# 1 Analysis Motivation

An alcohol company that owns a chain of stores across Russia recently had ran a successful wine promotional in Saint Petersburg. Due to a limited budget, the company will not be able to run a promotion in all regions. Hence, the marketing team has tasked me to identify ten other regions to run a successful wine promotion. To do this, we will implement a clustering model that will return regions with similar purchasing habits as St Petersburg. These regions will then be analysed further based on their recent wine sales to pinpoint the best ten regions for a future wine promotion.

# 2 Exploratory Data Analysis

The dataset used in this project can be found at the DataCamp careerhub-data repository. Specifically, we will be focusing on the "Alcohol Consumption in Russia" dataset. This dataset consists of the following seven columns.

| Description | Type | Non-Null Count | Data Type |
|---|---|---|---|
| Year | numerical | 1615 | int64 |
| Region | categorical | 1615 | object |
| Wine | numerical | 1552 | float64 |
| Beer | numerical | 1557 | float64 |
| Vodka | numerical | 1554 | float64 |
| Champagne | numerical | 1552 | float64 |
| Brandy | numerical | 1549 | float64 |

Table 1: Type and Missing Data

In the reaming parts of this section we will go through the motions of exploring data. This will give us a deeper understanding of the data and help uncover any issues within it. So let us get started. Firstly, the shape of the data is $(1615 \times 7)$, and thus very small, as opposed to most data sets. Secondly, we will run some basic statistics on the numerical data using the describe method from pandas.

|  | year | wine | beer | vodka | champagne |
|---|---|---|---|---|---|
| count | 1615.000000 | 1552.000000 | 1557.000000 | 1554.000000 | 1552.000000 |
| mean | 2007.000000 | 5.628144 | 51.206148 | 11.818694 | 1.313177 |
| std | 5.478922 | 2.813208 | 25.372821 | 5.128806 | 0.797956 |
| min | 1998.000000 | 0.100000 | 0.400000 | 0.050000 | 0.100000 |
| 25% | 2002.00000 | 3.575000 | 32.400000 | 8.300000 | 0.800000 |
| 50% | 2007.000000 | 5.400000 | 49.970000 | 11.500000 | 1.200000 |
| 75% | 2012.000000 | 7.377500 | 67.400000 | 15.000000 | 1.665000 |
| max | 2016.000000 | 18.100000 | 207.300000 | 40.600000 | 5.560000 |

Table 2: Initial Statistics with Describe

From Table 1, we can already see that there is some missing data that will need to be dealt with. In fact,

the only complete columns are year and region. Before treating the missing data, let's visualize the data using histograms as shown in Figure 1. It appears that the data is mostly right skewed. Where champagne and brandy have longer tails. On the other hand wine, beer and vodka are closer to normal.



Figure 1: Histograms

# 3 Treating Missing Data

Now we will handle the missing data using one of my personal favorite modules in python, missingo. Using this module we can visualize the distribution of missing data. Note in the following figures, the white "chuncks" refer to missing data. In Figure 2, the missing data is scattered throughout the entire data set. Whereas, in Figure 3 the data is arranged by region, and it is clear the missing data occurs in four different regions. After some further analysis the regions of interests are: Chechen Republic, Republic of Crimea, Sevastopol and Republic of Ingushetia.

Figure 2: Location of Missing Data



Figure 3: Location of Missing Data Grouped by Region

In order to handle the missing data, any fully incomplete rows will be dropped. We can retain other rows of missing data by imputing the mean of its corresponding region. With this in place, we now have a complete data set. However, we must treat outliers. It is standard practice to treat outliers but more importantly they can pose a massive threat in terms of the quality of our future clustering model.

# 4  Treating Outliers

Treating outliers can be tricky. We will first use box plots to visualize any data which reaches far beyond the mean. These points will be the main suspects for the outliers. Secondly we will evaluate the z-scores as

a double checking method. Lastly we will pinpoint where the outliers are located within the data and deal with them accordingly.



Figure 4: Box plot of each alcohol type

Immediately, from Figure 4 the top three data points for beer and the top data point for vodka are suspected outliers. Since the other three alcohol types are difficult to interpret due to scaling. We will create a separate plot for these categories.



Figure 5: Box plot of wine, champagne, and brandy

From Figure 5 the suspected outliers are the top five data points for wine, top five data points for champagne, and the top data point for brandy. Next we shall check the z-scores for the top 5 data points. Recall the z-score of a data point signifies how many standard deviations away from the mean the data point is. The z-score can be positive or negative number, indicating the data point is above or below the mean respectively. Comparing the list of suspected outliers to Table 3, our suspicions can be put to rest. The data points discussed in the prior paragraph are far beyond the mean, and should be imputed in some way. Each of the data points mentioned will be filled with the mean alcohol sales during the year where the outlier occurred.

| | wine | beer | vodka | champagne | brandy |
|---|---|---|---|---|---|
| 0 | 4.442437 | 6.187809 | 5.641310 | 5.327700 | 4.431710 |
| 1 | 4.228541 | 5.981261 | 3.933907 | 4.800513 | 4.006786 |
| 2 | 4.014644 | 5.516529 | 3.933907 | 4.750304 | 3.931799 |
| 3 | 3.800748 | 3.633765 | 3.894657 | 4.750304 | 3.931799 |
| 4 | 3.586852 | 3.093563 | 3.757280 | 4.624784 | 3.931799 |
| 5 | 3.337306 | 3.057814 | 3.502150 | 4.624784 | 3.931799 |
| 6 | 3.266008 | 2.995453 | 3.462900 | 4.499263 | 3.681844 |
| 7 | 3.266008 | 2.970429 | 3.462900 | 4.499263 | 3.681844 |
| 8 | 3.230358 | 2.954540 | 3.266647 | 4.373742 | 3.431889 |
| 9 | 3.159059 | 2.930708 | 3.247021 | 4.373742 | 3.431889 |

Table 3: Top 10 z-scores (descending)

# 5 K-Means Clustering

Now that the data has been processed we can work toward picking our model. As the section title suggests we will be utilizing the K-means clustering algorithm. This is an unsupervised machine learning algorithm which aims at clustering regions into similar groups. In this case, "similar" refers to regions which share analogous historical wine sales. Since we know that St Petersburg had a successful wine promotion, we can check its cluster for other potential regions for a future wine promotion. However, we have a limited budget, so we may need to dig deeper beyond just selecting ten random regions in St Petersburg's cluster.

## 5.1 Prepping the Data

In order to use the K-Means algorithm we need to have data with low variance, and similar means. A quick fix for this is to normalize our wine sales data. That is to do the following:

1. Create a pivot table with regions as the row headers and years as the column headers

2. Fill the pivot table with wine sales

3. Divide each year column by the max wine sale for that year.

By doing so, all the wine sales data is within the range [0,1], and is now fitted for the K-Means algorithm. Moreover, normalized data is great for heat maps. In Figure 6 we can see the overall trend of wine sales for each region. Notice, the darker the cell, the closer to the max it is.

Wine Sales Rate

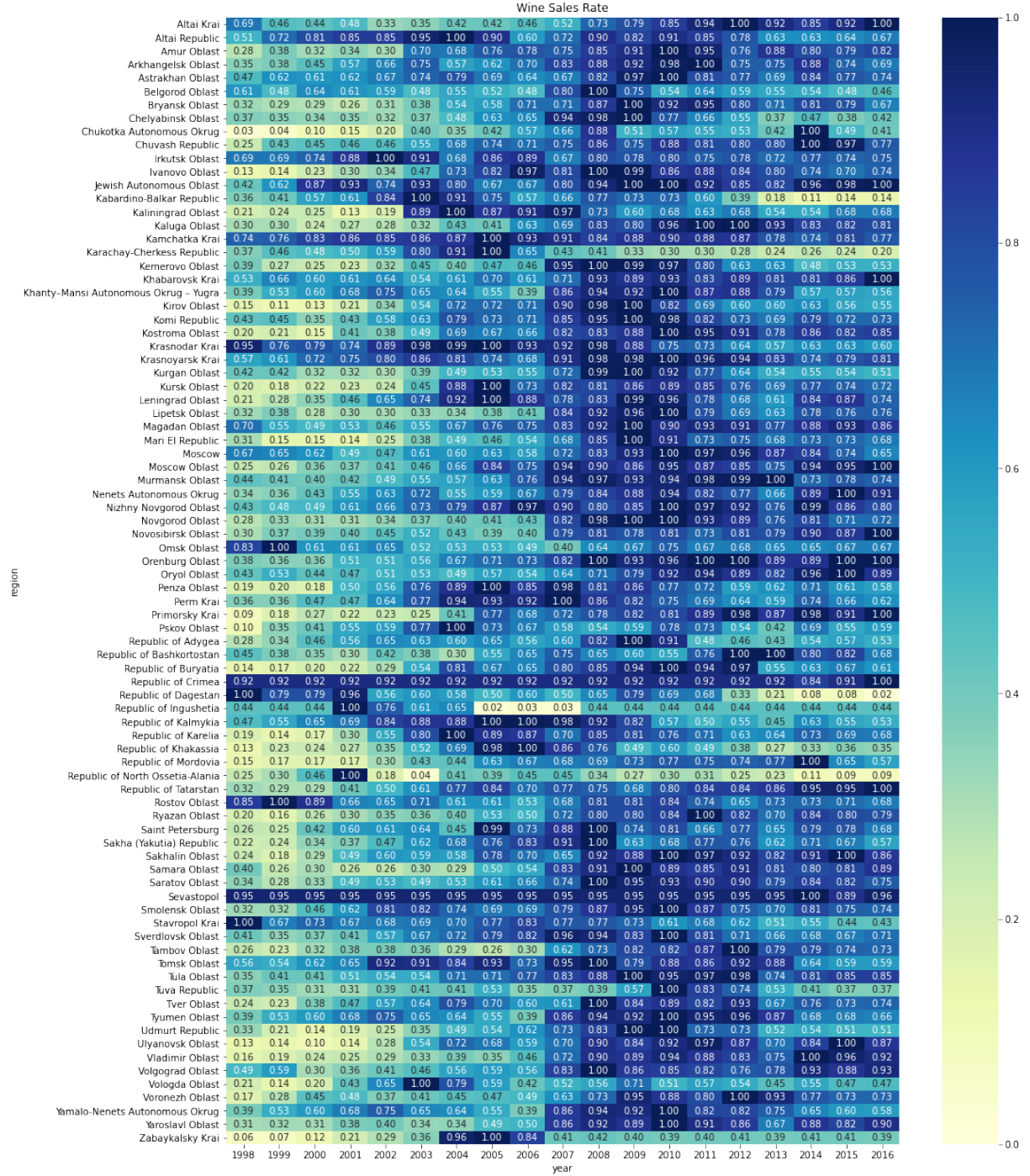| region | 1998 | 1999 | 2000 | 2001 | 2002 | 2003 | 2004 | 2005 | 2006 | 2007 | 2008 | 2009 | 2010 | 2011 | 2012 | 2013 | 2014 | 2015 | 2016 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Altai Krai | 0.69 | 0.46 | 0.44 | 0.48 | 0.33 | 0.35 | 0.42 | 0.42 | 0.46 | 0.52 | 0.73 | 0.79 | 0.85 | 0.94 | 1.00 | 0.92 | 0.85 | 0.92 | 1.00 |
| Altai Republic | 0.51 | 0.72 | 0.81 | 0.85 | 0.85 | 0.95 | 1.00 | 0.90 | 0.60 | 0.72 | 0.90 | 0.82 | 0.91 | 0.85 | 0.78 | 0.63 | 0.63 | 0.64 | 0.67 |
| Amur Oblast | 0.28 | 0.38 | 0.32 | 0.34 | 0.30 | 0.70 | 0.68 | 0.76 | 0.78 | 0.75 | 0.85 | 0.91 | 1.00 | 0.95 | 0.76 | 0.88 | 0.80 | 0.79 | 0.82 |
| Arkhangelsk Oblast | 0.35 | 0.38 | 0.45 | 0.57 | 0.66 | 0.75 | 0.57 | 0.62 | 0.70 | 0.83 | 0.88 | 0.92 | 0.98 | 1.00 | 0.75 | 0.75 | 0.88 | 0.74 | 0.69 |
| Astrakhan Oblast | 0.47 | 0.62 | 0.61 | 0.62 | 0.67 | 0.74 | 0.79 | 0.69 | 0.64 | 0.67 | 0.82 | 0.97 | 1.00 | 0.81 | 0.77 | 0.69 | 0.84 | 0.77 | 0.74 |
| Belgorod Oblast | 0.61 | 0.48 | 0.64 | 0.61 | 0.59 | 0.48 | 0.55 | 0.52 | 0.48 | 0.80 | 1.00 | 0.75 | 0.54 | 0.64 | 0.59 | 0.55 | 0.54 | 0.48 | 0.46 |
| Bryansk Oblast | 0.32 | 0.29 | 0.29 | 0.26 | 0.31 | 0.38 | 0.54 | 0.58 | 0.71 | 0.71 | 0.87 | 1.00 | 0.92 | 0.95 | 0.80 | 0.71 | 0.81 | 0.79 | 0.67 |
| Chelyabinsk Oblast | 0.37 | 0.35 | 0.34 | 0.35 | 0.32 | 0.37 | 0.48 | 0.63 | 0.65 | 0.94 | 0.98 | 1.00 | 0.77 | 0.66 | 0.55 | 0.37 | 0.47 | 0.38 | 0.42 |
| Chukotka Autonomous Okrug | 0.03 | 0.04 | 0.10 | 0.15 | 0.20 | 0.40 | 0.35 | 0.42 | 0.57 | 0.66 | 0.88 | 0.51 | 0.57 | 0.55 | 0.53 | 0.42 | 1.00 | 0.49 | 0.41 |
| Chuvash Republic | 0.25 | 0.43 | 0.45 | 0.46 | 0.46 | 0.55 | 0.68 | 0.74 | 0.71 | 0.75 | 0.86 | 0.75 | 0.88 | 0.81 | 0.80 | 0.80 | 1.00 | 0.97 | 0.77 |
| Irkutsk Oblast | 0.69 | 0.69 | 0.74 | 0.88 | 1.00 | 0.91 | 0.68 | 0.86 | 0.89 | 0.67 | 0.80 | 0.78 | 0.80 | 0.75 | 0.78 | 0.72 | 0.77 | 0.74 | 0.75 |
| Ivanovo Oblast | 0.13 | 0.14 | 0.23 | 0.30 | 0.34 | 0.47 | 0.73 | 0.82 | 0.97 | 0.81 | 1.00 | 0.99 | 0.86 | 0.88 | 0.84 | 0.80 | 0.74 | 0.70 | 0.74 |
| Jewish Autonomous Oblast | 0.42 | 0.62 | 0.87 | 0.93 | 0.74 | 0.93 | 0.80 | 0.67 | 0.67 | 0.80 | 0.94 | 1.00 | 1.00 | 0.92 | 0.85 | 0.82 | 0.96 | 0.98 | 1.00 |
| Kabardino-Balkar Republic | 0.36 | 0.41 | 0.57 | 0.61 | 0.84 | 1.00 | 0.91 | 0.75 | 0.57 | 0.66 | 0.77 | 0.73 | 0.73 | 0.60 | 0.39 | 0.18 | 0.11 | 0.14 | 0.14 |
| Kaliningrad Oblast | 0.21 | 0.24 | 0.25 | 0.13 | 0.19 | 0.89 | 1.00 | 0.87 | 0.91 | 0.97 | 0.73 | 0.60 | 0.68 | 0.63 | 0.68 | 0.54 | 0.54 | 0.68 | 0.68 |
| Kaluga Oblast | 0.30 | 0.30 | 0.24 | 0.27 | 0.28 | 0.32 | 0.43 | 0.41 | 0.63 | 0.69 | 0.83 | 0.80 | 0.96 | 1.00 | 1.00 | 0.93 | 0.83 | 0.82 | 0.81 |
| Kamchatka Krai | 0.74 | 0.76 | 0.83 | 0.86 | 0.85 | 0.86 | 0.87 | 1.00 | 0.93 | 0.91 | 0.84 | 0.88 | 0.90 | 0.88 | 0.87 | 0.78 | 0.74 | 0.81 | 0.77 |
| Karachay-Cherkess Republic | 0.37 | 0.46 | 0.48 | 0.50 | 0.59 | 0.80 | 0.91 | 1.00 | 0.65 | 0.43 | 0.41 | 0.33 | 0.30 | 0.30 | 0.28 | 0.24 | 0.26 | 0.24 | 0.20 |
| Kemerovo Oblast | 0.39 | 0.27 | 0.25 | 0.23 | 0.32 | 0.45 | 0.40 | 0.47 | 0.46 | 0.95 | 1.00 | 0.99 | 0.97 | 0.80 | 0.63 | 0.63 | 0.48 | 0.53 | 0.53 |
| Khabarovsk Krai | 0.53 | 0.66 | 0.60 | 0.61 | 0.64 | 0.54 | 0.61 | 0.70 | 0.61 | 0.71 | 0.93 | 0.89 | 0.93 | 0.83 | 0.89 | 0.81 | 0.81 | 0.86 | 1.00 |
| Khanty-Mansi Autonomous Okrug - Yugra | 0.39 | 0.53 | 0.60 | 0.68 | 0.75 | 0.65 | 0.64 | 0.55 | 0.39 | 0.86 | 0.94 | 0.92 | 1.00 | 0.87 | 0.88 | 0.79 | 0.57 | 0.57 | 0.56 |
| Kirov Oblast | 0.15 | 0.11 | 0.13 | 0.21 | 0.34 | 0.54 | 0.72 | 0.72 | 0.71 | 0.90 | 0.97 | 1.00 | 0.82 | 0.69 | 0.60 | 0.60 | 0.63 | 0.56 | 0.55 |
| Komi Republic | 0.43 | 0.45 | 0.35 | 0.43 | 0.58 | 0.63 | 0.79 | 0.73 | 0.71 | 0.85 | 0.95 | 1.00 | 0.98 | 0.82 | 0.73 | 0.69 | 0.79 | 0.72 | 0.73 |
| Kostroma Oblast | 0.20 | 0.21 | 0.15 | 0.41 | 0.38 | 0.49 | 0.69 | 0.67 | 0.66 | 0.82 | 0.83 | 0.88 | 1.00 | 0.95 | 0.91 | 0.78 | 0.86 | 0.82 | 0.85 |
| Krasnodar Krai | 0.95 | 0.76 | 0.79 | 0.74 | 0.89 | 0.98 | 0.99 | 1.00 | 0.93 | 0.92 | 0.98 | 0.88 | 0.75 | 0.73 | 0.64 | 0.57 | 0.63 | 0.63 | 0.60 |
| Krasnoyarsk Krai | 0.57 | 0.61 | 0.72 | 0.75 | 0.80 | 0.86 | 0.81 | 0.74 | 0.68 | 0.91 | 0.98 | 0.98 | 1.00 | 0.96 | 0.94 | 0.83 | 0.74 | 0.79 | 0.81 |
| Kurgan Oblast | 0.42 | 0.42 | 0.32 | 0.32 | 0.30 | 0.39 | 0.49 | 0.53 | 0.55 | 0.72 | 0.99 | 1.00 | 0.92 | 0.77 | 0.64 | 0.54 | 0.55 | 0.54 | 0.51 |
| Kursk Oblast | 0.20 | 0.18 | 0.22 | 0.23 | 0.24 | 0.45 | 0.88 | 1.00 | 0.76 | 0.73 | 0.82 | 0.81 | 0.86 | 0.89 | 0.85 | 0.76 | 0.69 | 0.77 | 0.72 |
| Leningrad Oblast | 0.21 | 0.28 | 0.35 | 0.46 | 0.65 | 0.74 | 0.92 | 1.00 | 0.88 | 0.78 | 0.83 | 0.99 | 0.96 | 0.78 | 0.68 | 0.61 | 0.84 | 0.87 | 0.74 |
| Lipetsk Oblast | 0.32 | 0.38 | 0.28 | 0.30 | 0.30 | 0.33 | 0.34 | 0.38 | 0.41 | 0.84 | 0.92 | 0.96 | 1.00 | 0.79 | 0.69 | 0.63 | 0.78 | 0.76 | 0.76 |
| Magadan Oblast | 0.70 | 0.55 | 0.49 | 0.53 | 0.46 | 0.55 | 0.67 | 0.76 | 0.75 | 0.83 | 0.92 | 1.00 | 0.90 | 0.93 | 0.91 | 0.77 | 0.88 | 0.93 | 0.86 |
| Mari El Republic | 0.31 | 0.15 | 0.15 | 0.14 | 0.25 | 0.38 | 0.49 | 0.46 | 0.54 | 0.68 | 0.85 | 1.00 | 0.91 | 0.73 | 0.75 | 0.68 | 0.73 | 0.73 | 0.68 |
| Moscow | 0.67 | 0.65 | 0.62 | 0.49 | 0.47 | 0.61 | 0.60 | 0.63 | 0.58 | 0.72 | 0.83 | 0.93 | 1.00 | 0.97 | 0.96 | 0.87 | 0.84 | 0.74 | 0.65 |
| Moscow Oblast | 0.25 | 0.26 | 0.36 | 0.37 | 0.41 | 0.46 | 0.66 | 0.84 | 0.75 | 0.94 | 0.90 | 0.86 | 0.95 | 0.87 | 0.85 | 0.75 | 0.94 | 0.95 | 1.00 |
| Murmansk Oblast | 0.44 | 0.41 | 0.40 | 0.42 | 0.49 | 0.55 | 0.57 | 0.63 | 0.76 | 0.94 | 0.97 | 0.93 | 0.94 | 0.98 | 0.99 | 1.00 | 0.73 | 0.78 | 0.74 |
| Nenets Autonomous Okrug | 0.34 | 0.36 | 0.43 | 0.55 | 0.63 | 0.72 | 0.55 | 0.59 | 0.67 | 0.79 | 0.84 | 0.88 | 0.94 | 0.82 | 0.77 | 0.66 | 0.89 | 1.00 | 0.91 |
| Nizhny Novgorod Oblast | 0.43 | 0.48 | 0.49 | 0.61 | 0.66 | 0.73 | 0.79 | 0.87 | 0.97 | 0.90 | 0.80 | 0.85 | 1.00 | 0.97 | 0.92 | 0.76 | 0.99 | 0.86 | 0.80 |
| Novgorod Oblast | 0.28 | 0.33 | 0.31 | 0.31 | 0.34 | 0.37 | 0.40 | 0.41 | 0.43 | 0.82 | 0.98 | 1.00 | 1.00 | 0.93 | 0.89 | 0.76 | 0.81 | 0.71 | 0.72 |
| Novosibirsk Oblast | 0.30 | 0.37 | 0.39 | 0.40 | 0.45 | 0.52 | 0.43 | 0.39 | 0.40 | 0.79 | 0.81 | 0.78 | 0.81 | 0.73 | 0.81 | 0.79 | 0.90 | 0.87 | 1.00 |
| Omsk Oblast | 0.83 | 1.00 | 0.61 | 0.61 | 0.65 | 0.52 | 0.53 | 0.53 | 0.49 | 0.40 | 0.64 | 0.67 | 0.75 | 0.67 | 0.68 | 0.65 | 0.65 | 0.67 | 0.67 |
| Orenburg Oblast | 0.38 | 0.36 | 0.36 | 0.51 | 0.51 | 0.56 | 0.67 | 0.71 | 0.73 | 0.82 | 1.00 | 0.93 | 0.96 | 1.00 | 1.00 | 0.89 | 0.89 | 1.00 | 1.00 |
| Oryol Oblast | 0.43 | 0.53 | 0.44 | 0.47 | 0.51 | 0.53 | 0.49 | 0.57 | 0.54 | 0.64 | 0.71 | 0.79 | 0.92 | 0.94 | 0.89 | 0.82 | 0.96 | 1.00 | 0.89 |
| Penza Oblast | 0.19 | 0.20 | 0.18 | 0.50 | 0.56 | 0.76 | 0.89 | 1.00 | 0.85 | 0.98 | 0.81 | 0.86 | 0.77 | 0.72 | 0.59 | 0.62 | 0.71 | 0.61 | 0.58 |
| Perm Krai | 0.36 | 0.36 | 0.47 | 0.47 | 0.64 | 0.77 | 0.94 | 0.93 | 0.92 | 1.00 | 0.86 | 0.82 | 0.75 | 0.69 | 0.64 | 0.59 | 0.74 | 0.66 | 0.62 |
| Primorsky Krai | 0.09 | 0.18 | 0.27 | 0.22 | 0.23 | 0.25 | 0.41 | 0.77 | 0.68 | 0.72 | 0.78 | 0.82 | 0.81 | 0.89 | 0.98 | 0.87 | 0.98 | 0.91 | 1.00 |
| Pskov Oblast | 0.10 | 0.35 | 0.41 | 0.55 | 0.59 | 0.77 | 1.00 | 0.73 | 0.67 | 0.58 | 0.54 | 0.59 | 0.78 | 0.73 | 0.54 | 0.42 | 0.69 | 0.55 | 0.59 |
| Republic of Adygea | 0.28 | 0.34 | 0.46 | 0.56 | 0.65 | 0.63 | 0.60 | 0.65 | 0.56 | 0.60 | 0.82 | 1.00 | 0.91 | 0.48 | 0.46 | 0.43 | 0.54 | 0.57 | 0.53 |
| Republic of Bashkortostan | 0.45 | 0.38 | 0.35 | 0.30 | 0.42 | 0.38 | 0.30 | 0.55 | 0.65 | 0.75 | 0.65 | 0.60 | 0.55 | 0.76 | 1.00 | 1.00 | 0.80 | 0.82 | 0.68 |
| Republic of Buryatia | 0.14 | 0.17 | 0.20 | 0.22 | 0.29 | 0.54 | 0.81 | 0.67 | 0.65 | 0.80 | 0.85 | 0.94 | 1.00 | 0.94 | 0.97 | 0.55 | 0.63 | 0.67 | 0.65 |
| Republic of Crimea | 0.92 | 0.92 | 0.92 | 0.92 | 0.92 | 0.92 | 0.92 | 0.92 | 0.92 | 0.92 | 0.92 | 0.92 | 0.92 | 0.92 | 0.92 | 0.92 | 0.84 | 0.91 | 1.00 |
| Republic of Dagestan | 1.00 | 0.79 | 0.79 | 0.96 | 0.56 | 0.60 | 0.58 | 0.50 | 0.60 | 0.50 | 0.65 | 0.79 | 0.69 | 0.68 | 0.33 | 0.21 | 0.08 | 0.08 | 0.02 |
| Republic of Ingushetia | 0.44 | 0.44 | 0.44 | 1.00 | 0.76 | 0.61 | 0.65 | 0.02 | 0.03 | 0.03 | 0.44 | 0.44 | 0.44 | 0.44 | 0.44 | 0.44 | 0.44 | 0.44 | 0.44 |
| Republic of Kalmykia | 0.47 | 0.55 | 0.65 | 0.69 | 0.84 | 0.88 | 0.88 | 1.00 | 1.00 | 0.98 | 0.92 | 0.82 | 0.57 | 0.50 | 0.55 | 0.45 | 0.63 | 0.55 | 0.53 |
| Republic of Karelia | 0.19 | 0.14 | 0.17 | 0.30 | 0.55 | 0.80 | 1.00 | 0.89 | 0.87 | 0.70 | 0.85 | 0.81 | 0.76 | 0.71 | 0.63 | 0.64 | 0.73 | 0.69 | 0.68 |
| Republic of Khakassia | 0.13 | 0.23 | 0.24 | 0.27 | 0.35 | 0.52 | 0.69 | 0.98 | 1.00 | 0.86 | 0.76 | 0.49 | 0.60 | 0.49 | 0.38 | 0.27 | 0.33 | 0.36 | 0.35 |
| Republic of Mordovia | 0.15 | 0.17 | 0.17 | 0.17 | 0.30 | 0.43 | 0.44 | 0.63 | 0.67 | 0.68 | 0.69 | 0.73 | 0.77 | 0.75 | 0.74 | 0.77 | 1.00 | 0.76 | 0.57 |
| Republic of North Ossetia-Alania | 0.25 | 0.30 | 0.46 | 1.00 | 0.18 | 0.04 | 0.41 | 0.39 | 0.45 | 0.45 | 0.34 | 0.27 | 0.30 | 0.31 | 0.25 | 0.23 | 0.11 | 0.09 | 0.09 |
| Republic of Tatarstan | 0.32 | 0.29 | 0.29 | 0.41 | 0.50 | 0.61 | 0.77 | 0.84 | 0.70 | 0.77 | 0.75 | 0.68 | 0.80 | 0.84 | 0.84 | 0.86 | 0.95 | 0.95 | 1.00 |
| Rostov Oblast | 0.85 | 1.00 | 0.89 | 0.66 | 0.65 | 0.71 | 0.61 | 0.61 | 0.53 | 0.68 | 0.81 | 0.81 | 0.84 | 0.74 | 0.65 | 0.73 | 0.73 | 0.71 | 0.68 |
| Ryazan Oblast | 0.20 | 0.16 | 0.26 | 0.30 | 0.35 | 0.36 | 0.40 | 0.53 | 0.50 | 0.72 | 0.80 | 0.80 | 0.84 | 1.00 | 0.82 | 0.70 | 0.84 | 0.80 | 0.79 |
| Saint Petersburg | 0.26 | 0.25 | 0.42 | 0.60 | 0.61 | 0.64 | 0.45 | 0.99 | 0.73 | 0.88 | 1.00 | 0.74 | 0.81 | 0.66 | 0.77 | 0.65 | 0.79 | 0.78 | 0.68 |
| Sakha (Yakutia) Republic | 0.22 | 0.24 | 0.34 | 0.37 | 0.47 | 0.62 | 0.68 | 0.76 | 0.83 | 0.91 | 1.00 | 0.63 | 0.68 | 0.77 | 0.76 | 0.62 | 0.71 | 0.67 | 0.57 |
| Sakhalin Oblast | 0.24 | 0.18 | 0.29 | 0.49 | 0.60 | 0.59 | 0.58 | 0.78 | 0.70 | 0.65 | 0.92 | 0.88 | 1.00 | 0.97 | 0.92 | 0.82 | 0.91 | 1.00 | 0.86 |
| Samara Oblast | 0.40 | 0.26 | 0.30 | 0.26 | 0.26 | 0.30 | 0.29 | 0.50 | 0.54 | 0.83 | 0.91 | 1.00 | 0.89 | 0.85 | 0.91 | 0.81 | 0.80 | 0.81 | 0.89 |
| Saratov Oblast | 0.34 | 0.28 | 0.33 | 0.49 | 0.53 | 0.49 | 0.53 | 0.61 | 0.66 | 0.74 | 1.00 | 0.95 | 0.93 | 0.90 | 0.90 | 0.79 | 0.84 | 0.82 | 0.75 |
| Sevastopol | 0.95 | 0.95 | 0.95 | 0.95 | 0.95 | 0.95 | 0.95 | 0.95 | 0.95 | 0.95 | 0.95 | 0.95 | 0.95 | 0.95 | 0.95 | 0.95 | 1.00 | 0.89 | 0.96 |
| Smolensk Oblast | 0.32 | 0.32 | 0.46 | 0.62 | 0.81 | 0.82 | 0.74 | 0.69 | 0.69 | 0.79 | 0.87 | 0.95 | 1.00 | 0.87 | 0.75 | 0.70 | 0.81 | 0.75 | 0.74 |
| Stavropol Krai | 1.00 | 0.67 | 0.73 | 0.67 | 0.68 | 0.69 | 0.70 | 0.77 | 0.83 | 0.77 | 0.77 | 0.73 | 0.61 | 0.68 | 0.62 | 0.51 | 0.55 | 0.44 | 0.43 |
| Sverdlovsk Oblast | 0.41 | 0.35 | 0.37 | 0.41 | 0.57 | 0.67 | 0.72 | 0.79 | 0.82 | 0.96 | 0.94 | 0.83 | 1.00 | 0.81 | 0.71 | 0.66 | 0.68 | 0.67 | 0.71 |
| Tambov Oblast | 0.26 | 0.23 | 0.32 | 0.38 | 0.38 | 0.36 | 0.29 | 0.26 | 0.30 | 0.62 | 0.73 | 0.82 | 0.82 | 0.87 | 1.00 | 0.79 | 0.79 | 0.74 | 0.73 |
| Tomsk Oblast | 0.56 | 0.54 | 0.62 | 0.65 | 0.92 | 0.91 | 0.84 | 0.93 | 0.73 | 0.95 | 1.00 | 0.79 | 0.88 | 0.86 | 0.92 | 0.88 | 0.64 | 0.59 | 0.59 |
| Tula Oblast | 0.35 | 0.41 | 0.41 | 0.51 | 0.54 | 0.54 | 0.71 | 0.71 | 0.77 | 0.83 | 0.88 | 1.00 | 0.95 | 0.97 | 0.98 | 0.74 | 0.81 | 0.85 | 0.85 |
| Tuva Republic | 0.37 | 0.35 | 0.31 | 0.31 | 0.39 | 0.41 | 0.41 | 0.53 | 0.35 | 0.37 | 0.39 | 0.57 | 1.00 | 0.83 | 0.74 | 0.53 | 0.41 | 0.37 | 0.37 |
| Tver Oblast | 0.24 | 0.23 | 0.38 | 0.47 | 0.57 | 0.64 | 0.79 | 0.70 | 0.60 | 0.61 | 1.00 | 0.84 | 0.89 | 0.82 | 0.93 | 0.67 | 0.76 | 0.73 | 0.74 |
| Tyumen Oblast | 0.39 | 0.53 | 0.60 | 0.68 | 0.75 | 0.65 | 0.64 | 0.55 | 0.39 | 0.86 | 0.94 | 0.92 | 1.00 | 0.95 | 0.96 | 0.87 | 0.92 | 0.68 | 0.66 |
| Udmurt Republic | 0.33 | 0.21 | 0.14 | 0.19 | 0.25 | 0.35 | 0.49 | 0.54 | 0.62 | 0.73 | 0.83 | 1.00 | 1.00 | 0.73 | 0.73 | 0.52 | 0.54 | 0.51 | 0.51 |
| Ulyanovsk Oblast | 0.13 | 0.14 | 0.10 | 0.14 | 0.28 | 0.54 | 0.72 | 0.68 | 0.59 | 0.70 | 0.90 | 0.84 | 0.92 | 0.97 | 0.87 | 0.70 | 0.84 | 1.00 | 0.87 |
| Vladimir Oblast | 0.16 | 0.19 | 0.24 | 0.25 | 0.29 | 0.33 | 0.39 | 0.35 | 0.46 | 0.72 | 0.90 | 0.89 | 0.94 | 0.88 | 0.83 | 0.75 | 1.00 | 0.96 | 0.92 |
| Volgograd Oblast | 0.49 | 0.59 | 0.30 | 0.36 | 0.41 | 0.46 | 0.56 | 0.59 | 0.56 | 0.83 | 1.00 | 0.86 | 0.85 | 0.82 | 0.76 | 0.78 | 0.93 | 0.88 | 0.93 |
| Vologda Oblast | 0.21 | 0.14 | 0.20 | 0.43 | 0.65 | 1.00 | 0.79 | 0.59 | 0.42 | 0.52 | 0.56 | 0.71 | 0.51 | 0.57 | 0.54 | 0.45 | 0.55 | 0.47 | 0.47 |
| Voronezh Oblast | 0.17 | 0.28 | 0.45 | 0.48 | 0.37 | 0.41 | 0.45 | 0.47 | 0.49 | 0.63 | 0.73 | 0.95 | 0.88 | 0.80 | 1.00 | 0.93 | 0.77 | 0.73 | 0.73 |
| Yamalo-Nenets Autonomous Okrug | 0.39 | 0.53 | 0.60 | 0.68 | 0.75 | 0.65 | 0.64 | 0.55 | 0.39 | 0.86 | 0.94 | 0.92 | 1.00 | 0.82 | 0.82 | 0.75 | 0.65 | 0.60 | 0.58 |
| Yaroslavl Oblast | 0.31 | 0.32 | 0.31 | 0.38 | 0.40 | 0.34 | 0.34 | 0.49 | 0.50 | 0.86 | 0.92 | 0.89 | 1.00 | 0.91 | 0.86 | 0.67 | 0.88 | 0.82 | 0.90 |
| Zabaykalsky Krai | 0.06 | 0.07 | 0.12 | 0.21 | 0.29 | 0.36 | 0.96 | 1.00 | 0.84 | 0.41 | 0.42 | 0.40 | 0.39 | 0.40 | 0.41 | 0.39 | 0.41 | 0.41 | 0.39 |

Figure 6: Heat Map of Regional Wine Sales

## 5.2 Elbow Method

A major drawback to the K-Means algorithm is having to specify the number of clusters, $k$. The elbow method will help remedy this issue. Essentially, the elbow method works by running the K-Means algorithm $n$ number of times. Then comparing the mean inertia at each iteration of $n$. Inertia in this context refers to

6

the distance from each instance (data point) to its closest cluster. Note as the number of clusters approach the number of instances, the inertia tends to 0. Thus, it does not suffice to take the higher number of clusters to minimize inertia. Rather you want to pick the number of clusters, $k$ around the point where you have inertia but $k + 1$ does not lower inertia by a substantial amount. In this project, we will select $n = 29$. Therefore, the K-Means algorithm was ran 29 times, keeping track of the mean inertia at each iteration. In Figure 7 we can see the graph of the number of clusters $k$ versus the mean inertia. Since the drop-off of inertia is minuscule from $k = 4$ to $k = 5$, we will take $k = 4$ to be the number of clusters for this analysis.



Figure 7: Line Plot: Number of Clusters versus Mean Inertia

## 5.3 Post K-Means Analysis

Now that we have the four clusters, we should check which regions are in the same cluster as St Petersburg. Again, we are interested in this cluster since the wine promotion was success full in St Petersburg, and we want to recreate this success. In Figure 8 we can see the overall trend of wine sales for regions in each cluster. For example, cluster 0 has a progressive increase in wine sales, peaking around 2008, with fluctuations after 2008. It turns out that St Petersburg is apart of cluster 0, so this cluster will be important. In fact, we could end our analysis here by randomly taking ten regions from this cluster and running the wine promotion there. Recall however, our budget is small, so let's get more detailed.

In figure 9 we can see the different regression lines for each region. Of course the regions with promising, positive slopes will be of most interest. However, we need some way of ranking regions to pick the top ten regions most likely to have a successful wine campaign. By scoring regions based on the slope of their

7

Figure 8: Heat Map of Wine Sales Trend per Cluster

regression line, and average wine sales over the recent (four) years, we can arrive at the following ten regions: Republic of Karelia, Kirov Oblast, Leningrad Oblast, Sverdlovsk Oblast, Chukotka Autonomous Okrug, Kaliningrad Oblast, Pskov Oblast, St Petersburg, Kursk Oblast and Penza Oblast.



Figure 9: Linear Regression of regions in St Petersburg's Cluster

# 6    Hierarchical Clustering

Out of curiosity, a hierarchical clustering algorithm was run on the recent wine sales to see if it would produce a similar result as the K-Means clustering algorithm that was presented earlier. The Hierarchical Clustering produced the deprogram as seen in Figure 12.



Figure 10: Dendrogram

From here we can chose the number of clusters. Selecting four to be the number of clusters, led to St Petersburg being in a cluster of 37 other regions. All top regions picked from the K-Means cluster appeared in the Hierarchical cluster, and thus no new information was gained. However, changing the cluster count to 8, we get a more refined St Petersburg cluster of size seven. In the end there were four regions which were in both the K-Means cluster and the Hierarchical cluster. These four regions are: Chukotla Autonomous Okrug, Kaliningrad Oblast, Kursk Oblast and (unsurprisingly) St Petersburg.

# 7    Conclusion

In the end we have ten regions of interest: Republic of Karelia, Kirov Oblast, Leningrad Oblast, Sverdlovsk Oblast, **Chukotka Autonomous Okrug**, **Kaliningrad Oblast**, Pskov Oblast, **St Petersburg**, **Kursk Oblast** and Penza Oblast. Note the boldface regions matched with the Hierarchical clustering and K-Means clustering. Hence, we emphasize those as the highest priority regions. Lastly, it is important to mention any aspect of the report which could be improved. First, when faced with missing data or outliers the solution involved imputing reasonable means. However, this is not always the best solution, and other

patching methods may prove more adequate. Secondly, transforming the wine data so that it more accurately represents a normal distribution could lead to more precise clustering. Finally we could improve the K-Means clustering by testing different values of $k$ and running the model on new data to check its accuracy. We conclude this report with two final graphs. In the following two figures we are plotting the slope of the linear regression versus the average wine sales.
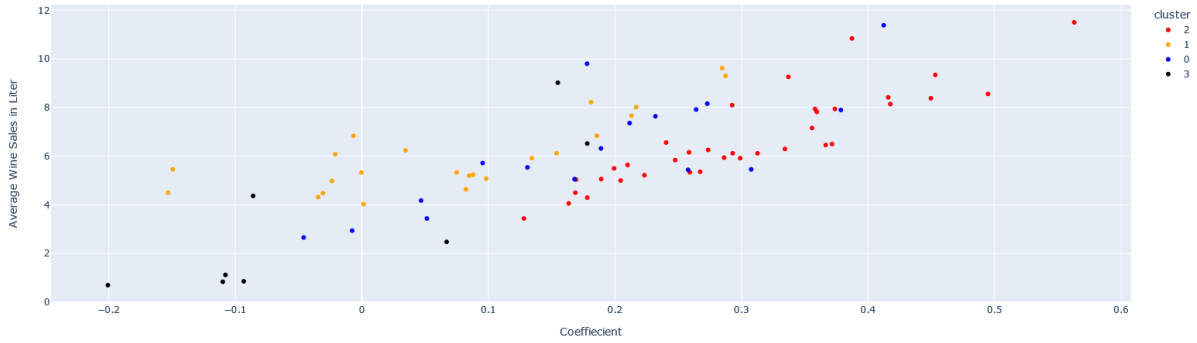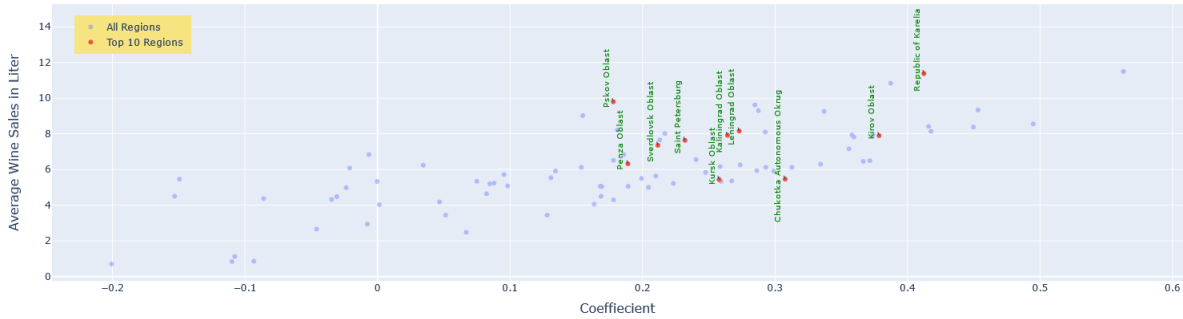


Figure 11: Dendrogram



Figure 12: Dendrogram

10