht mit dem CSV-Decoder der Übung acht in Verbindung zu bringen, um so eine Importierung des CSV-Datei zu ermöglichen. Nachdem dies erfolgreich umgesetzt worden ist, wurde für den Scatterplot noch die Buttons eingebaut um diesen entsprechend Interaktiv zu machen. Dafür waren weitere Anpassungen am CSV-Decoder und am Code nötig, welche auch anschließend bei der Entwicklung der parallelen Koordinaten einfließen könnte. Bei der Baumhierarchie wurde musste wenig angepasst werden, da bereits in der Übung ein JSON Decoder enthalten war. Hierbei war es nur wichtig die entsprechenden

Daten bereitzustellen. Zuletzt wurden die Darstellungen noch einmal in ihrer Darstellung überarbeitet. Die Entwicklung dieser Darstellungen erfolgte dabei im Wesentlichen nacheinander wobei sich einige Entwicklungen überlagerten, um die Zeit der Fehlersuche zu überbrücken. So wurde zuerst der Scatterplot entwickelt anschließend die parallelen Koordinaten und zuletzt die Baumhierarchie.

Der Quellcode ist dabei in die verschiedenen Parts von Elm des Models, View und Update aufgeteilt. Dabei enthalten diese einzelnen entsprechenden Funktionen welche entsprechend von diesen verwendet werden können. Oberhalb dieser Definitionen sind immer die Importe der verschiedenen Elm Packages zu finden. In allen drei verschiedenen Darstellungen wurde dabei auf die Packages „Html“ und „TypedSVG“ zugegriffen. Weiterhin entsprechende auf die Module je nach Art der Quelldatei also auf das CSV oder JSON Package. Weitere Packages wurden entsprechenden nach Bedarf in der Darstellung importiert. Die Daten werden innerhalb des Quellcodes innerhalb eines Records abgespeichert welcher innerhalb der Main zu finden ist. Auf diesen Record anschließend zugegriffen, wenn etwas zu aktualisieren gibt, um auf die entsprechenden Daten zuzugreifen.

Die Umsetzungen der Darstellungen waren die besonderen Herausforderungen den CSV-Decoder zu erweitern und mit den Übungen zu verbinden. Bei der Erweiterung des CSV-Decoders lag die Herausforderung darin mehr als nur 4 Spalten auf einmal zu verarbeiten da das Tupel von Elm nur maximal drei verschiedene Items akzeptieren. Dies wurde mithilfe eines neu geschriebenen Decoders umgangen. Bei der Verbindung des CSV-Decoders und der Übung gab es primär Probleme mit den Funktionen und den Datentypen, welche aber mit verschiedenen Anpassungen der Funktionen angepasst werden konnten. Die eigentliche Darstellung wurde jedoch fast so wie aus der Übung übernommen und hat somit nicht viel in Anspruch genommen. Selbiges gilt für die Baumhierarchie, welche fast ohne größere Anpassungen der Übung übernommen werden konnte.

## 5. Anwendungsfälle

Nachfolgend werden für die drei Visualisierungen jeweils ein Anwendungsfall vorgestellt. In diesem Anwendungsfall wird die entsprechende Darstellung verwendet, um verschiedene Informationen aus dieser zu gewinnen, welche für die verschiedenen Zielgruppen relevant sein könnte.

### 5.1 Anwendung Visualisierung Eins

In dem ersten Anwendungsfall wird der Preis mit dem Körper eines Weines verglichen. Dies erfolgt mithilfe eines Scatterplots und in Abbildung 4 zu erkennen.

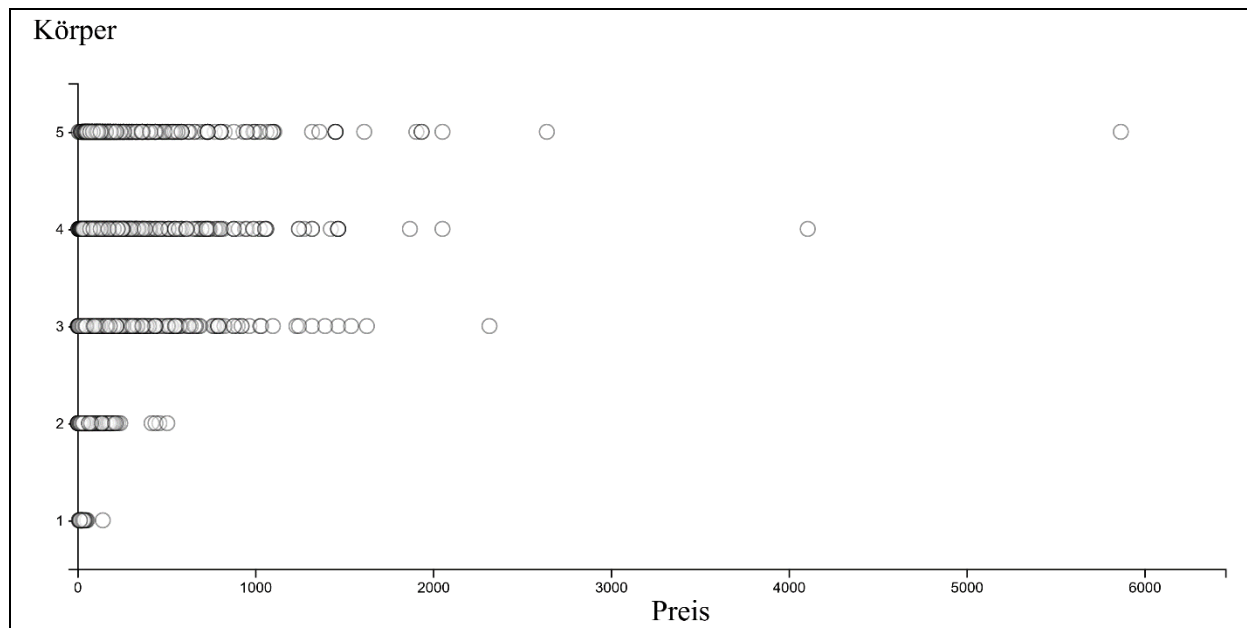


Abbildung 4: Scatterplot welcher die Eigenschaften Preis und Körper gegenüberstellt (Quelle: eigene Darstellung)

In Abbildung 4 ist zu erkennen das es bei den niedrigeren Preise Weine in allen Körpergrößen gibt. Jedoch ist klar zu erkenne das mit steigender Körpergröße der Preis auch immer weiter ansteigen kann. So ist besitzt der teuerste Wein, auch die größte Körpergröße. Weiterhin steigt die Verteilung der Weine auf der Preisachse mit steigender Körpergröße an, bis dieser sich bei einer Körpergröße von 3 bis 5 gleichbleibend ist. Weswegen davon auszugehen könnte das mit steigendem Preis immer eine gewisse Körpergröße einher gehen würde.

Diese Darstellung könnte so von einem Weininteressierten verwendet werden, um herauszufinden, wie er unter einem guten Preis-Leistungsverhältnis einen Wein mit einem großen Körper erhält. Hierbei sollte dieser Zielgruppe auffallen das er verschiedenen Körpergrößen auch zu einem geringen Preis bekommt. Jedoch je mehr die Zielgruppe beriet ist an Geld für den Wein auszugeben, desto wahrscheinlicher ist es das diese Zielgruppe einen Wein mit einer Größen Körpergröße erhält.

Als alternative zu diesen Diagrammen hätte es eine Zeitreihendiagramm in diesem Projekt gegeben. Da jedoch bei diesen Daten keinen zeitlichen Hintergrund gibt, ist es nicht möglich dieses Problem über Zeitreihendiagramme darzustellen.

## 5.2 Anwendung Visualisierung Zwei

Bei diesem Anwendungsfall werden mithilfe der parallelen Koordinaten die verschiedenen Eigenschaften des Körpers, Gerbstoffe, Süße und Säuregehalt des Weines gegenübergestellt. Diese Darstellung ist dabei in Abbildung 5 zu erkennen.



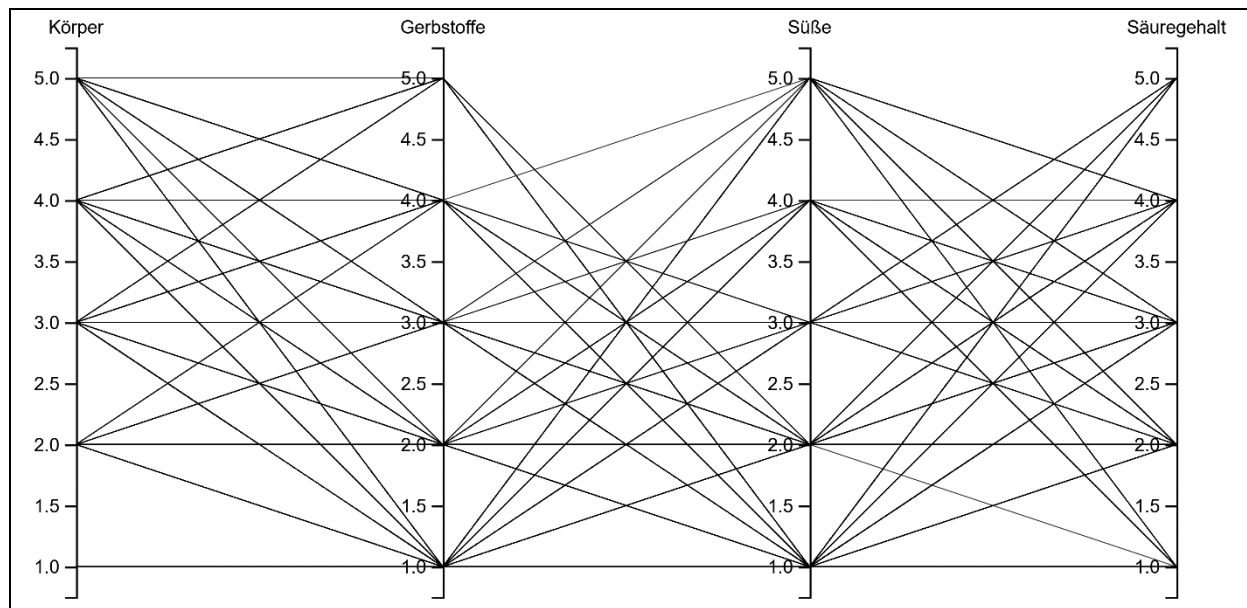


Abbildung 5: Gegenüberstellung der Eigenschaften Körper, Gerbstoffe, Süße und Säuregehalt mithilfe der Parallelen Koordinaten (Quelle: eigene Darstellung)

Dabei ist in Abbildung 5 zu erkennen, dass es keine besonderen Muster zwischen diesen verschiedenen Eigenschaften gibt. Dementsprechend gibt es viele verschiedene Möglichkeiten, die verschiedenen Eigenschaften in unterschiedlicher Intensität in einem Wein wiederzufinden. Auffälligkeiten an den Daten gibt es allerdings schon. So kann ein Wein, welcher die maximale Süße von fünf besitzt, nur einen Säuregehalt von vier besitzen. Umgekehrt ist es ähnlich: so kann ein Wein mit einem Säuregehalt von fünf maximal eine Süße von drei besitzen. Somit kann es keinen Wein geben, welcher auf beiden Skalen eine fünf besitzt. Dies hängt mit dem Gärungsprozess der Weintraube zusammen, wie bereits im Kapitel Anwendungshintergrund erklärt worden ist. Eine weitere Besonderheit besteht zwischen der Süße und den Gerbstoffen. So kann ein Wein, welcher einen hohen Gerbstoffgehalt von fünf besitzt, maximal eine Süße von zwei erreichen. Dies setzt sich aber nicht in den unteren Stufen fort. Dort ist es bereits ab einem Gerbstoffgehalt von vier möglich, die volle Süße von fünf zu erreichen. Jedoch gibt es dementsprechend auch wie bei der Säure und der Süße keinen Wein, welcher eine fünf in den Gerbstoffen und in der Süße hat. Umgekehrt verhält es sich bei der Beziehung zwischen den Körper und den Gerbstoffen. Hier besitzt ein Wein mit der Körpergröße eins auch nur einen Gerbstoffgehalt von eins. Jedoch besitzen die restlichen Körpergrößen alle verschiedene Gerbstoffgehalte.

Dieser Anwendungsfall wäre etwas für die Zielgruppe der Weininteressierten. Dabei könnten diese versuchen herauszufinden, welche Zusammenhänge es zwischen den verschiedenen Geschmäckern bei den Weinen gibt. Dabei sollte dieser Zielgruppe auffallen, dass es pauschal bis auf wenige Ausnahmen alles möglich zwischen bei der Kombination der verschiedenen Geschmäcker sein sollte. Dabei gibt es jedoch beispielsweise die Einschränkung, dass es keinen Wein geben kann, der eine Süße von fünf und einen Säuregehalt von fünf hat.

Bei dieser Visualisierung gäbe es verschiedene Alternativen zu dieser Diagrammart. Diese wären Scatterplots, Projektion und Selektion, K-Means und Datentinte. Da die K-Means und Datentinte Darstellungen waren, in denen etwas weggelassen worden ist, fallen diese heraus, da der Weininteressierte alle Weine miteinander vergleichen möchte, um entsprechende Trends oder Muster zu erkennen. Bei den Scatterplots und Projektion und Selektion würde es verschiedene Tabellen geben, in denen die verschiedenen Eigenschaften gegenübergestellt werden. Dies würde zwar zu einer besseren Analyse der Daten führen, da jede Eigenschaft direkt gegenübergestellt wird, was

aktuell bei den parallelen Koordinaten mit dem Säuregehalt und dem Körper nicht geschieht, jedoch würde dies gleichzeitig auch zu einem gewissen Informationsverlust einhergehen. Dieser Informationsverlust würde daher kommen, dass die verschiedenen Eigenschaften einzeln gegenübergestellt werden würden und somit Trends oder Muster unter den Muster nicht so einfach zu erkennen gewesen wären.

### 5.3 Anwendung Visualisierung Drei

- Anwendungsfall für Baumhierarchie
- Heraussuchen verschiedener Asiatischer Weine
- Georgien hat die meisten Weine
- Israel die wenigsten
- Japan und Süd-Korea gleich viele
- Spezifischer Anwendungsfall -> wo Muster da sind oder nicht was es zu was Besonderes macht
- Relevanz für die Zielgruppe
- Möglichkeit Umsetzung mit anderen Diagrammen

Innrehalb des letzten Anwendungsfall werden die verschiedenen Weine von Asien genauer analysiert. Ein Ausschnitt aus dem Baumdiagramm welche die Weine Asien darstellt ist in Abbildung 6 zu erkennen.

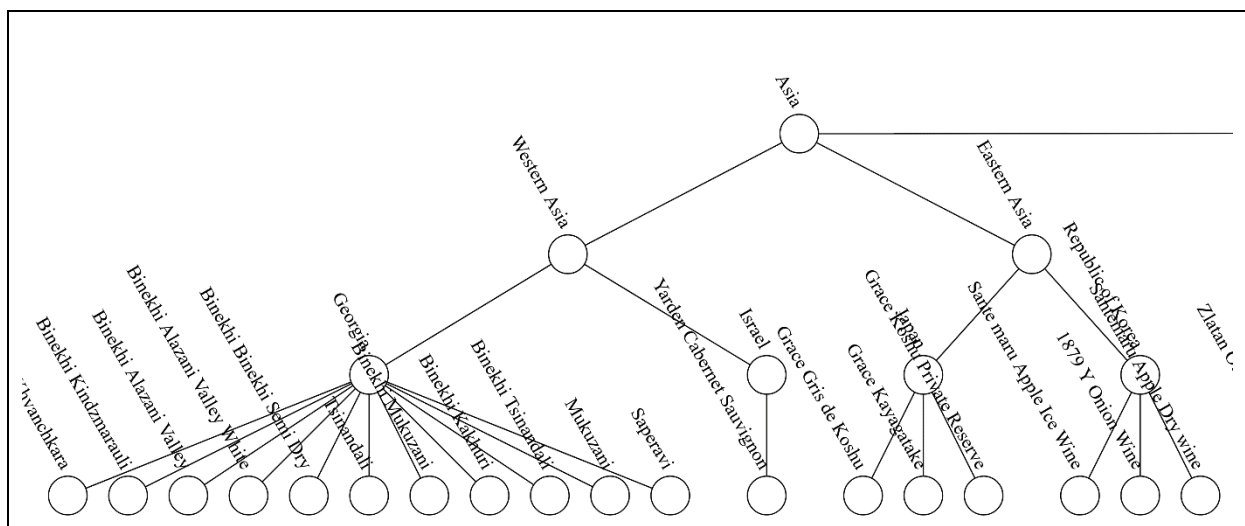


Abbildung 6: Ausschnitt mit Asiatischen Weinen aus der Baumhierarchie (Quelle: eigene Darstellung)

In Abbildung 6 sind die verschiedenen Weinamen zu erkennen, welcher der Datensatz an Asiatischen Weinen zu bieten hatte. Dabei ist auffällig, dass Georgien die meisten Weinarten mit elf Stück besitzt. Die niedrigste Anzahl besitzt Israel mit nur einer Weinart. Die beiden anderen Länder Japan und Südkorea sind mit einer Weinarten Anzahl von drei im Datensatz vorhanden. Auffällig ist dabei, dass es trotz der vielen Länder in Asien nur 4 Länder in diesem Datensatz vertreten sind, welche entsprechende Weinarten besitzen.

Diese Darstellung könnte von einem Weineinkäufer verwendet werden, welcher gerne auf dem heimischen Markt asiatische Weine verkaufen wollte. So könnte er mithilfe dieser Darstellung herausfinden, welche verschiedenen Weine es aus dem Asiatischen Raum gibt und aus welchem Land diese kommen, um so Transportkosten oder eine weitere Spezialisierung seines Geschäfts vorzunehmen.

Als Alternative für diese Darstellung gab es die Möglichkeit die verschiedenen Weinarten mithilfe der Graphen Technik darzustellen. Da es bei diesen Daten um eine einfache Hierarchiebeziehung handelt und somit nur eine Beziehung zwischen den einzelnen Kreisen dargestellt werden muss, ist diese Graphen Technik überflüssig.

## 6. Verwandte Arbeiten

Die hierbei verwandten Arbeiten beschäftigen sich mit der Visualisierung von Medizinischen Daten. Dies wird dabei bei beiden Daten unter anderem mit einem Scatterplot realisiert. Nachfolgend werden die beiden Arbeiten vorgestellt und Unterschiede zwischen dieser und der vorgestellten Arbeiten hervorgehoben.

Bei der ersten Arbeit handelt es sich um die Visualisierung von Patientendaten, welche an Multiplen Sklerose leiden. Ziel dieser Arbeit war es dabei die Darstellungen dabei mithilfe eines Scatterplots darzustellen und dafür sollte die Benutzeroberfläche so einfach wie möglich sein und Interaktionen der Nutzenden zulassen. Um dies zu realisieren wurden die Daten dabei bereits im vornherein mithilfe von verschiedenen Klassen eingeteilt. Dies sollte der Übersichtlichkeit der Daten weiterhelfen. Zusätzlich dazu sollte es in der Darstellung möglich sein verschiedene Statistische Größen wie beispielsweise den Median sich zu berechnen lassen, was entsprechend realisiert werden musste. Die Daten werden dabei mithilfe der Middleware „ColdFusion“ aufbereitet. Weiterhin werden die verschiedenen Daten mithilfe dieser Software dargestellt. Dabei wurden in dieser Arbeit ein Scatterplot, Balkendiagramm und ein Histogramm erstellt. Eine Interaktionsmöglichkeit gab es dabei ausschließlich mit dem Balkendiagramm. Dabei war es möglich über ein Dropdown Menü möglich sich alle Eigenschaften der Datenbank sich anzuzeigen lassen und eine entsprechende Auswahl zu treffen. Dies ist in der Darstellung sowohl für die X- als auch die Y-Achse möglich. Darüber hinaus ist es möglich nach anschließender Auswahl sich ein Histogramm zu erzeugen lassen. Bei dem Scatterplot als Ergebnis wurde darüber hinaus noch die Konfidenzintervalle und eine Regressionsgerade hinzugefügt. [14]

Die Gemeinsamkeiten zu dieser Arbeit liegen dabei in der Zielstellung, so sollte auch in dieser Arbeit eine möglichst eine Benutzeroberfläche erstellt, welche eine Interaktion mit den Nutzenden zulässt. Darüber hinaus wurde sich auch für eine Darstellung der Daten als Scatterplot entschieden. Jedoch sind die Interaktionsmöglichkeiten bei dieser Arbeit weitreichender als bei der hier vorliegenden. So werden ist es möglich mithilfe der Interaktionsmöglichkeit bei den Balkendiagramm sich ein Histogramm erstellen zu lassen. Weiterhin wurde sich für eine andere Interaktionsmöglichkeit entschieden. Anstatt der Buttons, welcher in dieser Arbeit verwendet wurden, wurde auf ein Dropdown Menü gesetzt, um den unter die Auswahl der verschiedenen Achsen zu überlassen. Zusätzlich wurden auch in beiden Projekten unterschiedliche Sprachen geschrieben. [14]

In der zweiten verwandten Arbeit geht es um die Visualisierung von verschiedenen Daten rund um Diabetespatienten. Dabei fokussierte sich diese Arbeit besonders auf die Analyse der Nutzbarkeit des Scatterplots. Dabei sollte untersucht werden in wie mit diesem Scatterplot umgegangen wird. Der Scatterplot besitzt dabei verschiedene Interaktionsmöglichkeiten. So ist es möglich die verschiedenen Achsen mithilfe von einem Dropdown Menü anzupassen und die Daten können nach dem Wünschen der Nutzenden angepasst werden. Darüber hinaus ist es möglich auf die verschiedenen Daten herein- und herauszoomen. Weiterhin ist es möglich Hilfe des Berührens des Punktes mit der Maus mehr Informationen zu der entsprechenden Untersuchung angezeigt zu bekommen. Da es sich bei den Daten um Daten handelt, welche über eine zeitliche Dimension besitzen ist es in dem diesem Scatterplot weiterhin möglich mithilfe der Bedienung eines Videorecorders zwischen den verschiedenen Zeitpunkten hin und her zu bewegen, und so

verschiedene Änderungen zu erkennen. Um dies zu erleichtern, gibt es weiterhin eine Funktion, um die Daten mithilfe einer Spur verfolgen zu können. [15]

Der große Unterschied zwischen dieser und der eben vorgestellten Arbeit ist, dass der Scatterplot um viele Funktionen erweitert wurde, um die Bedienung für die Nutzenden zu vereinfachen. So gibt es zwar in beiden die Möglichkeit die Achsen mithilfe von Buttons oder Dropdown Menüs zu verändern. Weiterhin ist es möglich in beiden genauere Informationen über den Datensatz zu erhalten, indem dieser mithilfe der Maus berührt wird. Jedoch gibt es durch die zeitliche Dimension viel mehr Möglichkeiten diese Daten anders in diesem Scatterplot zu realisieren und somit mehr Interaktionsmöglichkeiten einzubauen. Dies wurde beispielsweise durch die Art Videofernbedienung umgesetzt. Aber auch andere Funktionen fehlen in dieser Arbeit welche nicht Zeitlich abhängig sind, wie Beispiele das Zoomen auf eine andere Ebene, um die Daten beispielsweise besser erkennen zu können. Zusammenfassend ist zu sagen das diese Arbeit mehr Darstellungen als die vorgestellte Arbeit besitzt, jedoch die vorgestellte dabei viel mehr Interaktionsmöglichkeiten bietet, um die Nutzungserfahrung zu einfach und angenehm zu gestalten wie nur möglich. [15]

## 7. Zusammenfassung und Ausblick

Mithilfe dieser Projektarbeit ist es möglich diesen Weindatensatz, welcher die Grundlage für dieses Projekt bildet, einfach und mithilfe verschiedenster Darstellungen zu analysieren. Diese Darstellungen besitzen dabei einfache Benutzeroberfläche, um mit diesen zu interagieren und es auf die entsprechenden Wünsche der Nutzenden anzupassen. Dadurch ist eine einfache individuell angepasste Analyse dieser verschiedenen Daten möglich. Somit ist es den verschiedenen Nutzenden möglich verschieden Trends, Muster oder auch verschieden Zusammenhänge mit diesen Daten zu finden, nachzuvollziehen oder zu reproduzieren.

Auch die Zielgruppen, für welche diese Arbeit geschrieben worden ist, erhalten durch diese Visualisierungen verschiedene Vorteile. So ist es dem Weininteressierten nun möglich mithilfe des Scatterplots verschieden Weine zu finden, welche seinem Geschmack entsprechen oder auch komplett neue Sorten entdecken. Darüber hinaus ist es beispielsweise auch möglich mithilfe der parallelen Koordinaten und der Baumhierarchie verschieden Zusammenhänge besser verstehen zu können oder ganz neue Erkenntnisse zu erhalten. Auch die Zielgruppe der Weinexperten erhält verschieden Vorteile, so können diese mithilfe der verschiedenen Visualisierungen neue Sorten entdecken oder auch ihre aktuellen Sorten besser einschätzen, um so ein noch besseres Wissen an Weinen anzusammeln. Weiterhin sollte es möglich sein mithilfe dieser verschiedenen Darstellungen ihr Weinwissen besser darzustellen und nachzuvollziehen. Innerhalb der Gruppe der Weinverkäufer ist es möglich neue Weine für die Kunden zu finden und diesen damit eine noch bessere Beratung zu geben. Weiterhin ist mithilfe des Baumdiagrammen auch eine Spezialisierung auf ein gewisses Land oder Kontinent denkbar, um eine Spezialisierung des eigenen Ladens zu erreichen.

Eine Erweiterung dieses Projektes bei den Visualisierungen wäre vor allem bei der Interaktivität der einzelnen Darstellungen denkbar. So wäre es möglich eine Filterfunktion für die Daten einzubauen, oder auch eine Hervorhebung der Daten welche aktuell verwendet ausgewählt sind, welche sich über alle Darstellungen ziehen würde. Aber auch könnten weitere Darstellungen hinzugefügt werden, um eine Analyse der verschiedenen Daten noch einfacher zu gestalten. So wäre beispielsweise ein Balkendiagramm denkbar in dem die einzelnen Eigenschaften eines Weines aufgeführt werden, um so eine Übersicht über diese zu erhalten. Weiterhin sollte bei steigender Anzahl der Datenmengen über eine Kategorisierung der Daten nachgedacht werden um so die Übersichtlichkeit innerhalb der Darstellungen nicht zu verlieren. Bei dem Baumhierarchie wäre es denkbar die einzelnen Bäume erst per Mausklick zu öffnen um so eine noch bessere Sichtbarkeit zu beisetzen. Bei den parallelen

Koordinaten sollten die Daten hervorgehen werden, über welche der Mauszeiger aktuell liegt, umso noch mehr Einzelinformationen dem Nutzenden bereitzustellen.

Bei der Datenebene ist es wichtig der Datensatz immer weiter erweitert wird, um final alle Weine der Welt abzudecken. Dabei ist es jedoch nicht nur wichtig alle Weine enthalten sind, jedoch auch alle Eigenschaften über diesen Wein bekannt sind, um so sicherzustellen das diese Daten auch untereinander verglichen werden können. Diese Vollständigkeit der Daten ist beim aktuellen Stand noch sehr dürftig, könnte aber durch eine vervollständigung der einzelne Datensätze sich schnell verbessern. Darüber hinaus wäre es denkbar, wenn dieser Datensatz für Weinexperten oder Weineinkäufer/Verkäufer mehr als nur eine Indikation geben soll, die Daten weiter zu konkretisieren. So sollten die 1 bis 5 Skala, welche es beispielsweise bei der Süße, Säuregehalt und Körper gibt, durch eine bessere Messskala ersetzt werden. Diese könnte den entsprechenden Zielgruppen mehr Informationen übermitteln als nur ein wert zwischen ein und fünf.

## Literaturverzeichnis

- [1] P. Roca, *State of the Vitivinicultural World in 2020*, o. O., **20.04.2021**.
- [2] M. Yi, *A Complete Guide to Scatter Plots*, <https://chartio.com/learn/charts/what-is-a-scatter-plot/>, **2019**.
- [3] S. Few, *Multivariate Analysis Using Parallel Coordinates*, [http://www.perceptualedge.com/articles/b-eye/parallel\\_coordinates.pdf](http://www.perceptualedge.com/articles/b-eye/parallel_coordinates.pdf), **2006**.
- [4] American Society for Quality, *What is a Tree Diagram? Systemic or Hierarchy Analysis | ASQ*, <https://asq.org/quality-resources/tree-diagram>, **o. J.**
- [5] L. Beilmann, *Alkoholgehalt in Wein - Das solltest du unbedingt Wissen!*, <https://wein-fuer-laien.de/weinwissen/alkoholgehalt-im-wein/>, **o. J.**
- [6] M. Teufel, *Die richtige Trinktemperatur für Wein - warum wichtig?*, <https://swisscave.com/de/swisscave-blog/post/die-richtige-trinktemperatur-fur-wein-warum-wichtig>, **2021**.
- [7] Brogsitter Weinversand, *Süße*, <https://www.brogsitter.de/weinlexikon/suesse/#>, **o. J.**
- [8] Weinkenner GmbH, *Die Säure | Weinkenner.de*, <https://www.weinkenner.de/die-saeure/>, **2011**.
- [9] Brogsitter Weinversand, *Körper*, <https://www.brogsitter.de/weinlexikon/koerper/>, **o. J.**
- [10] Brogsitter Weinversand, *Gerbstoffe*, <https://www.brogsitter.de/weinlexikon/gerbstoffe/>, **o. J.**
- [11] Vineyard99, *Wein & Wissen: Weinjahrgang – Geheimnis der Weinalterung*, <https://www.vineyard99.de/weinjahrgang-und-weinalterung/>, **2020**.
- [12] dev7halo, *Wine Information*, <https://www.kaggle.com/dev7halo/wine-information>, **2021**.
- [13] Curran Kelleher, *World Countries Hierarchy*, <https://gist.github.com/curran/1dd7ab046a4ed32380b21e81a38447aa/>, **2015**.
- [14] M. Petrova, *Web-basierte dynamische Visualisierung klinischer Daten*, <https://www.nm.informatik.uni-muenchen.de/common/pub/Fopras/petr02/PDF-Version/petr02.pdf>, **2002**.
- [15] U. Fels, *Usability-Analyse des Programms Animated Scatter Plot*, <https://repositum.tuwien.at/bitstream/20.500.12708/4082/2/Fels%20Ulrich%20-%202015%20-%20Usability%20Analyse%20des%20Programms%20Animated%20Scatter%20Plot.pdf>, **o. J.**

## Anhang