

FEDERAL STATE AUTONOMOUS EDUCATIONAL INSTITUTION
OF HIGHER EDUCATION

NATIONAL RESEARCH UNIVERSITY

HIGHER SCHOOL OF ECONOMICS

Faculty of Computer Sciences

Department of Data Analysis and Artificial Intelligence

Educational program “Data Science”

MODERN DATA ANALYSIS

HOMEWORK 2020

Govorova Irina

Telnov Sergey

Shenker Anastasia

Yakovlev Daniil

Moscow, 2020

Dataset description

For the implementation of this project, a “Star dataset” was chosen, since it approaches well for the task of clustering, represents an interesting topic and gives a possibility to explore some enthralling issues.

The source of the “Star dataset” is a system for organizing competitions for data research and a social network for data processing and machine learning specialists Kaggle¹.

“Star dataset” consists of 241 observations and the following features:

- Temperature (quantitative)
- Luminosity (quantitative)
- Radius (quantitative)
- Absolute Magnitude (quantitative)
- Star Type (categorical)
- Star Color (categorical)
- Spectral Class (categorical)

Temperature

“Temperature” is an absolute temperature of the star measured in Kelvins.

Luminosity

“Luminosity” is a relative luminosity of a star, which is calculated as $\frac{L}{L_0}$, where L is a luminosity of a star and $L_0 = 3.828 \cdot 10^{26}$ Watts is an average luminosity of the Sun. Luminosity of the star is the total amount of electromagnetic energy emitted by a star per unit of time. The luminosity of the main sequence stars can be approximately calculated by the formula:

$$L = 4\pi R^2 \cdot \sigma T^4$$

R – star radius

T – star photosphere temperature

σ – Stefan - Boltzmann constant

¹ <https://www.kaggle.com/deepu1109/star-dataset?select=6+class+csv.csv>

Radius

“Radius” is a relative radius of a star, which is calculated as $\frac{R}{R_0}$, where R is a radius of a star and $R_0 = 6.9551 \cdot 10^8$ m is an average radius of the Sun.

Absolute Magnitude

“Absolute Magnitude” is a value of a star absolute magnitude measured in M_v . Absolute magnitude is a measure of the luminosity of a celestial object. For stars it is defined as the apparent magnitude of an object if it was located at a distance of 10 parsecs from the observer and would not experience any interstellar or atmospheric absorption.

Star Type

“Star Type” is for one of the following types: Red Dwarf, Brown Dwarf, White Dwarf, Main Sequence, Supergiants, Hypergiants.

A red dwarf is the smallest and coolest kind of star on the main sequence. Red dwarfs are by far the most common type of star in the Milky Way, at least in the neighborhood of the Sun, but because of their low luminosity, individual red dwarfs cannot be easily observed.

A brown dwarf is a type of substellar object that has a mass between the most massive gas giant planets and the least massive stars.

A white dwarf is a stellar core remnant composed mostly of electron-degenerate matter. A white dwarf is very dense: its mass is comparable to that of the Sun, while its volume is comparable to that of Earth. A white dwarf's faint luminosity comes from the emission of stored thermal energy.

Main-sequence stars or dwarf stars are the most numerous true stars in the universe, and include the Earth's Sun. The main sequence is a continuous and distinctive band of stars that appears on plots of stellar color versus brightness.

Supergiants are among the most massive and most luminous stars.

A hypergiant is a very rare type of star that has an extremely high luminosity, mass, size and mass loss because of their extreme stellar winds.

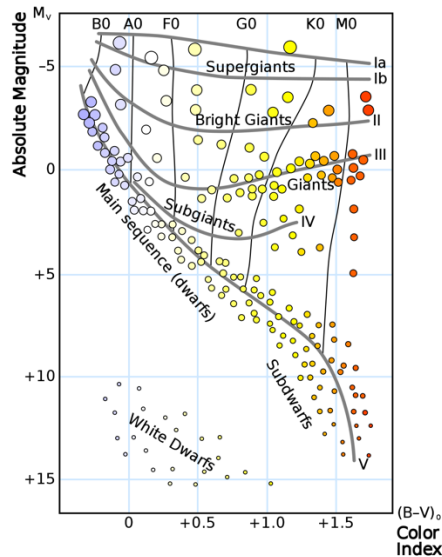


Fig. 1. The Hertzsprung–Russell diagram relates stellar classification with absolute magnitude, luminosity, and surface temperature.

Star Color

“Star Color” is a color of a star, which depends on its temperature and can be White, Red, Blue, Blue-White, Yellow, Yellow-Orange, Orange-Red, etc.

Spectral Class

“Spectral Class” is one of the following spectral type: O, B, A, F, G, K, M according to the Morgan–Keenan (MK) classification system from the hottest (O type) to the coolest (M type).

K-Means clustering

For K-Means clustering all 4 quantitative features were considered:

- «Temperature»
- «Luminosity»
- «Radius»
- «Absolute magnitude»

We chose all these features because they represent different characteristics of stars and they are not linearly dependent. All these four features are equally significant for category definition.

We standardized all columns and then applied K-Means algorithm to split the data into 4 clusters and into 7 clusters. K-Means was applied 10 times with following parameters: maximum number of iterations 500, random initial cluster centers, relative tolerance 0.0001, EM-style algorithm and different random state, which equals the number of iteration. From 10 results of clustering we chose the best according to the parameter «inertia»: the less inertia, the better clustering.

In the following table, you can see results of K-Means initializations, where «i» is a number of initialization and a value of random state.

Table 1. Results of K-Means initializations

	i = 0	i = 1	i = 2	i = 3	i = 4	i = 5	i = 6	i = 7	i = 8	i = 9
4 clusters	281.5	257.21	283.93	277.05	264.86	282.05	281.51	281.51	257.21	256.9
7 clusters	140.63	147.51	136.51	137.01	137.01	137.01	198.68	137.06	192.15	137.01

The best result for 4 clusters was received in the initialization with $i = 9$. The best result for 7 clusters was received in the initialization with $i = 2$.

There is a more thorough interpretation of all four clusters is given below.

Relative difference

In this section, relative difference for all clusters and all quantitative features are counted. Relative difference shows how much cluster mean of a feature differs from grand mean of the same feature.

In the following table, you can see relative differences for every feature and for every cluster of 4-cluster K-Means.

Table 2. Relative differences for every feature and for every cluster of 4-cluster K-Means

	Temperature	Luminosity	Radius	Absolute magnitude
Relative difference for cluster “0”	-53.53	134.18	476.36	-327.05
Relative difference for cluster “1”	171.34	423.36	337.89	-288.69
Relative difference for cluster “2”	-39.18	-99.99	-99.85	192.67
Relative difference for cluster “3”	71.66	59.37	-81.71	-201.15

In the following table, you can see relative differences for every feature and for every cluster of 7-cluster K-Means.

Table 3. Relative differences for every feature and for every cluster of 7-cluster K-Means

	Temperature	Luminosity	Radius	Absolute magnitude
Relative difference for cluster “0”	-64.35	-100	-99.91	241.32
Relative difference for cluster “1”	-9.93	-84.26	-96.54	-94.11
Relative difference for cluster “2”	-51.24	130.11	467.12	-325.33
Relative difference for cluster “3”	171	56.89	-88.37	-213.82
Relative difference for cluster “4”	2.32	216.21	-76.95	249.08
Relative difference for cluster “5”	196.75	390.51	449.73	-299.48
Relative difference for cluster “6”	53.71	-100	-100	172.46

The absolute value of relative difference in these tables shows the percent of deviation from the grand mean of the feature. The greater the value, the greater the deviation. The sign shows the direction of deviation. Minus means that the mean of

the cluster is less than the grand mean and plus means that the mean of the cluster is greater than the grand mean.

Quetelet indeces

Quetelet index shows what part of category belongs to a cluster. It doesn't show an exact part of category but show it comparably: the less absolute value the bigger part of cluster belongs to a category.

Quetelet index is calculated by the formula:

$$q_{kv} = 100 \left[\frac{p_{kv}}{p_k p_v} - 1 \right],$$

where p_{kv} – the proportion of objects belonging to a cluster and a category simultaneously, p_k – the proportion of objects belonging to a cluster, p_v – the proportion of objects belonging to a category.

In the following tables, you can see quetelet indeces for each category of features «Star type», «Star color» and «Star Class» for 4-cluster splitting.

Table 4. Quetelet indeces for star types for 4-cluster K-Means

	0 – Red Dwarf	1 – Brown Dwarf	2 – White Dwarf	3 – Main Sequence	4 – SuperGiants	5 – HyperGiants
Quetelet index for cluster “0”	-100	-100	-100	-100	-100	-100
Quetelet index for cluster “1”	-100	-100	-100	-100	71.43	328.57
Quetelet index for cluster “2”	77.78	77.78	64.44	-20	-100	-100
Quetelet index for cluster “3”	-100	-100	-70.97	112.9	248.39	-90.32

Table 5. Quetelet indeces for star colors for 4-cluster K-Means

	Red	Blue White	White	Yellowish White	Blue white	Pale yellow orange	Blue
Quetelet index for cluster “0”	69.955	-100	136.45	-100	-100	-100	-84.95
Quetelet index for cluster “1”	-100	-100	-100	-100	-100	-100	211.69
Quetelet index for cluster “2”	26.98	60	26.98	77.78	77.78	77.78	-67.68
Quetelet index for cluster “3”	-68.89	-61.29	-100	-100	-100	-100	139.3

	Blue-white	Whitish	yellow-white	Orange	White-Yellow	white	Blue
Quetelet index for cluster “0”	-68.17	-100	-100	727.59	-100	-100	-100
Quetelet index for cluster “1”	163.74	-100	-100	-100	-100	-100	-100
Quetelet index for cluster “2”	-79.49	77.78	77.78	-100	77.78	77.78	77.78
Quetelet index for cluster “3”	167.99	-100	-100	-100	-100	-100	-100

	yellowish	Yellowish	Orange-Red	Blue white	Blue-White
Quetelet index for cluster “0”	-100	-100	-100	-100	-100
Quetelet index for cluster “1”	-100	-100	-100	-100	-100
Quetelet index for cluster “2”	77.78	77.78	77.78	77.78	77.78
Quetelet index for cluster “3”	-100	-100	-100	-100	-100

Table 6. *Quetelet indices for spectral classes for 4-cluster K-Means*

	M	B	A	F	O	K	G
Quetelet index for cluster “0”	64.03	-82.01	-12.89	-100	-79.31	175.86	727.59
Quetelet index for cluster “1”	-100	86.34	-100	-100	285.71	-100	-100
Quetelet index for cluster “2”	28.13	-14.98	12.28	77.78	-100	18.52	-100
Quetelet index for cluster “3”	-68.61	51.47	1.87	-100	190.32	-100	-100

The value -100 means that no object belongs to a cluster and a category simultaneously.

The greater the number (the smaller the absolute number), the more objects from a cluster belong to a category simultaneously and the more objects from a category belong to a cluster.

The value 0 is possible only when all objects belong to one category and to one cluster.

These tables show the differences in clusters.

Making sets V^+ and V^- .

V^+ is a set, which is formed by features and categories, values of which is greater than 35%.

V^- is a set, which is formed by features and categories, values of which is less than -35%.

Sets for 4-cluster split.

Cluster «0»

V^+ for cluster 0 consist of following numeric features and categories: «Luminosity», «Radius», «HyperGiants», «Red», «White», «Orange», «M», «K», «G».

V^- for cluster 0 consist of following numeric features and categories: «Temperature», «Absolute magnitude», «*Red Dwarf*», «*Brown Dwarf*», «*White Dwarf*», «*SuperGiants*», «*Blue White*», «*Yellowish White*», «*Blue white*», «*Pale yellow orange*», «Blue», «Blue-white», «*Whitish*», «*yellow-white*», «Orange», «*White-Yellow*», «*white*», «Blue», «*yellowish*», «*Yellowish*», «*Orange-Red*», «*Blue white*», «*Blue-White*», «B», «F», «O».

Cluster «0» is less than average in these features. But we want to highlight that most color categories are equal -100 (features in italics), which means that there is no objects of such cluster and this is true for every cluster.

In general cluster «0» is very much greater in sizes and has red, orange and blue colors.

Cluster «1»

V^+ for cluster 1 consist of following numeric features and categories: «Temperature», «Luminosity», «Radius», «SuperGiants», «HyperGiants», «Blue», «Blue-white», «B», «O».

V^- for cluster 1 consist of following numeric features and categories: «Absolute magnitude», «*Red Dwarf*», «*Brown Dwarf*», «*White Dwarf*», «*Main*

Sequence», «Red», «Blue White», «White», «Yellowish White», «Blue white», «Pale yellow orange», «Whitish», «yellow-white», «Orange», «White-Yellow», «white», «Blue», «yellowish», «Yellowish», «Orange-Red», «Blue white», «Blue-White», «M», «A», «F», «K», «G».

In general cluster «1» is greater in sizes and has more blue and blue-white colors.

Cluster «2»

V^+ for cluster 2 consist of following numeric features and categories: «Absolute magnitude», «Red Dwarf», «Brown Dwarf», «White Dwarf», «Blue White», «Yellowish White», «Blue white», «Pale yellow orange», «Whitish», «yellow-white», «White-Yellow», «white», «Blue», «yellowish», «Yellowish», «Orange-Red», «Blue white», «Blue-White», «F».

V^- for cluster 2 consist of following numeric features and categories: «Luminosity», «Radius», «SuperGiants», «HyperGiants», «Blue», «Blue-white», «Orange», «O», «G».

Cluster «2» is very much less in sizes and has more white, yellow and blue colors.

Cluster «3»

V^+ for cluster 3 consist of following numeric features and categories: «Temperature», «Luminosity», «Main Sequence», «SuperGiants», «Blue», «Blue-white», «B», «O».

V^- for cluster 3 consist of following numeric features and categories: «Radius», «Absolute magnitude», «Red Dwarf», «Brown Dwarf», «White Dwarf», «HyperGiants», «Red», «Blue White», «White», «Yellowish White», «Blue white», «Pale yellow orange», «Whitish», «yellow-white», «Orange», «White-Yellow», «white», «Blue», «yellowish», «Yellowish», «Orange-Red», «Blue white», «Blue-White», «M», «F», «K», «G».

Cluster «3» is greater in temperature and luminosity, has more blue colors, less in size and magnitude.

Sets for 7-cluster split.

In the following tables, you can see quetelet indices for each category of features «Star type», «Star color» and «Star Class» for 7-cluster splitting.

Table 7. Quetelet indices for star types for 7-cluster K-Means

	0 – Red Dwarf	1 – Brown Dwarf	2 – White Dwarf	3 – Main Sequence	4 – SuperGiants	5 – HyperGiants
Quetelet index for cluster “0”	-100	-100	-100	-100	-100	-100
Quetelet index for cluster “1”	-100	-100	-100	-100	71.43	328.57
Quetelet index for cluster “2”	77.78	77.78	64.44	-20	-100	-100
Quetelet index for cluster “3”	-100	-100	-70.97	112.9	248.39	-90.32
Quetelet index for cluster “4”	-100	-100	-100	-100	500	-100
Quetelet index for cluster “5”	-100	-100	-100	-100	-45.45	445.45
Quetelet index for cluster “6”	-100	-100	500	-100	-100	-100

Table 8. Quetelet indices for star colors for 7-cluster K-Means

	Red	Blue White	White	Yellowish White	Blue white	Pale yellow orange	Blue
Quetelet index for cluster “0”	88.38	-73.63	50.71	-12.09	75.82	163.74	-100
Quetelet index for cluster “1”	-86.61	-100	-100	-100	-100	-100	-72.73
Quetelet index for cluster “2”	64.29	-100	128.57	-100	-100	-100	-85.45
Quetelet index for cluster “3”	-68.89	-61.29	-100	-100	-100	-100	139.3
Quetelet index for cluster “4”	-37.5	-100	-100	-100	-100	-100	209.9
Quetelet index for cluster “5”	-100	-100	-100	-100	-100	-100	177.69
Quetelet index for cluster “6”	-100	644.83	18.23	451.72	175.86	-100	80.56

	Blue-white	Whitish	yellow-white	Orange	White-Yellow	white	Blue
Quetelet index for cluster “0”	-100	-100	-100	-100	163.74	-12.09	-100
Quetelet index for cluster “1”	275	650	650	-100	-100	-100	-100
Quetelet index for cluster “2”	-38.46	-100	-100	700	-100	-100	-100
Quetelet index for cluster “3”	167.99	-100	-100	-100	-100	-100	-100
Quetelet index for cluster “4”	-100	-100	-100	-100	-100	-100	-100
Quetelet index for cluster “5”	235.66	-100	-100	-100	-100	-100	-100
Quetelet index for cluster “6”	-100	-100	-100	-100	-100	451.72	727.59

	yellowish	Yellowish	Orange-Red	Blue white	Blue-White
Quetelet index for cluster “0”	-100	-100	-100	-100	-100
Quetelet index for cluster “1”	650	650	650	-100	650
Quetelet index for cluster “2”	-100	-100	-100	-100	-100
Quetelet index for cluster “3”	-100	-100	-100	-100	-100
Quetelet index for cluster “4”	-100	-100	-100	-100	-100
Quetelet index for cluster “5”	-100	-100	-100	-100	-100
Quetelet index for cluster “6”	727.59	-100	-100	727.59	-100

Table 9. Quetelet indices for spectral classes for 7-cluster K-Means

	M	B	A	F	O	K	G
Quetelet index for cluster “0”	90.08	-100	-2.83	-37.94	-100	-100	-100
Quetelet index for cluster “1”	-86.49	-2.17	294.74	252.94	-62.5	400	-100
Quetelet index for cluster “2”	58.56	-65.22	-15.79	-100	-80	166.67	700
Quetelet index for cluster “3”	-100	104.16	-100	-100	265.22	-100	-100
Quetelet index for cluster “3”	-36.94	-100	-100	-100	325	-100	-100
Quetelet index for cluster “3”	-100	137.15	-100	-100	227.27	-100	-100

Quetelet index for cluster “3”	-100	331.78	-100	143.41	-100	-100	-100
--------------------------------	------	--------	------	--------	------	------	------

Cluster «0»

V^+ for cluster 0 consist of following numeric features and categories: «Absolute magnitude», «Red Dwarf», «Brown Dwarf», «Red», «White», «Blue white», «Pale yellow orange», «White-Yellow», «M»,

V^- for cluster 0 consist of following numeric features and categories: «Temperature», «Luminosity», «Radius», «Main Sequence», «SuperGiants», «HyperGiants», «Blue White», «Blue», «Blue-white», «Whitish», «yellow-white», «Orange», «Blue», «yellowish», «Yellowish», «Orange-Red», «Blue white», «Blue-White», «B», «O», «K», «G».

Cluster «0» is less in sizes, temperature and luminosity, has more red, orange and white colors.

Cluster «1»

V^+ for cluster 1 consist of following numeric features and categories: «Main Sequence», «Blue-white», «Whitish», «yellow-white», «yellowish», «Yellowish», «Orange-Red», «Blue-White», «A», «F», «K».

V^- for cluster 1 consist of following numeric features and categories: «Luminosity», «Radius», «Absolute magnitude», «Red Dwarf», «Brown Dwarf», «White Dwarf», «HyperGiants», «Red», «Blue White», «White», «Yellowish White», «Blue white», «Pale yellow orange», «Blue», «Orange», «White-Yellow», «white», «Blue», «Blue white», «M», «O», «G».

Cluster 1 has more size and more yellow colors.

Cluster «2»

V^+ for cluster 2 consist of following numeric features and categories: «Luminosity», «Radius», «HyperGiants», «Red», «White», «Orange», «M», «K», «G».

V^- for cluster 2 consist of following numeric features and categories: «Temperature», «Absolute magnitude», «*Red Dwarf*», «*Brown Dwarf*», «*White Dwarf*», «*Main Sequence*», «*SuperGiants*», «*Blue White*», «*Yellowish White*», «*Blue white*», «*Pale yellow orange*», «*Whitish*», «*yellow-white*», «*White-Yellow*», «*white*», «*Blue*», «*yellowish*», «*Yellowish*», «*Orange-Red*», «*Blue white*», «*Blue-White*», «B», «F».

Cluster 2 has much greater luminosity and size and more red and orange.

Cluster «3»

V^+ for cluster 3 consist of following numeric features and categories: «Temperature», «Luminosity», «Main Sequence», «SuperGiants», «Blue», «Blue-white», «B», «F».

V^- for cluster 3 consist of following numeric features and categories: «Radius», «Absolute magnitude», «*Red Dwarf*», «*Brown Dwarf*», «*White Dwarf*», «*HyperGiants*», «Red», «Blue White», «*White*», «*Yellowish White*», «*Blue white*», «*Pale yellow orange*», «*Whitish*», «*yellow-white*», «Orange», «*White-Yellow*», «*white*», «*Blue*», «*yellowish*», «*Yellowish*», «*Orange-Red*», «*Blue white*», «*Blue-White*», «M», «A», «O», «K», «G».

Cluster 3 has much higher temperature and luminosity and more blue colors.

Cluster «4»

V^+ for cluster 4 consist of following numeric features and categories: «Luminosity», «SuperGiants», «Blue», «O».

V^- for cluster 4 consist of following numeric features and categories: «Radius», «Absolute magnitude», «*Red Dwarf*», «*Brown Dwarf*», «*White Dwarf*», «*Main Sequence*», «*HyperGiants*», «Red», «Blue White», «*White*», «*Yellowish*

White», «*Blue white*», «*Pale yellow orange*», «*Blue-white*», «*Whitish*», «*yellow-white*», «*Orange*», «*White-Yellow*», «*white*», «*Blue*», «*yellowish*», «*Yellowish*», «*Orange-Red*», «*Blue white*», «*Blue-White*», «*M*», «*B*», «*A*», «*F*», «*K*», «*G*».

Cluster 4 has a very much greater size and luminosity and more blue colors.

Cluster «5»

V^+ for cluster 5 consist of following numeric features and categories: «Temperature», «Luminosity», «Radius», «HyperGiants», «Blue», «Blue-white», «B», «O».

V^- for cluster 5 consist of following numeric features and categories: «Absolute magnitude», «*Red Dwarf*», «*Brown Dwarf*», «*White Dwarf*», «*Main Sequence*», «*Red*», «*Blue White*», «*White*», «*Yellowish White*», «*Blue white*», «*Pale yellow orange*», «*Whitish*», «*yellow-white*», «*Orange*», «*White-Yellow*», «*white*», «*Blue*», «*yellowish*», «*Yellowish*», «*Orange-Red*», «*Blue white*», «*Blue-White*», «*M*», «*A*», «*F*», «*K*», «*G*».

Cluster 4 has a very much greater size, luminosity and temperature and more blue colors.

«Temperature», «Luminosity», «Radius», «Absolute magnitude», «Red Dwarf», «Brown Dwarf», «White Dwarf», «Main Sequence», «SuperGiants», «HyperGiants», «Red», «Blue White», «White», «Yellowish White», «Blue white», «Pale yellow orange», «Blue», «Blue-white», «Whitish», «yellow-white», «Orange», «White-Yellow», «white», «Blue», «yellowish», «Yellowish», «Orange-Red», «Blue white», «Blue-White», «M», «B», «A», «F», «O», «K», «G».

Cluster «6»

V^+ for cluster 6 consist of following numeric features and categories: «Temperature», «Absolute magnitude», «White Dwarf», «Blue White», «Yellowish White», «Blue white», «Blue», «white», «Blue», «Blue white»,

V^- for cluster 6 consist of following numeric features and categories: «Luminosity», «Radius», «Red Dwarf», «Brown Dwarf», «Main Sequence», «SuperGiants», «HyperGiants», «Red», «Pale yellow orange», «Blue-white», «Whitish», «yellow-white», «Orange», «White-Yellow», «yellowish», «Yellowish», «Orange-Red», «Blue-White».

Cluster 6 has less size, higher temperature and more blue and white colors.

Conceptualization

Now let us have a look on the results of clustering according to the types of stars. We will assume that a type of star belongs to a cluster, which contains most of the objects of this type.

In the following graph, we can see distribution of different types of stars for 4 clusters.

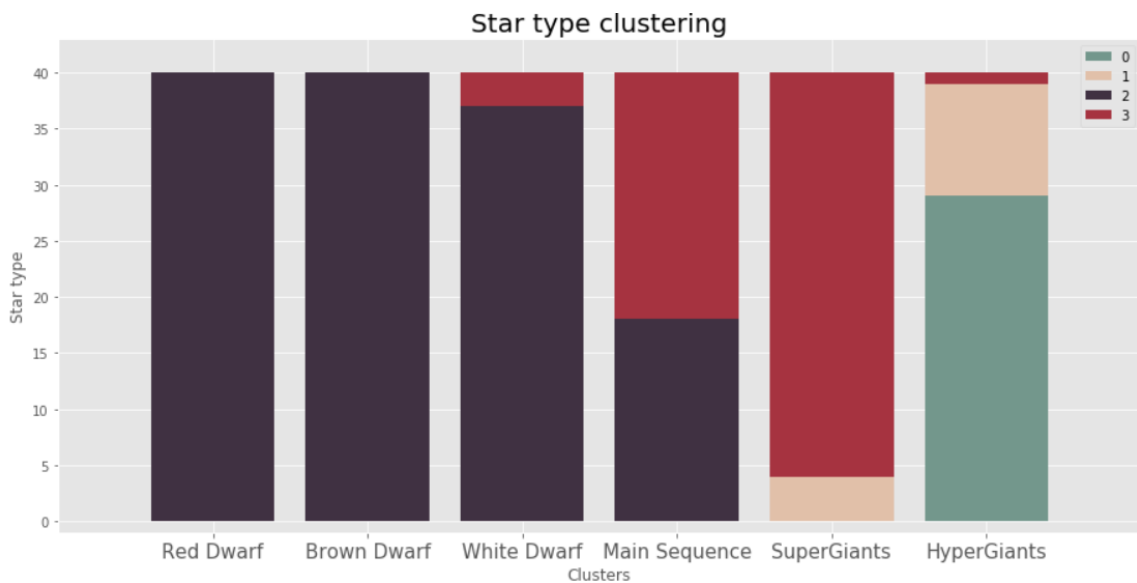


Fig. 2. Distribution of different types of stars for 4 clusters

According to our assumption, we can say that «Red Dwarf», «Brown Dwarf» and «White Dwarf» are in the cluster «2», «SuperGiants» and «HyperGiants» are in the cluster «3». «HyperGiants» are in the cluster «0». The cluster «1» contains some stars of «SuperGiants» type and some stars of «HyperGiants».

Certainly, objects are not splitted perfect. There are few objects of «WhiteDwarf» and «HyperGiants» belong to the cluster «3» and almost a half of objects of «Main Sequence» belong to the cluster «2».

In the following graph, we can see distribution of different types of stars for 4 clusters.

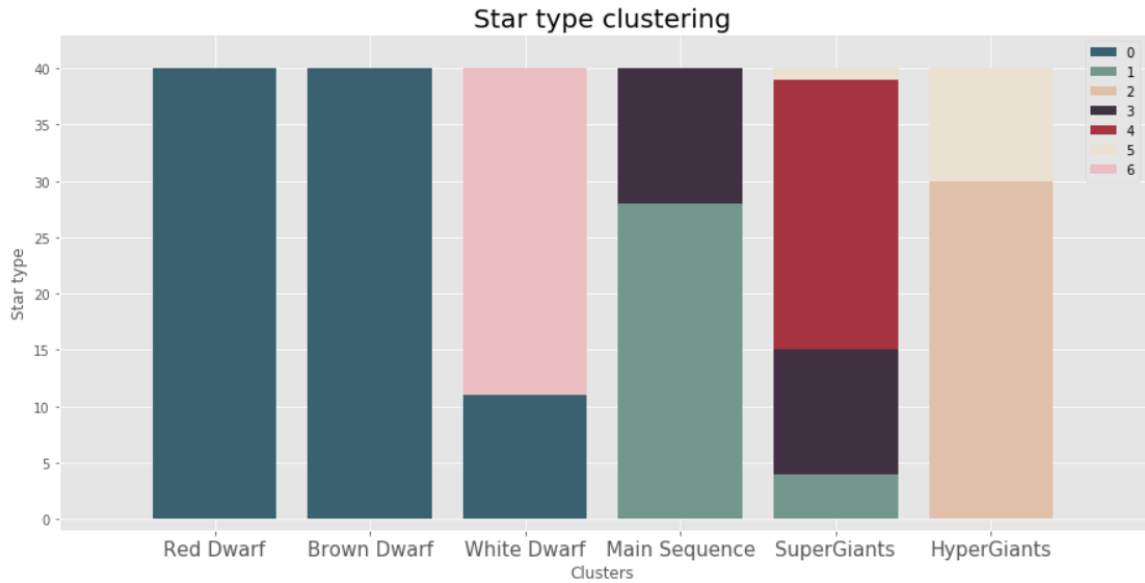


Fig. 3. Distribution of different types of stars for 4 clusters

According to our assumption, we can say that «Red Dwarf», «Brown Dwarf» are in the cluster «0», «White Dwarf» is in cluster «6», «Main Sequence» is in cluster «1» «SuperGiants» are in the cluster «4». «HyperGiants» are in the cluster «2». The cluster «5» contains some stars of «HyperGiants» and «SuperGiants» and cluster «3» contains some stars of «Main Sequence» and «SuperGiants».

K-Means in both cases (4 clusters and 7 clusters) split stars according to their type. We think that 4-cluster splitting is better. Because it is more interpretable with given features in dataset for people without astronomic education. It splits dwarf stars, stars with average / a little bit bigger size and giant stars, while 7-cluster splitting select more subtypes of stars which are not selected in usual classifications.

Also 4-cluster splitting cope with big number of color categories better. There are a lot of the same colors which are written differently (for example: «Blue White»,

«Blue-white», «Blue white», «Blue-White») but in fact they are the same, little number of clusters let to ignore this thing.

Bootstrap

For bootstrap the feature “Temperature” and 4-cluster partition with a number of initialization and a value of random state equal to 9 from previous part are used. For each application of bootstrap non-pivotal and pivotal versions are considered.

The first step is to find the 95% confidence interval for feature’s grand mean by using bootstrap.

As can be seen from the following table, the average temperature of the considered stars is between 9 206.58 K and 11 714.13 K for non-pivotal version and between 9 308.09 and 11 681.74 for pivotal version according to 95% confidence interval.

Table 10. Grand mean and 95% confidence intervals for “Temperature”

	Value
Grand mean	10 497.46
Non-pivotal 95% confidence interval	9 206.58 – 11 714.13
Pivotal 95% confidence interval	9 308.09 – 11 681.74

The second step is to take 2 clusters of the feature and for them to compare the within-cluster means using bootstrap. Cluster 0 and cluster 1 are considered.

As can be seen from the following table, the 95% confidence interval for the difference of within-cluster means is from 19 084.64 K to 27 937.23 K for non-pivotal version and from 19 186.28 K to 27 864.09 K for pivotal version.

Table 11. 95% confidence intervals for within-cluster means difference

	Value
Non-pivotal 95% confidence interval	19 084.64 – 27 937.23

Pivotal 95% confidence interval 19 186.28 – 27 864.09

So, the temperature in cluster 1 is in average at 23 510.935 K (23 525.23 K for pivotal version) higher than in cluster 0.

The next step is to compare the grand mean of the cluster with the within-cluster mean for the feature by using bootstrap. The grand mean is considered for the cluster 1.

As can be seen from the following table, the difference of the cluster 1 grand mean and within-cluster mean for the feature “Temperature” is between 13 646.25 K and 22 495.87 K for non-pivotal version and between 13 577.53 K and 22 457.15 K for pivotal version.

Table 12. 95% confidence intervals for the difference between cluster 1 grand mean and within-cluster mean for the feature

	Value
Non-pivotal 95% confidence interval	13 646.25 – 22 495.87
Pivotal 95% confidence interval	13 577.53 – 22 457.15

So, the average temperature in cluster 1 is in average at K 18 071.04 (18 017.34 K for pivotal version) higher than the average temperature in other clusters.

Contingency Table

For contingency table 3 nominal features were considered. One of them was taken from nominal features of

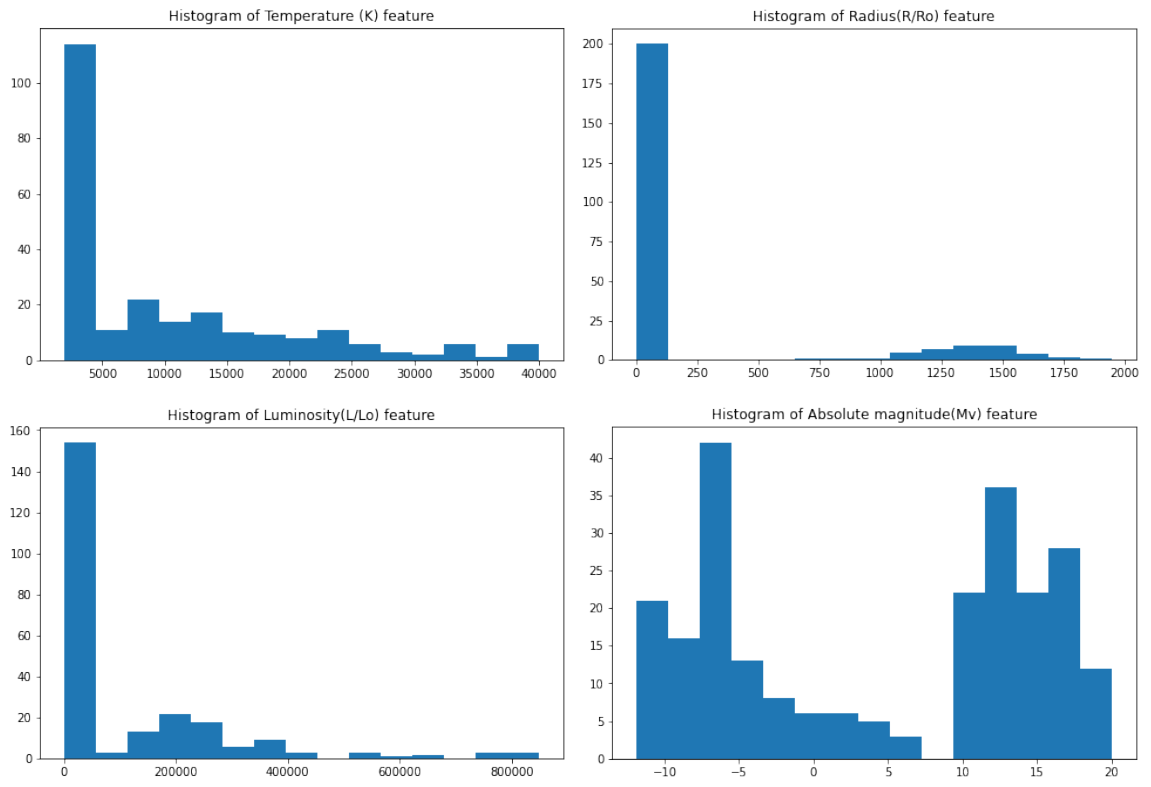


Fig. 4. Histograms of features

PCA/SVD

For PCA and SVD the same features as for K-Means were chosen: «Temperature», «Luminosity», «Radius», «Absolute Magnitude», as they represent different linearly independent characteristics, all of which are important.

Standardizing data

We made three standardizations of the data: z-scoring, range and rank. Z-scoring is the same as StandardScaler from sklearn.preprocessing module. It transforms feature into equivalent with 0 as the mean and 1 as standard deviation. Range standardization is done by subtracting the mean (as in z-scoring standardization) and dividing by range (difference between maximum and minimum of the feature). Rank standardization is done by subtracting mean and dividing by range. Rank standardization image a feature to a segment $[0,1]$ and it is an equivalent of MinMaxScaler from sklearn.preprocessing module.

SVD

We made SVD for original data and for standardized datasets. We cannot show whole matrices of SVD because their shape is not appropriate for this. But in the following screenshots you can see these matrices partially.

SVD on original data:


```

SVD
U: [[ 3.91972718e-05 -1.75769413e-02 -8.99915824e-04 ... 4.32169233e-02
      4.21049371e-02 1.64754571e-01]
     [ 3.88645000e-05 -1.74279859e-02 -8.89844263e-04 ... -1.75505225e-01
      -1.25519068e-01 1.00027540e-02]
     [ 3.32174398e-05 -1.48957226e-02 -7.53627895e-04 ... 7.36611090e-02
      5.94190503e-02 1.88707634e-01]
     ...
     [ 1.66081594e-01 7.68328498e-02 -7.03178257e-02 ... -9.60287711e-01
      2.83796643e-02 1.56399942e-02]
     [ 1.25156666e-01 4.30850083e-02 -5.95435180e-02 ... 2.83465548e-02
      -9.78944404e-01 1.67258531e-02]
     [ 9.15455888e-02 -1.47110341e-01 -1.94199041e-01 ... 1.65308022e-02
      1.66469061e-02 -9.31105108e-01]]]
Sigma: [3.23575983e+06 1.74397555e+05 6.86024479e+03 1.54256782e+02]
V*: [[ 4.13398236e-02 9.99143699e-01 1.69918151e-03 -1.90918876e-05]
      [-9.99143144e-01 4.13363114e-02 2.04829346e-03 -3.02725394e-04]
      [-1.97746978e-03 1.78251192e-03 -9.99988764e-01 3.92234827e-03]
      [-2.93922696e-04 2.45976481e-05 3.92298707e-03 9.99992262e-01]]

```

Fig. 5. SVD on original data

SVD on Z-scored standardized data:

```

SVD
U: [[ 0.06132141 0.02284989 -0.01668461 ... 0.01864612 0.0138139
      -0.17797746]
     [ 0.06246649 0.02289292 -0.01710579 ... -0.17109863 -0.12027931
      0.00577292]
     [ 0.06797651 0.02496603 -0.02031495 ... -0.06176361 -0.05318097
      -0.21918571]
     ...
     [-0.13339922 0.09919886 -0.05945816 ... -0.96444296 0.02338664
      0.01509854]
     [-0.10608223 0.07437359 -0.03232511 ... 0.02294347 -0.98285746
      0.00591266]
     [-0.15327085 -0.04132412 0.20715114 ... 0.01446281 0.00600947
      -0.90531603]]]
Sigma: [24.07137206 15.02023728 9.46385723 8.08683657]
V*: [[ -0.35018343 -0.55933789 -0.47477107 0.58232734]
      [-0.82161262 -0.00452564 0.56898115 -0.03453503]
      [ 0.37607837 -0.75509623 0.53314117 -0.06446119]
      [ 0.24675763 0.34198002 0.40818302 0.80965855]]

```

Fig. 6. SVD on Z-scored standardized data

SVD on range standardized data:

```

SVD
U: [[ 0.06566131  0.02565575  0.02068869 ...  0.01864612  0.0138139
      -0.17797746]
     [ 0.06735571  0.02584818  0.02427366 ... -0.17109863 -0.12027931
      0.00577292]
     [ 0.07515804  0.02852901  0.03827494 ... -0.06176361 -0.05318097
      -0.21918571]
     ...
     [-0.12373692  0.09462852  0.03597184 ... -0.96444296  0.02338664
      0.01509854]
     [-0.10530785  0.06975847 -0.0080703 ...  0.02294347 -0.98285746
      0.00591266]
     [-0.14378221 -0.03725076  0.25980513 ...  0.01446281  0.00600947
      -0.90531603]]
Sigma: [6.60736293 3.83106847 2.4030212 2.09064323]
V*: [[-0.29874304 -0.39982598 -0.46366863 0.73205408]
      [-0.81632001 -0.03596368 0.5763485 0.01227448]
      [0.46959659 -0.01569408 0.65145473 0.59568404]
      [-0.1544451 0.91575078 -0.16864323 0.330313 ]]

```

Fig. 7. SVD on range standardized data

SVD on rank standardized data:

```

SVD
U: [[ 0.08623102  0.03685018 -0.03085358 ...  0.0370722  0.03770005
      0.1709604 ]
     [ 0.08767107  0.03759248 -0.03159861 ... -0.17877728 -0.12768141
      0.01040809]
     [ 0.09371058  0.04164006 -0.03644315 ...  0.07046201  0.05906135
      0.18897713]
     ...
     [ 0.02093518 -0.1395772 -0.11346408 ... -0.96095803  0.02867502
      0.01436116]
     [ 0.01696782 -0.11252599 -0.07887478 ...  0.02867547 -0.97867616
      0.01492314]
     [ 0.05139389 -0.19891292 -0.00734302 ...  0.01461407  0.01502313
      -0.92038504]]
Sigma: [9.73832479 6.05617744 3.77203748 2.14904906]
V*: [[ 0.29448451  0.09911757  0.07199447  0.9477718 ]
      [-0.56235793 -0.52845282 -0.57413342 0.27360921]
      [0.69253414 -0.05494345 -0.7021478 -0.1560966 ]
      [0.342679 -0.84136469 0.41494106 -0.05000475]]

```

Fig. 8. SVD on rank standardized data

The goal of SVD is to find contributions of each component. In the following table you can see natural contribution of each feature in original dataset and in standardized datasets:

Table 13. Natural contribution of each feature

	Temperature	Luminosity	Radius	Absolute magnitude
Original data	1.0407e+13	3.0415e+10	4.7063e+07	2.3752e+04
Ranked std data	94.835	36.6773	14.2283	4.6184
Ranged std data	43.6573	14.677	5.7745	4.3708
Z-Scored std data	579.431	225.6075	89.5646	65.3969

It shows absolute value of contribution.

In the following table you can see percent contribution of each feature in datasets:

Table 14. Percent contribution of each feature

	Temperature	Luminosity	Radius	Absolute magnitude
Original data	99.71	0.2897	0.000448	0.000000227
Ranked std data	63.07	24.39	9.46	3.072
Ranged std data	63.75	21.43	8.43	6.38
Z-Scored std data	60.36	23.5	9.33	6.81

The first component («Temperature») contributes the most in each dataset.

In the following table you can see data scatter of each dataset:

Table 15. Data scatter

Data Scatter	
Original data	10 500 603 295 793.56
Ranked std data	150.359
Ranged std data	68.48

Range standardized data has the least data scatter. Original data has the biggest data scatter because it is not standardized at all, each feature has values of different orders of magnitude.

Hidden ranking factor

Hidden factor is a value which characterize an object according to original features. Hidden factor is a value, which is highly correlated with original features.

Hidden ranking factor is rank standardized hidden factor. It has value from 0 to 100.

In the following tables you can see stars with top 5 hidden ranking factors for original data.

Table 16. Stars with top 5 hidden ranking factors for original data

	Temperature	Luminosity	Radius	Absolute magnitude	Star type	Star color	Spectral Class
233	27 739	849 420.0	1 252.0	-7.590	5	Blue-white	B
236	30 839	834 042.0	1 194.0	-10.630	5	Blue	O
101	40 000	813 000.0	14.0	-6.230	4	Blue	O
227	10 930	783 930.0	25.0	-6.224	4	Blue	O
229	21 738	748 890.0	92.0	-7.346	4	Blue	O

Their hidden ranking factors are in the following table.

Table 17. Top 5 hidden ranking factors for original data

Hidden ranking factor	
233	1
236	0.982
101	0.958

227	0.922
229	0.882

Contribution of this hidden factor is 99.7%

In the following tables you can see stars with top 5 hidden ranking factors for Z-score standardized data.

Table 18. Stars with top 5 hidden ranking factors for Z-score standardized data

	Temperature	Luminosity	Radius	Absolute magnitude	Star type	Star color	Spectral Class
4	1939	0.000138	0.1030	20.06	0	Red	M
122	3218	0.000540	0.1100	20.02	0	Red	M
128	2856	0.000896	0.0782	19.56	0	Red	M
188	2778	0.000849	0.1120	19.45	0	Red	M
121	3531	0.000930	0.0976	19.94	0	Red	M

Their hidden ranking factors are in the following table.

Table 19. Top 5 hidden ranking factors for Z-score standardized data

Hidden ranking factor	
4	1
122	0.992
128	0.9908
188	0.9903
121	0.9902

Contribution of this hidden factor is 60.36%

In the following tables you can see stars with top 5 hidden ranking factors for range standardized data.

Table 20. Stars with top 5 hidden ranking factors for range standardized data

	Temperature	Luminosity	Radius	Absolute magnitude	Star type	Star color	Spectral Class
4	1939	0.000138	0.1030	20.06	0	Red	M
122	3218	0.000540	0.1100	20.02	0	Red	M
121	3531	0.000930	0.0976	19.94	0	Red	M
128	2856	0.000896	0.0782	19.56	0	Red	M
125	3225	0.000760	0.1210	19.63	0	Red	M

Their hidden ranking factors are in the following table.

Table 21. Top 5 hidden ranking factors for range standardized data

Hidden ranking factor	
4	1
122	0.994
121	0.991
128	0.989
4	1

Contribution of this hidden factor is 63.75%

In the following tables you can see stars with top 5 hidden ranking factors for rank standardized data.

Table 22. Stars with top 5 hidden ranking factors for rank standardized data

	Temperature	Luminosity	Radius	Absolute magnitude	Star type	Star color	Spectral Class
121	3531	0.000930	0.0976	19.94	0	Red	M
122	3218	0.000540	0.1100	20.02	0	Red	M
4	1939	0.000138	0.1030	20.06	0	Red	M
125	3225	0.000760	0.1210	19.63	0	Red	M
128	2856	0.000896	0.0782	19.56	0	Red	M

Their hidden ranking factors are in the following table.

Table 23. Top 5 hidden ranking factors for rank standardized data

Hidden ranking factor	
121	1
122	0.9999
4	0.9898
125	0.9866
128	0.9809

Contribution of this hidden factor is 63.07%

Ranks of the same object from different dataset are not equal. Hidden factors of standardized datasets are quite similar. Objects of M class, red color, and 0 star type has the highest hidden factors. In the original dataset stars of other types have the highest hidden factors.

Data visualization

Using two first components (principal components) we made data visualization in coordinates of these two components.

Objects of different colors belong to different clusters, labels of clusters are mentioned in graph legend.

In the following graphs you can see visualizations of PCAs of Z-scored standardized data and range standardized data:

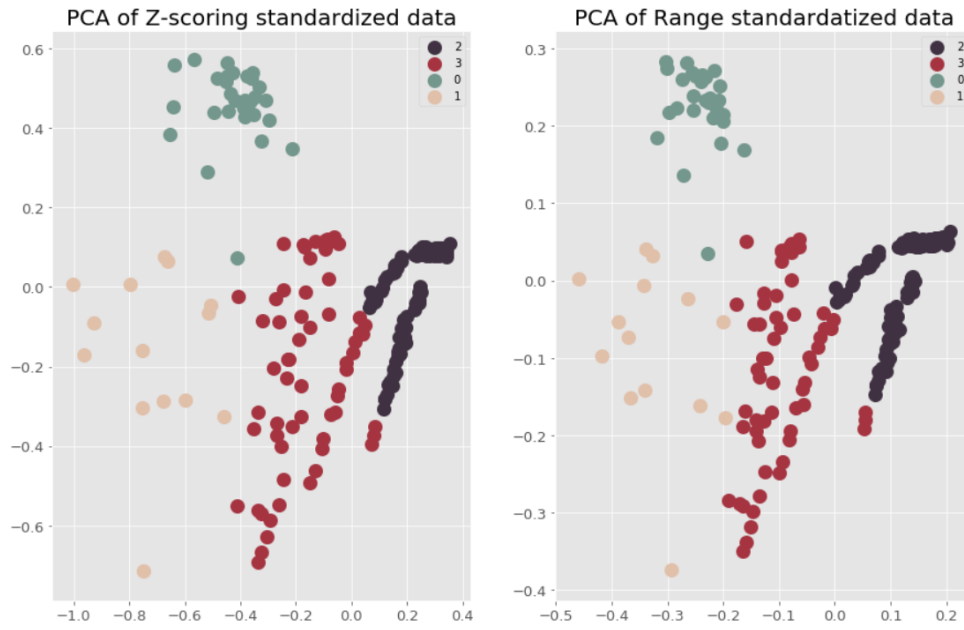


Fig. 9. PCAs of Z-scored standardized (left) and range standardized (right) data

Visualizations are quite similar, they just differ a bit in scale. Clusters, which were obtained by K-Means are well separated in both graphs. Also, PCAs have different axes scales.

After that we applied conventional PCA for original data and for Z-score standardized data. Also we applied PCA from sklearn.decomposition module.

Now let's compare the results.

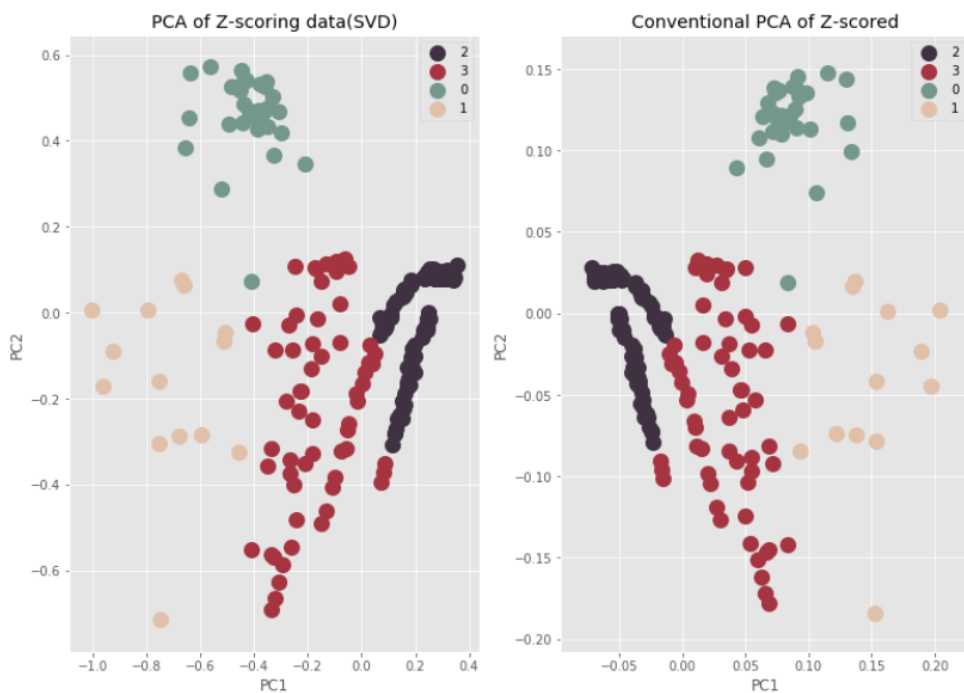


Fig. 10. PCA (left) and conventional PCA (right) of Z-scored standardized data

Graphs are mirrored and have different axes scales. Clusters are well separated in both graphs.

In the following graphs you can see visualizations of PCAs of range standardized data and original dataset.

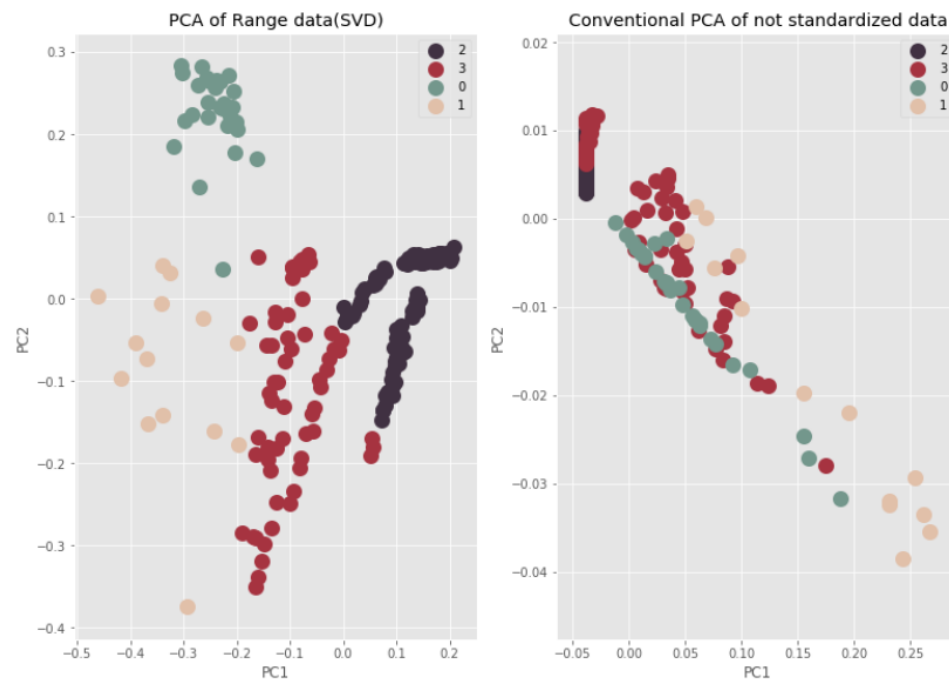


Fig. 11. PCAs of range standardized (left) and original (right) data

Visualization of PCA of range standardized data is quite similar with previous visualizations. Clusters are well separated.

Visualization of PCA of original data differs. Clusters are not separated; this visualization is not interpretable. This fact means that standardization is obligatory for clustering and other methods of data analysis.

In the following graph, you can see sklearn PCA visualization:

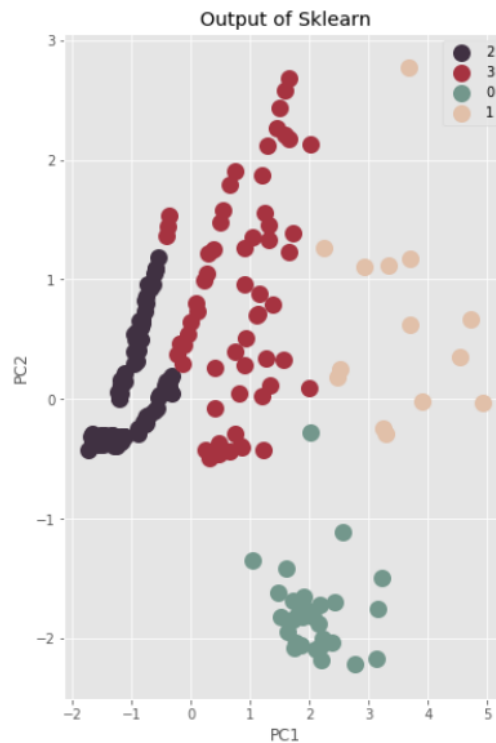


Fig. 12. sklearn PCA

The graph is mirrored visualization of PCA of range standardized data.

It is hard to say exactly which method of standardization is better, because they have very similar results. Z-score standardization can be easily interpreted. If the value equals 0, it means that the value is the mean. If the value is less or bigger than 0, it means that value differs from the mean, and absolute value shows how much and in what side it differs from the mean. In our point of view Z-scoring is better in the context of interpreting values in a dataset. But in fact, it doesn't influence on results greatly.

Correlation Coefficient

Before evaluating a correlation coefficients, 2 features of the “Star dataset” with the most similar to linear scatterplot were found. Scatterplots of all quantitative features (“Temperature”, “Luminosity”, “Radius”, “Absolute magnitude”) are shown on the graph below.

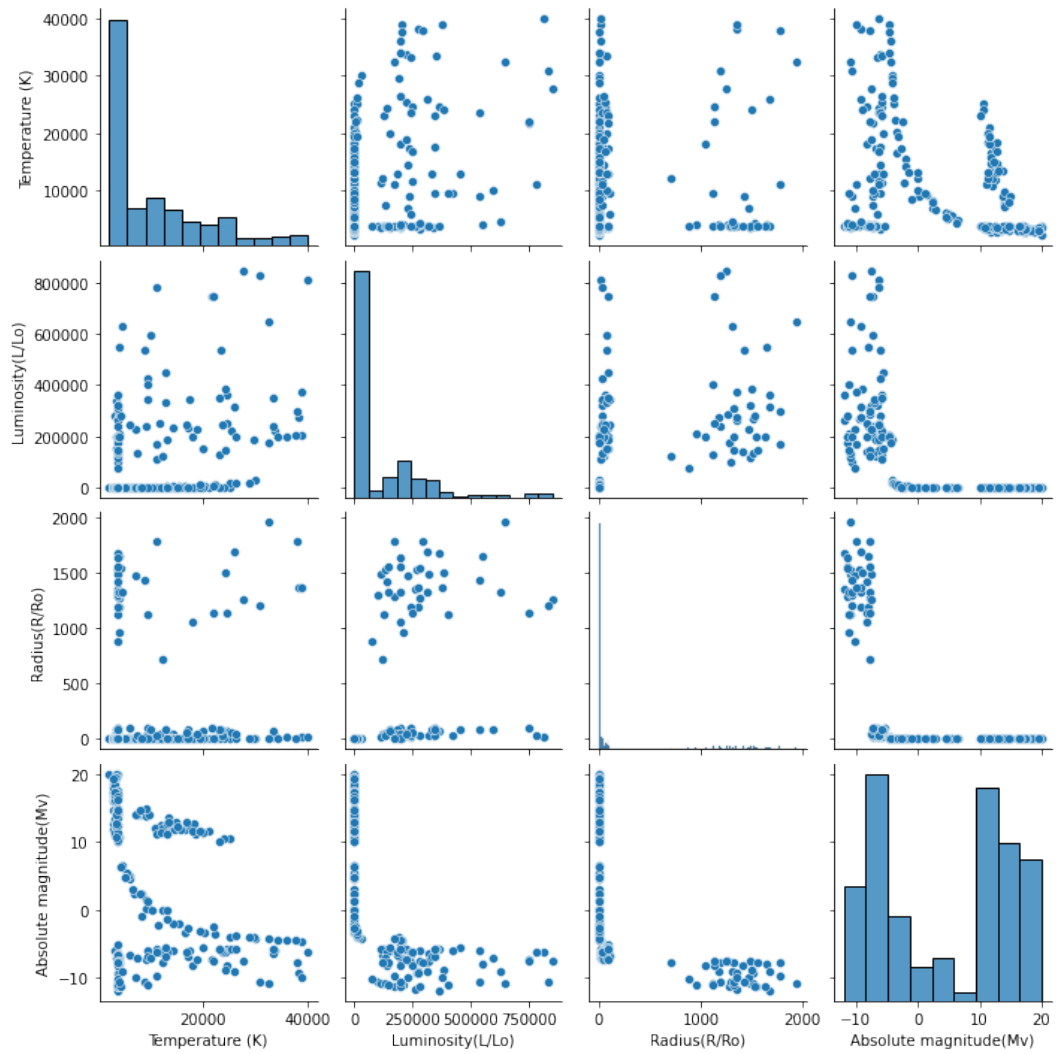


Fig. 13. Scatterplots of the quantitative features

According to the visual analysis of the graph, “Temperature” and “Absolute magnitude” have more linear-like scatterplot than other pairs.

On the graph below the scatterplot for selected pair is shown.

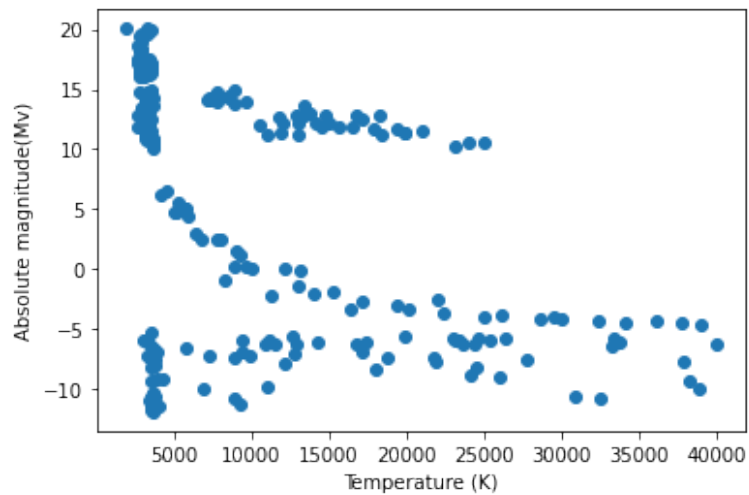


Fig. 14. Scatterplot for “Temperature” and “Absolute magnitude”

There is a linear trend called main sequence and two clusters above and below the line. Linear regression is suitable, since a slender linear-like dependency is visible. Presumably, the quality of the linear regression will not be high, as the main sequence looks a little hyperbolic on the left side, and two clusters above and below the line will not be reflected the linear model and will increase the error significantly.

Linear regression

After fitting of a linear regression, the following equation is obtained:

$$Y = -0.0004634X + 9.247$$

where Y is an “Absolute magnitude” and X is a “Temperature”.

On the following graph, there is a straight for this equation.

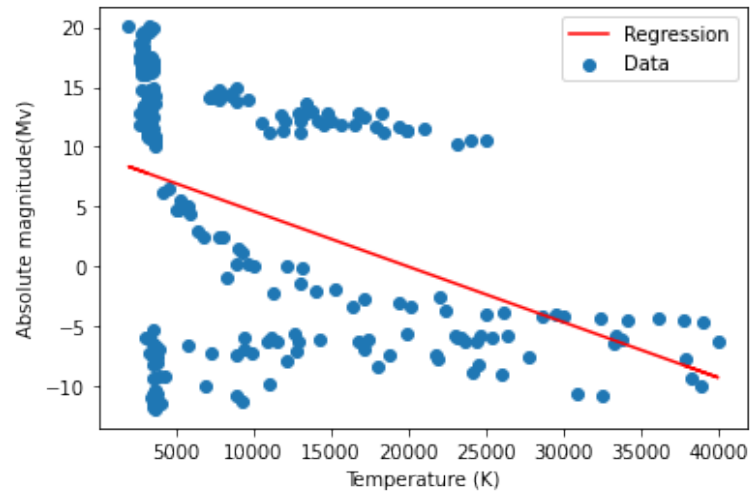


Fig. 15. Linear regression for “Temperature” and “Absolute magnitude”

It can be seen, that the slope is negative, as it should be for the main sequence of stars. It means that the absolute magnitude negatively depends on the temperature of the star. The greater the temperature, the lower the absolute magnitude.

Correlation coefficients

Correlation coefficients for all four features (“Temperature”, “Luminosity”, “Radius”, “Absolute magnitude”) were calculated and represented as a heatmap of the correlation matrix, which is on the following graph.

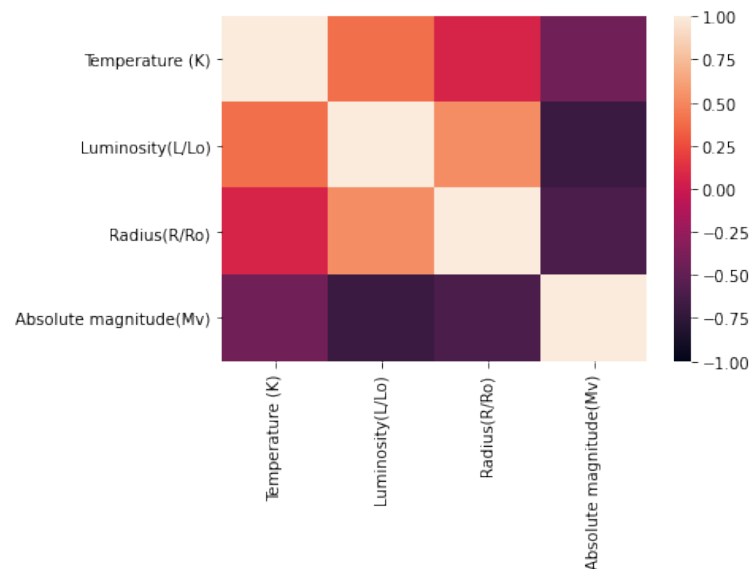


Fig. 16. Heatmap of correlation coefficients

Correlation coefficient for features “Temperature” and “Absolute magnitude” is -0.4203, which means that the correlation between them is quite weak and negative.

The strongest positive correlation is between “Luminosity” and “Radius”. The strongest negative correlation is between “Luminosity” and “Absolute magnitude”. The weakest correlation is observed for “Temperature” and “Radius”.

Determinacy coefficients

On the following graph the heatmap for determinacy coefficients matrix is presented.

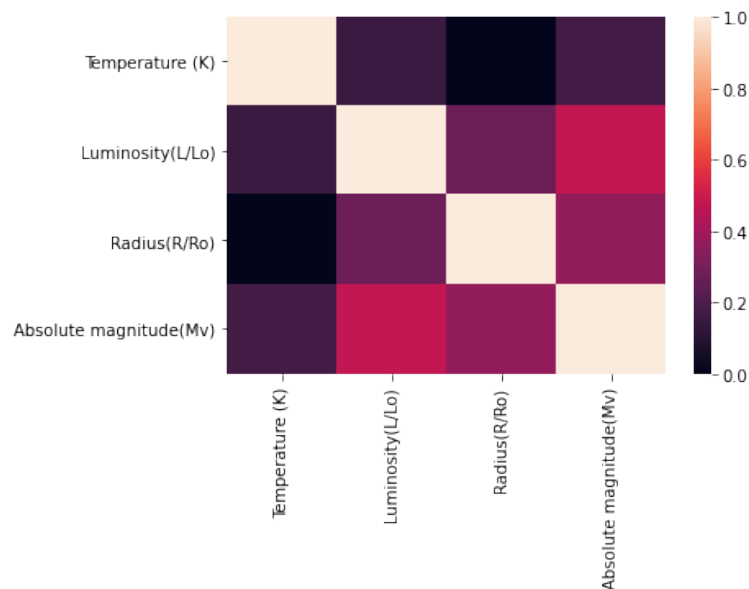


Fig. 17. Heatmap of determinacy coefficients

Determinacy coefficient for the features “Temperature” and “Absolute magnitude” is $R^2 = 0.1766$, which means that the variables explain the low proportion of variance of each other. This result is consistent with the presumption made before.

As known, the closer the value of the coefficient to 0, the weaker the dependence between the variables under consideration. In this case more than 80% of variance is unexplained because of two clusters of stars, stand out from the general pattern. The cluster above the line of the regression is the cluster of giants, the cluster below the line consists of white dwarfs. The stars in these clusters are characterized by extremely low dependence of absolute magnitude on temperature.

Upon the whole, “Absolute magnitude” and “Luminosity” have the highest determinacy coefficient, while “Temperature” and “Radius” have the lowest one.

Predictions

Let’s predict the absolute magnitude of the star by its temperature using estimated linear model.

We have three values of the temperature: 1 000, 25 000, 35000 Kelvins. After applying the linear equation, the result, presented on a graph, was obtained.

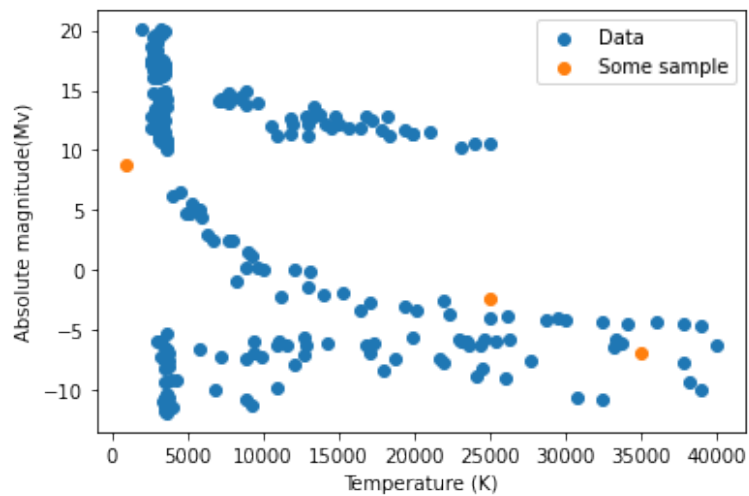


Fig. 18. Scatterplot with the predicted values

The model predicts the absolute magnitude by temperature quite well for the main sequence stars, but the prediction will always be wrong if the star is from another cluster (for example giants above the main sequence or white dwarfs below this line).

The mean relative absolute error MRAE of the model is 354.95 on the set of all points from the data. Such a big error is due to the fact that the data contains other types of stars in addition to the main sequence stars.