

Analysis of anime voice actor's social network and popularity.

Florencia Zanollo. Hideaki Takeda.

July 13, 2018

Contents

Introduction	1
1 Anime/Seiyuu Dataset	2
1.1 Wikidata	2
1.2 MyAnimeList	2
1.3 Data retrieved	3
2 Seiyuu Social Network	4
2.1 Node and edge definitions	4
2.2 Construction	4
2.3 Analysis	5
2.4 Conclusion	8
3 Analysis and prediction of seiyuu popularity	10
3.1 Correlation with only one feature	12
3.1.1 Why last 9 years of works has more correlation?	14
3.2 Correlation with multiple features	17
3.2.1 Only one category	19
3.2.2 Groups of two categories	20
3.2.3 Groups of three categories	21
3.2.4 All categories	22
3.3 Conclusion	23
Conclusion	24

Abstract

Although Social Network of actors is a relatively common object of investigation that has been addressed many times, we can say that seiyuu (anime voice actors) come from a very different industry with a distinct way to relate to each other. In this research we use Wikidata and MyAnimeList to collect seiyuu information and build a Social Network. Topics explored:

- Structure and characteristics of seiyuu Social Network.
- Understanding what properties have a main role describing and predicting popularity of seiyuu.
- Compare prediction performances between different machine learning algorithms (and different models).

Introduction

TODO

Previous work, including any relevant information learned during literature review (NND, books, something else?)

Social networks, small explanation

Anime/Seiyuu industry

Link open data, wikidata y my anime list

Chapter 1

Anime/Seiyuu Dataset

Anime –like any other animation projects– have voice actors to play roles of each character. They usually have multiple seasons or adaptations based on original content –which can be manga, games, visual novels, etc– and it’s not uncommon for a character to have always the same actor voicing it.

Since some time ago anime industry is growing bigger each year and so is seiyuu industry. Seiyuu can become very popular, with a great international fanbase and work in different areas other than voice acting, for example as singers, on theaters, etc.

1.1 Wikidata

Wikidata¹ is a collaboratively edited knowledge base intended to provide a common source of data which can be used by Wikimedia projects such as Wikipedia. The information is stored in RDF format, and can be retrieved in multiple ways, one of them being through a SPARQL endpoint.

Using Wikidata’s SPARQL endpoint we retrieved a list of seiyuu. This list contains all persons that have seiyuu as occupation, a total of 6472 entities were obtained². Gender, birthday and birthplace information was also fetched (last two were not used in the end because it was lacking in the majority of entities).

1.2 MyAnimeList

Wikidata information about seiyuu’s works is really incomplete that’s why MyAnimeList (MAL)³ was used to retrieve voice acting roles and anime information. MAL is a social networking and social cataloging application website with a large database on anime and manga that started in April 6, 2006. Users can make a list of currently watching, watched and/or favorite anime; score, review, comment and recommend similar ones. They can also comment about and favorite people working on the industry (voice actors, directors, editors, etc).

¹<http://wikidata.org/>

²There’s actually 7030 seiyuu in Wikidata but only 6472 of them have an English label (name)

³<https://myanimelist.net/>

Since only 59 of Wikidata’s seiyuu entities had MyAnimeList ID (MALID) property, a matching between Wikidata and MyAnimeList was done using seiyuu’s complete name to retrieve the ID for those who was missing. Successfully restoring 3033 MALIDs, giving a total of 3092 seiyuus with that property; 2956 of them having at least one work according to MAL so we are using this subset for our experiments.

Using Jikan API⁴ and MALID, seiyuu data, voice acting roles and more information about each anime was retrieved.

An issue to take into account is whether we unify all anime adaptations of the same intellectual property as one or take a single adaptation as a independent work. We chose the later because each adaptation has its own producer, score, popularity, among other information; it would be incorrect to say a seiyuu worked in a popular work when that adaptation didn’t have enough fame.

1.3 Data retrieved

All in all we were able to retrieve the following information for 2956 seiyuu and 7614 anime.

- For Seiyuu:
 - Name
 - Debut (this was obtained from oldest work aired date)
 - Gender
 - Popularity (member_favorites information of MAL)
 - Work (anime roles with anime information plus wheter is a main role or not)
- For Works (Anime):
 - Year that began airing
 - Favorites
 - Score (from 0 to 10, MAL user based)
 - Popularity (ranking over all MAL animes)
 - Members (how many MAL users have it on their list)
 - Genres

It’s important to notice that data such as popularity and scores are retrieved from MAL, which is user review based only; it may differ with actual awards winning or professional reviewing of works.

Further, this dataset is biased in favor of more recent anime and seiyuu, since it accounts for more complete data and with better quality. Oldest anime in this dataset is from 1960 having no record about previous ones. Majority of seiyuu’s debut are from 1988 which leads us to think information from thereon is more complete.

The data was stored using Virtuoso server to create a local SPARQL endpoint, mongodb was also used as an intermediate storage (before formatting data as RDF).

⁴<https://jikan.docs.apiary.io/#>

Chapter 2

Seiyuu Social Network

Social networks consist of a finite set of actors and the relations between them. Usually represented as a graph; with actors or organizations as set of nodes and a defined relation between them as set of edges. This structures are useful to analyze complex social interactions and communities.

2.1 Node and edge definitions

This social network is of a particular kind called *two-mode networks* which consists of a set of actors (seiyuu) and events (anime). So there exists two ways of viewing it, one will be from seiyuu perspective, using anime in common for edges; the other being from anime perspective, using seiyuu in common for edges. We chose the former since we found more interesting they being actual people and using other information about them such as debut and gender.

So our social network consists of voice actors as nodes and co-workship between them as edges. It's important to notice that this social network is time dependant since each seiyuu has a debut year and each anime has an aired time; giving us freedom to choose different time frames to observe.

Aside from being time dependant there exists different possible definitions of relationship or co-workship between seiyuu. One could say two actors know each other if they have worked in at least one job together, or maybe it requires more than one. There's also a time frame to define, relationship could take into account all works of both of them or only from certain years.

2.2 Construction

As a first approach Gephi was used to build the network. Since the graph was big enough to bring performance problems and we needed to build the edges dynamically (which couldn't be done in Gephi) NetworkX was used instead.

NetworkX was chosen because it's an easy yet powerful Python library, it doesn't get along with massive graphs but ours was not big enough to present a problem.

One can export the graph and open it on Gephi, for a more visual analysis.

And also we needed to build the edges dynamically because according to our definition they depend on the time frame we are looking at. For example, for at least 10 works in common, if two actors worked together in 9 jobs between 1960 and 1970 we shouldn't see an edge between them; but if they worked together again in 1971 then looking at 1960-1971 they should be connected.

2.3 Analysis

In this section we are going to compare and analyze two definitions of relationship for our social network in order to understand more about its structure and decide on a definition:

- at least 1 work in common
- at least 10 works in common

Both of them during the time frame between the first debut registered (1960) and the year of observation.

It is easy to tell at first glance that this social network is really interconnected. With merely 2956 nodes it has 395887 edges when only one work in common is required and 13629 edges when asking for 10 or more. It shows a tightly interconnected cluster surrounded by poorly or not connected nodes. This cluster represents 99% of the nodes of one work in common graph and 23% of 10 works in common. In terms of modularity we can see at least four clear communities in each graph, Fig 2.1 and Fig 2.2.

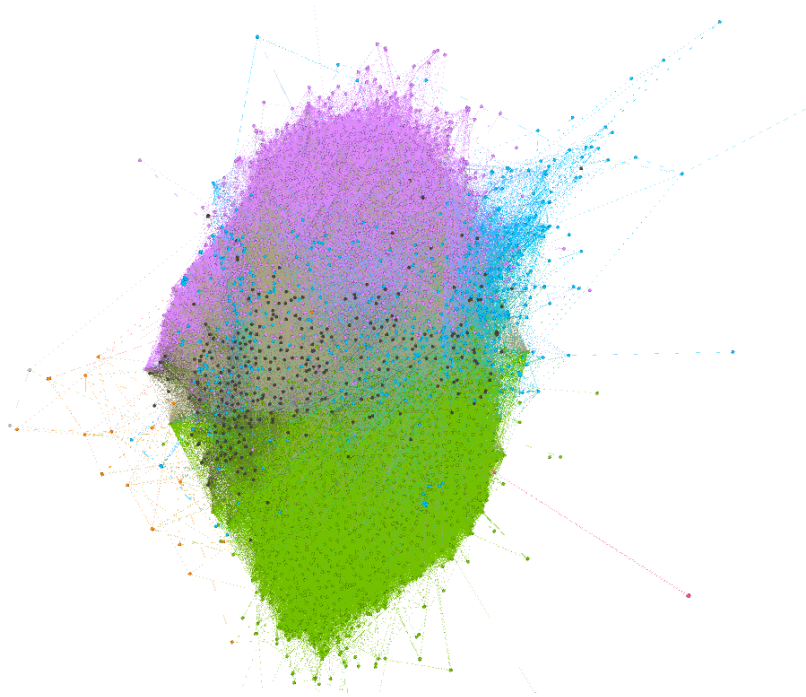


Figure 2.1: At least one work in common graph coloured by community.

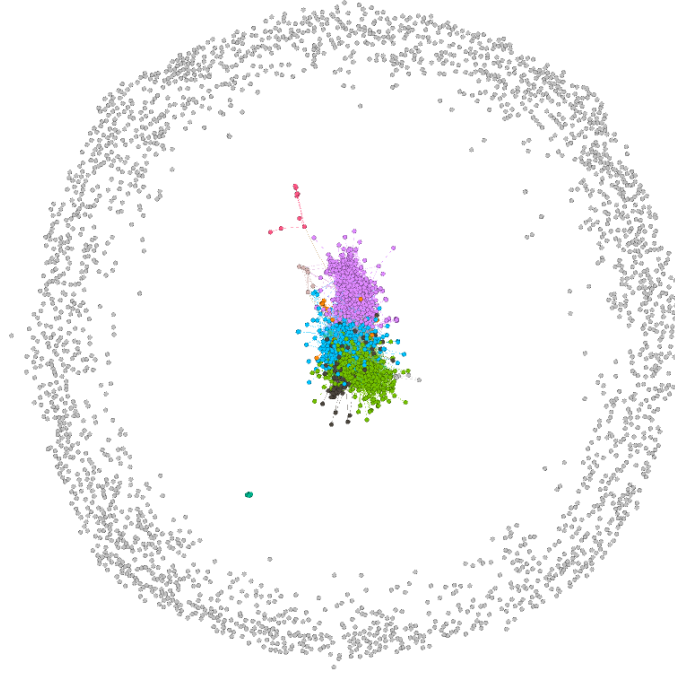


Figure 2.2: At least ten works in common graph coloured by community. Big cluster at the center, surrounded by loosely connected nodes.

Table 2.1 shows metrics about each graph. Requiring more works in common decreases average degree circumstantially but doesn't change much modularity or network diameter.

Table 2.1: Graph analysis

Graph	Avg Degree	Graph Density	Modularity
One work in common	267	0.09	0.2
Ten works in common	9	0.003	0.29

Graph	Network Diameter	Connected Components
One work in common	6	18
Ten works in common	7	2261

As proven by Fig. 2.3 growth of edges by year follows a similar distribution regardless of how many works in common are used to build the social network.

Fig. 2.4 shows that more than half of the nodes are from last 18 years (2000 to 2018), giving us an idea of how much seiyuu industry is growing.

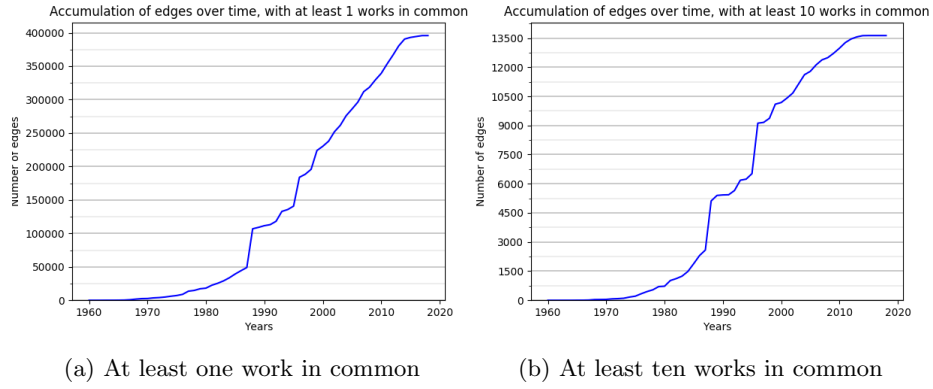


Figure 2.3: Growth of edges over time. For 1 and 10 works in common graphs

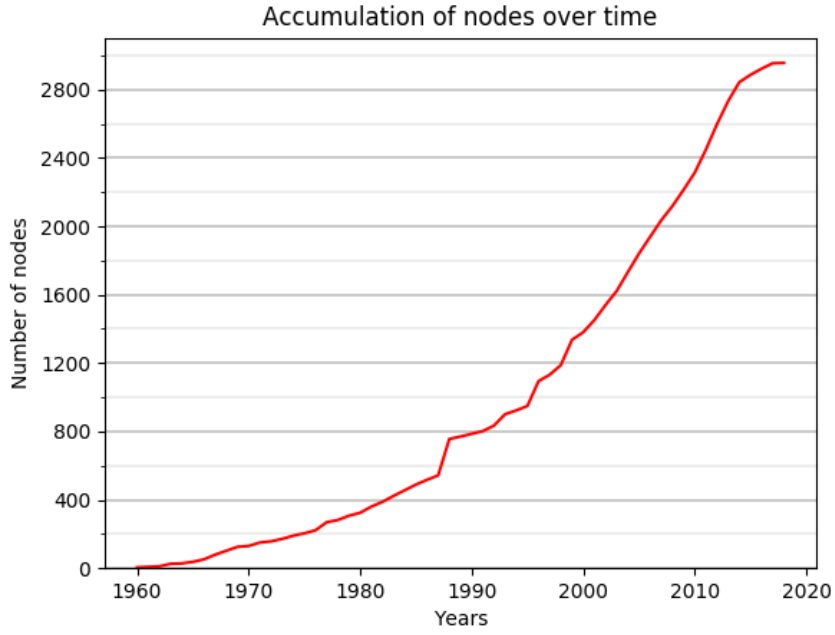


Figure 2.4: Growth of nodes over time.

Table 2.4 shows top 10 nodes, for degree and betweenness centrality for "at least 1 work in common" definition. And Table 2.7 does the same for "at least 10 works in common".

Name	Degree
Takehito Koyasu	1545
Akira Ishida	1488
Mamiko Noto	1422
Nobuo Tobita	1417
Daisuke Namikawa	1390
Nobuyuki Hiyama	1358
Rikiya Koyama	1331
Jrta Kosugi	1322
Keiji Fujiwara	1312
Kazuhiko Inoue	1309

Table 2.2: Top 10 degree

Name	Betweenness Centrality
Takehito Koyasu	49982.52
Akira Ishida	40221.50
Daisuke Namikawa	30448.43
Nobuo Tobita	29363.25
Mamiko Noto	29168.18
Rie Kugimiya	29122.31
Miyuki Sawashiro	28997.40
Kazuhiko Inoue	27693.88
Daisuke Ono	27034.592
Keiji Fujiwara	26802.69

Table 2.3: Top 10 Betweenness centrality

Table 2.4: At least one work in common

Name	Degree
Takehito Koyasu	311
Akira Ishida	273
Mamiko Noto	258
Daisuke Namikawa	232
Katsuyuki Konishi	229
Keiji Fujiwara	220
Junichi Suwabe	216
Toshiyuki Morikawa	215
Rie Kugimiya	213
Nobuyuki Hiyama	201

Table 2.5: Top 10 degree

Name	Betweenness Centrality
Takehito Koyasu	18489.44
Mamiko Noto	10988.96
Daisuke Namikawa	9570.48
Akira Ishida	8299.19
Rie Kugimiya	7560.16
Katsuyuki Konishi	7413.72
Kenichi Ogata	7160.54
Harumi Sakurai	6775.76
Keiji Fujiwara	5980.58
Yoshimasa Hosoya	5607.95

Table 2.6: Top 10 Betweenness centrality

Table 2.7: At least ten works in common

2.4 Conclusion

Both networks have fairly similar top 10s so it points to them having similar structure and connections among their nodes, aside from actual values.

From now on our definition for edges will be: *at least 10 works in common, during the time frame between the first debut registered (1960) and the year of observation*. Because requiring more jobs in common means less amount of edges, this leaves a more understandable graph and we verified it does without changing its structure so much.

There's also other interesting definitions of relationship, for example we can use only common works from the last x years or from all time. This options weren't explored; having into account our limited time.

Table 2.8 shows a little more information about seiyuu that appear on top 10s.

Table 2.8: More information about seiyuu

Name	Popularity	Debut	Gender	Birthyear
Takehito Koyasu	7235	1988	Male	1967
Akira Ishida	7612	1989	Male	1967
Mamiko Noto	7544	1988	Female	1980
Daisuke Namikawa	8304	1988	Male	1976
Katsuyuki Konishi	3702	1996	Male	1973
Keiji Fujiwara	2778	1986	Male	1964
Junichi Suwabe	10838	1996	Male	1972
Toshiyuki Morikawa	2455	1981	Male	1967
Rie Kugimiya	31668	1996	Female	1979
Jun Fukuyama	26811	1981	Male	1978
Kenichi Ogata	52	1974	Male	1942
Harumi Sakurai	341	2005	Female	1982
Yoshimasa Hosoya	4852	2006	Male	1982
Nobuo Tobita	139	1981	Male	1959
Nobuyuki Hiyama	1723	1988	Male	1967
Rikiya Koyama	2919	1996	Male	1963
Jrta Kosugi	114	1985	Male	1957
Miyuki Sawashiro	26501	1988	Female	1985
Kazuhiko Inoue	2445	1974	Male	1954
Daisuke Ono	24080	1996	Male	1978

Chapter 3

Analysis and prediction of seiyuu popularity

Popularity is an abstract criterion that must be defined as a numerical metric in order to be used for analysis and prediction. Since we are using MAL database and it has a social component, seems logic to use member_favorites as a representation of popularity. We can also get popularity and score of anime from opinions of the same set of users.

In terms of distribution *popularity* is highly unequal –as we can observe in Fig. 3.1– having a lot of seiyuu which are no member favorites and only a few who are favorite of more than 10000 members. It's good to keep in mind that users can favorite multiple seiyuu.

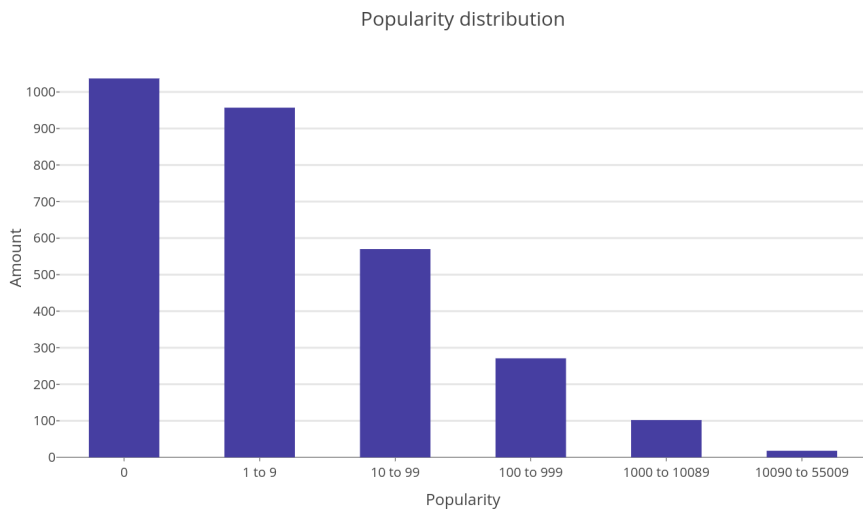


Figure 3.1: Amount of seiyuu with that popularity, divided into groups for better visualization.

Some metrics about popularity:

- Mean: 289.55
- Median: 2.0
- Max: 55018
- Min: 0
- 1037 values equal to zero
- Only 120 values bigger than 1000

TODO AGREGAR ALGUNOS DE LOS ROLES MS FAMOSOS DE CADA UNO ;)

Table 3.1: Top 10 popular seiyuu

Name	Popularity
Kana Hanazawa	56637
Hiroshi Kamiya	49685
Mamoru Miyano	43942
Rie Kugimiya	31668
Jun Fukuyama	26811
Miyuki Sawashiro	26501
Tomokazu Sugita	24449
Daisuke Ono	24080
Saori Hayami	18322
Aya Hirano	18094

3.1 Correlation with only one feature

Our first approach to explaining popularity was using Pearson correlation.

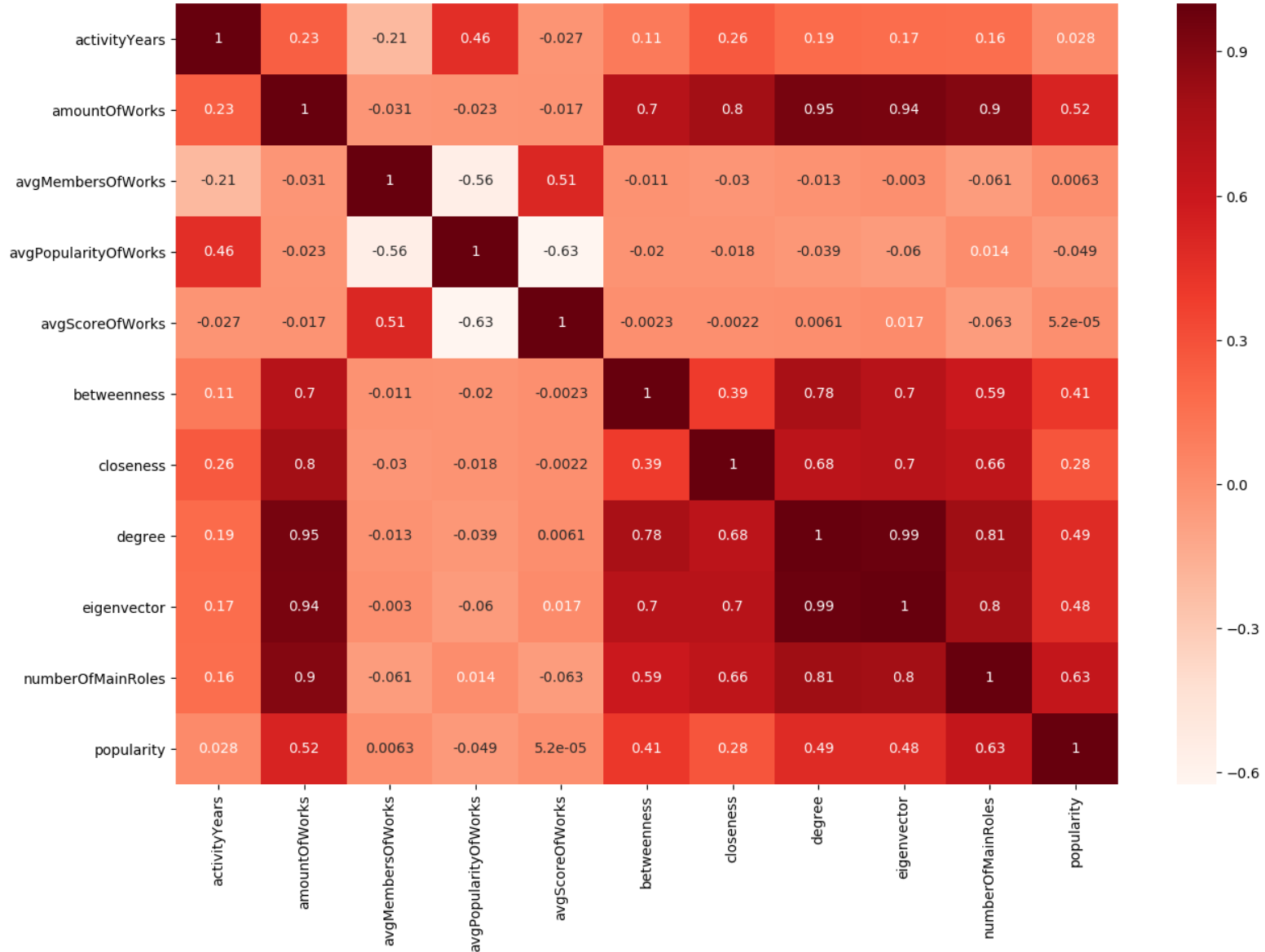


Figure 3.2: Pearson correlation between popularity and attribute of nodes (using all works).

As shown on Fig. 3.2 a fairly big correlation can be seen between popularity and amount of works. This attribute doesn't have the biggest correlation with popularity but "number of main roles" was added to the end of this investigation since we didn't had the data for doing so before. Number of main roles and

amount of works have a strong correlation with each other (0.9) but they have different influence over popularity, this means they provide distinct information.

We are showing only average of values for work's attributes (ex. favorites) mostly for better visualization but for predictions we use sum, mean, median and maximum.

Since our dataset is biased in favor of more modern anime we thought of correlate with more recent works only. But, how recent? Last 5, 10 or 20 years? Thus correlation between popularity and works from different data frames was analyzed, Fig. 3.3.

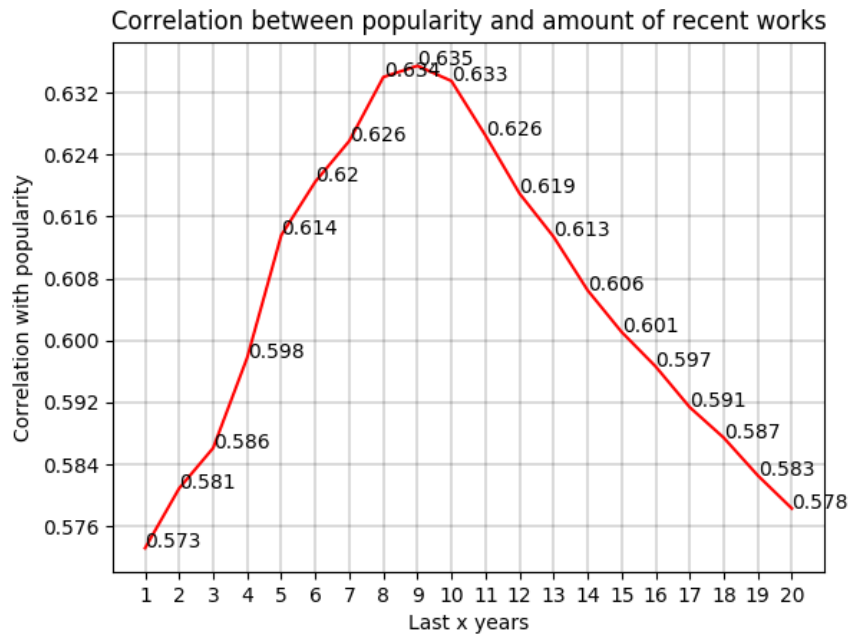


Figure 3.3: Last X years means works from 2018- X to present.

The best result was given by recent works from last 9 years. Therefore, this definition of recent works was used from there on.

Fig. 3.4 shows the result of running Pearson again but using information from last 9 years of works only. It's important to clarify that we didn't build the network using only last 9 years, so betweenness centrality and degree are exactly the same as before. We left "amount of works" attribute for easy comparison against "amount of recent works".

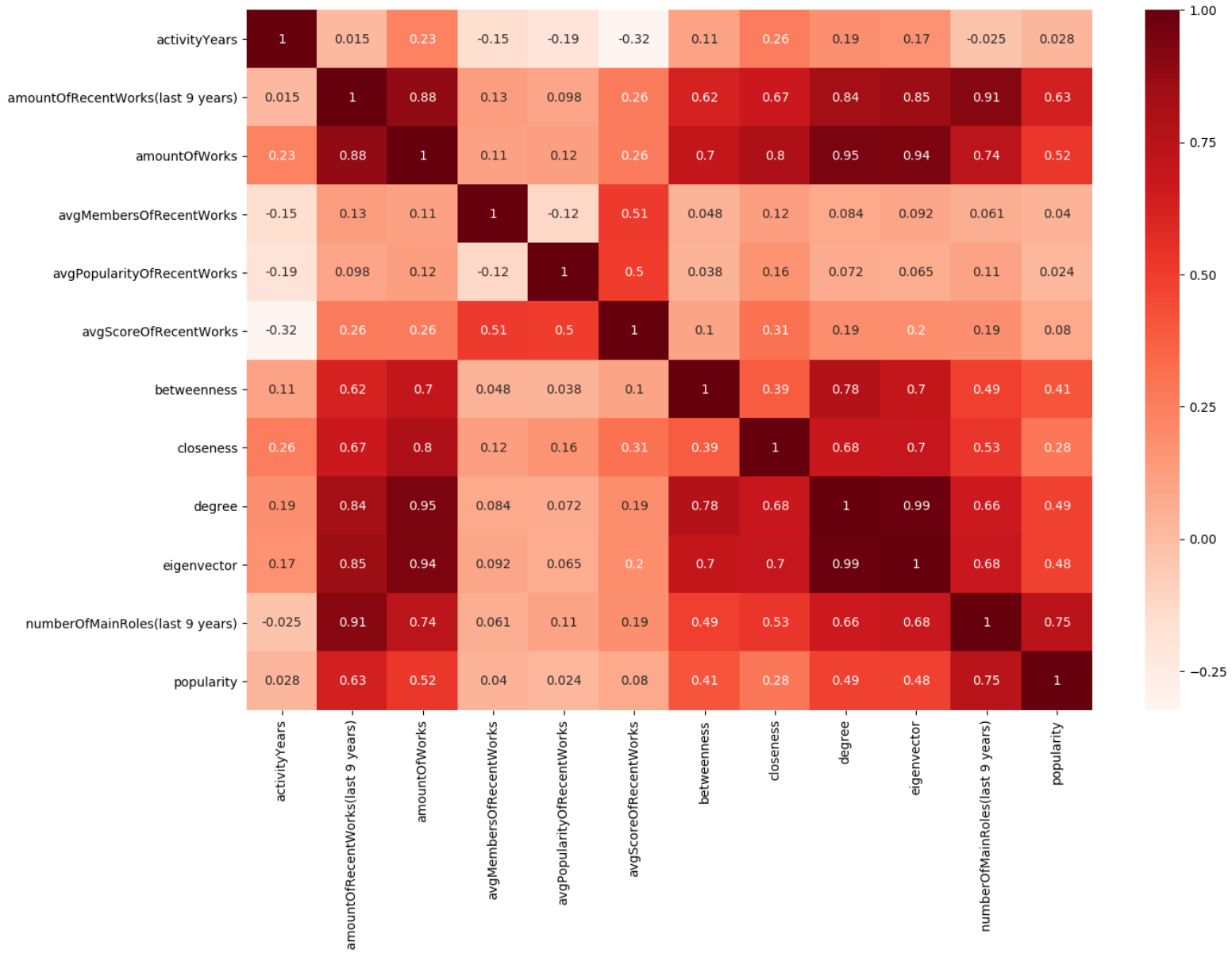
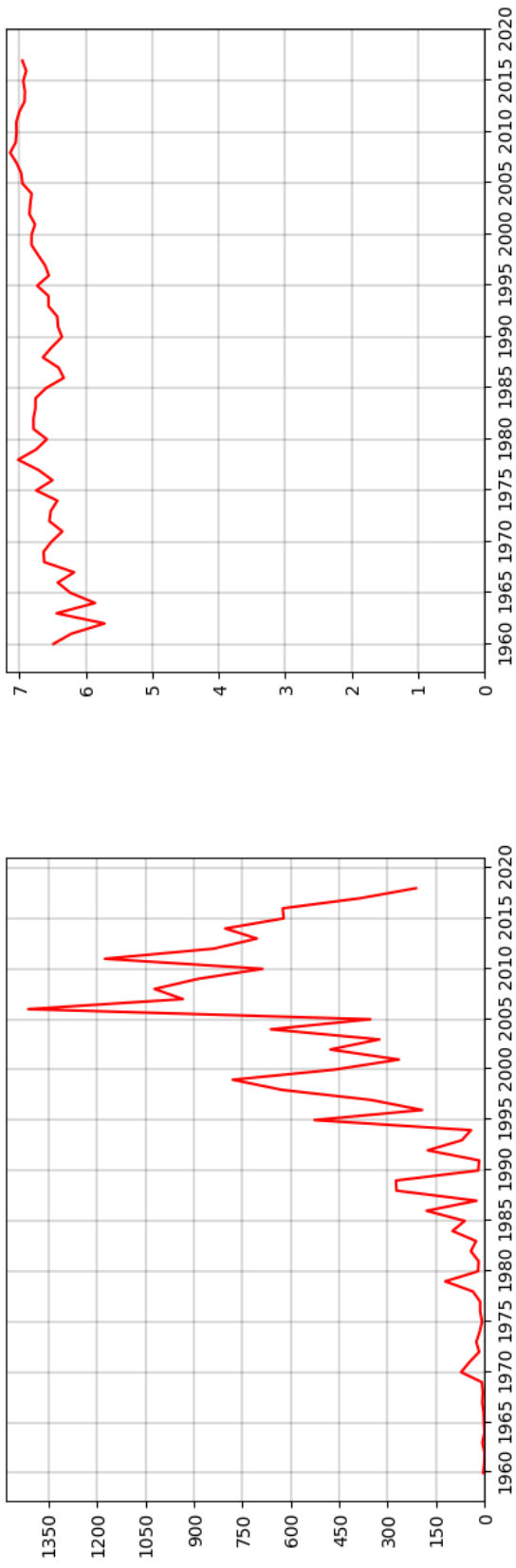


Figure 3.4: Pearson correlation between popularity and attribute of nodes (using recent works only).

3.1.1 Why last 9 years of works has more correlation?

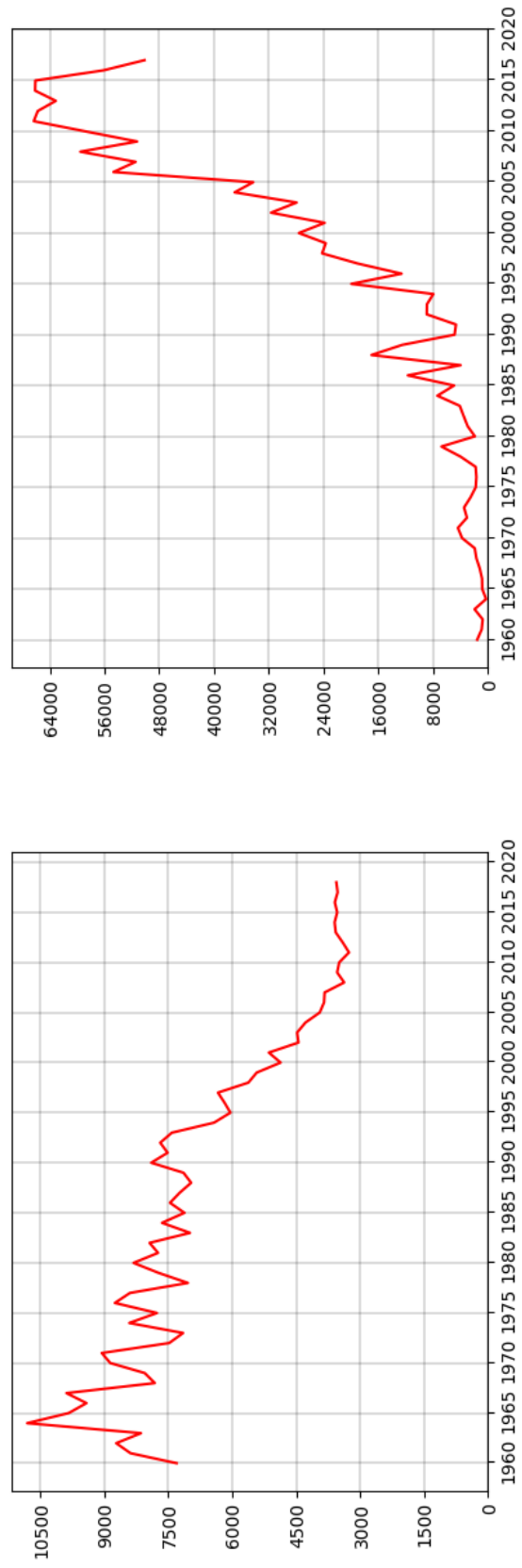
Graphics of some characteristics of works divided by years were made, trying to shed some light over why works from last 9 years were more "important".

Fig. 3.8a and 3.8b shows an improvement in average of scores and favorites. Biggest peak of favorites is on 2005, this may have to do with the start of MAL (2006), see Fig. 3.8a. We suppose as users started to use MAL they favorite

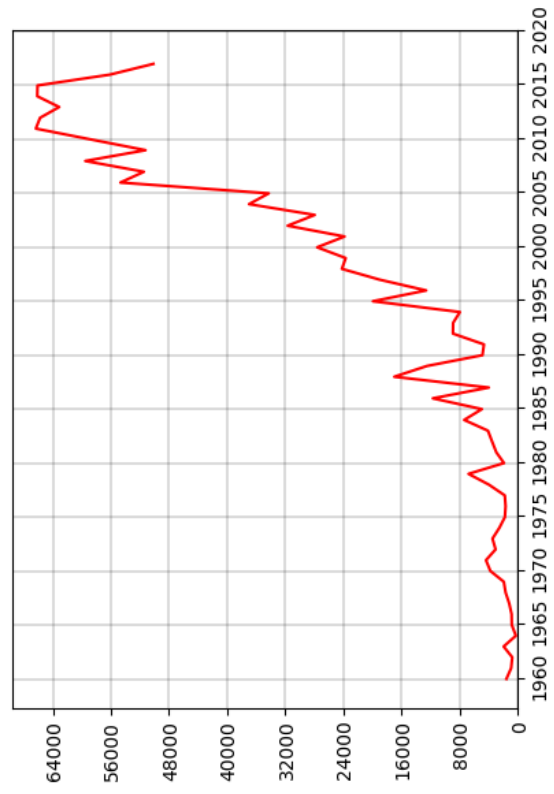


(a) Average amount of favorites per year.

(b) Average score per year.



(c) Average popularity per year.



(d) Average amount of members per year.

Figure 3.5: Averages of some attributes of anime divide by year.

anime they liked from that year and only some of the old ones; from then on they used the website often and favorite new works as they began airing.

Year 2018 was left out of Fig. 3.8b because, since 2018 is not finished yet, the average of scores was unusually small.

As Fig. 3.5c shows, popularity goes down in last years. But we need to consider that this metrics are from MAL and it could mean, for example, the parameter is in disuse; instead of representing how people feel about new anime.

We suppose the attribute "members" of an anime is taken from how many users have it on any of their lists. If that is correct Fig. 3.5d tells us feature of adding anime in watched / watching / plan to watch lists is really used. As for our experience on the web and social media we can confirm this is the most used feature of MAL. So this makes it one of the best metrics to measure "public" as another sense of popularity of anime.

As we can see on Fig. 3.6 anime industry is growing bigger each year, of course this is biased by the fact MAL will sure have every adaptation of last year but maybe not for anime from 1980.



Figure 3.6: Amount of works divided by years which were aired for the first time.

The majority of works are from 1990 to 2018 and half of them are distributed over the last 14 years (2014 to 2018) but as far as we can tell there isnt anything particular over the last 9 years nor on year 2009. Judging by amount of favorites per year it appears that users have been more active in recent years so this could be one of the reasons.

3.2 Correlation with multiple features

For this section Scikit-learn, a free software machine learning Python library, was used. The node attributes were divided into categories, leaving four distinct types:

- Personal data:
 - Debut
 - Gender
 - Activity years (2018-debut)
- Works data:
 - Amount
 - Top 5 genre
 - Favorites
 - Score
 - Popularity
 - Members
 - Number of main roles
- Recent works data:
 - Same as works but for only last 9 years
- Graph data:
 - Degree
 - Betweenness centrality
 - Closeness

Fitting and prediction experiments were run for each category, each combination of 2, 3 and all of them together; using 80% of seiyuu as train data and the rest as test. This was done for all following models:

- DecisionTreeRegressor
- DecisionTreeClassifier
- LinearRegression
- KNeighborsClassifier
- LinearDiscriminantAnalysis
- GaussianNB
- SVM

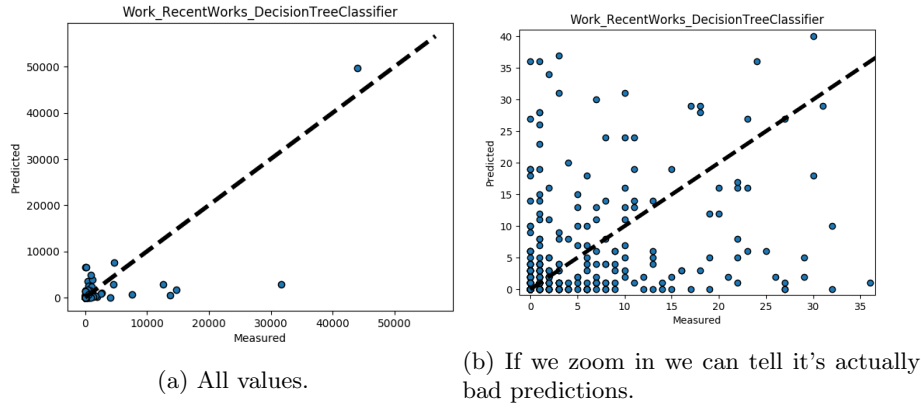


Figure 3.7: Scatter plot showing predicted / measured for Decision Tree Classifier using model Work + RecentWorks.

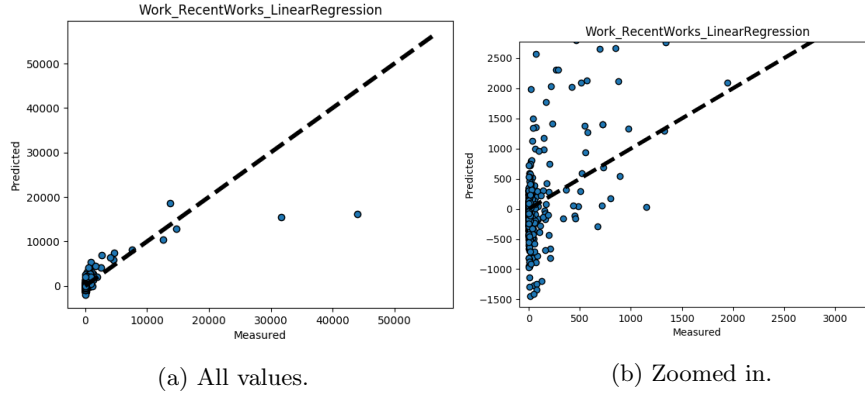


Figure 3.8: Scatter plot showing predicted / measured for Linear Regression using model Work + RecentWorks.

We used `r2_score`¹ for accuracy comparison. Problem is, since popularity variance is really high, we observed some good results in terms of `r2_score` values but particular predictions were a loof.

If we predict without error a seiyuu whose popularity is >50000 but we make "small" mistakes predicting seiyuu with <100 popularity then, is it a good prediction? and having into account more than $\frac{3}{4}$ of them have <100 popularity?

This was happening to our predictions. Since they were on spot for higher popularity values their prediction performance's metrics were good; but in fact they failed on small values, which is the big majority of them.

Fig. 3.7 and Fig. 3.8 shows some examples.

On the next subchapters we present `r2_score` values for each algorithm and

¹http://scikit-learn.org/stable/modules/model_evaluation.html#r2-score-the-coefficient-of-determination

each combination of categories. Along with information about which feature is more important, from the Decision Tree Classifier’s point of view (for only some of the models).

3.2.1 Only one category

Table 3.2: Only one category R2 score results

	Personal	Graph	Work	RecentWorks
DecisionTreeClassifier	-0.02	-3.26	-2.12	-0.59
DecisionTreeRegressor	-0.03	-1.05	-0.96	0.35
GaussianNB	-0.47	-0.57	0.01	0.03
KNeighborsClassifier	-0.02	0.04	0.07	0.10
LinearDiscriminantAnalysis	-0.02	-4.34	0.31	0.49
LinearRegression	-0.00	0.13	0.31	0.48
SVM	-0.02	-0.23	-0.02	0.02

Table. 3.2 reveals that Recent Works is the most performant if we are only using one category.

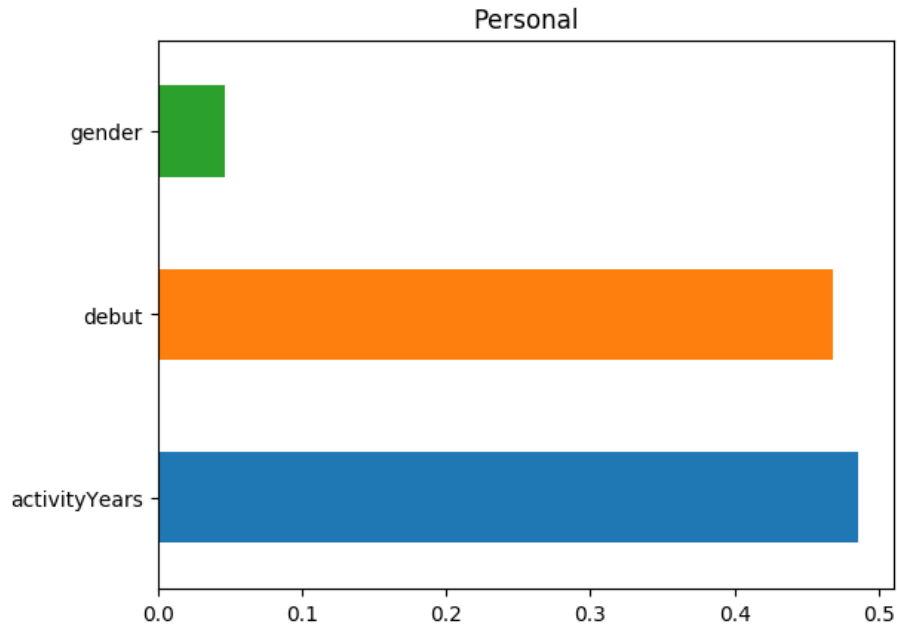


Figure 3.9: Feature information contribution for Personal data.

Fig. 3.9 shows that Decision Tree Classifier firstly uses activity years and debut to distinguish between most clases.

Decision Tree Classifier has bad performance when using only Recent works data but there’s something very interesting about Fig. 3.10, when using Work

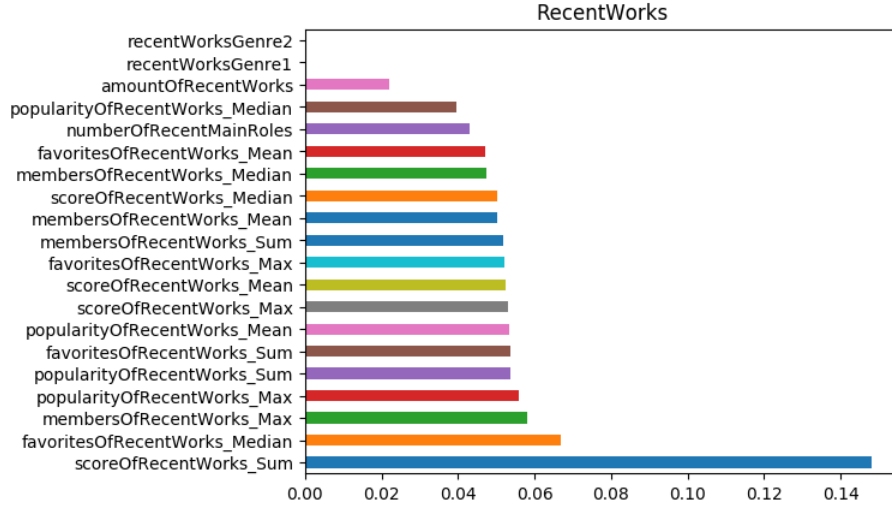


Figure 3.10: Feature information contribution for Recent works.

and Recent works "number of *recent* main roles" is not used to classify but here it is used; we suppose is because it doesn't provide new information against "number of mail roles" (see Fig. 3.11).

3.2.2 Groups of two categories

Table 3.3: Two categories R2 score results (R: recent works, P: personal, G: graph, W: work)

	P+G	P+W	P+R	G+W	G+R	W+R
DecisionTreeClassifier	0.11	-0.43	-1.86	-0.70	-2.51	0.57
DecisionTreeRegressor	-2.56	-1.15	-0.37	-0.21	-0.86	0.39
GaussianNB	-0.06	0.03	0.01	0.01	0.07	0.00
KNeighborsClassifier	0.10	0.15	0.03	0.08	0.25	0.08
LinearDiscriminantAnalysis	-0.35	0.43	-0.08	0.21	0.40	0.62
LinearRegression	0.38	0.57	0.38	0.52	0.33	0.63
SVM	-0.02	0.00	-0.04	-0.00	0.03	-0.02

Table. 3.3 shows ones of the best r^2 scores of all predictions. Decision Tree Classifier, Linear Discriminant Analysis and Linear Regression all did well when using Work + Recent Works model.

As Fig. 3.11 reveals number of main roles is a very important feature to classify seiyuu by popularity. We saw earlier that it has a strong correlation.



Figure 3.11: Feature information contribution for Work + Recent works model.

3.2.3 Groups of three categories

Table 3.4: Three categories R2 score results (R: recent works, P: personal, G: graph, W: work)

	P+G+W	P+G+R	P+W+R	G+W+R
DecisionTreeClassifier	-0.90	-8.46	-0.72	-0.73
DecisionTreeRegressor	-1.75	-2.12	-1.55	0.20
GaussianNB	0.01	0.05	0.04	0.02
KNeighborsClassifier	0.06	0.64	-0.04	0.18
LinearDiscriminantAnalysis	-0.16	-0.19	0.54	0.51
LinearRegression	-0.96	-1.63	0.09	0.38
SVM	-0.03	-0.04	-0.03	-0.02

According to Table. 3.4 three categories models have bad performance for most algorithm but it also has the highest (best) value of all. Is achieved by K Neighbors Classifier when using Personal + Graph + Recent works data.

Fig. 3.12 tells us, as we already saw before, that number of main roles is a very important feature. Also, notice that number of *recent* main roles does not appear.

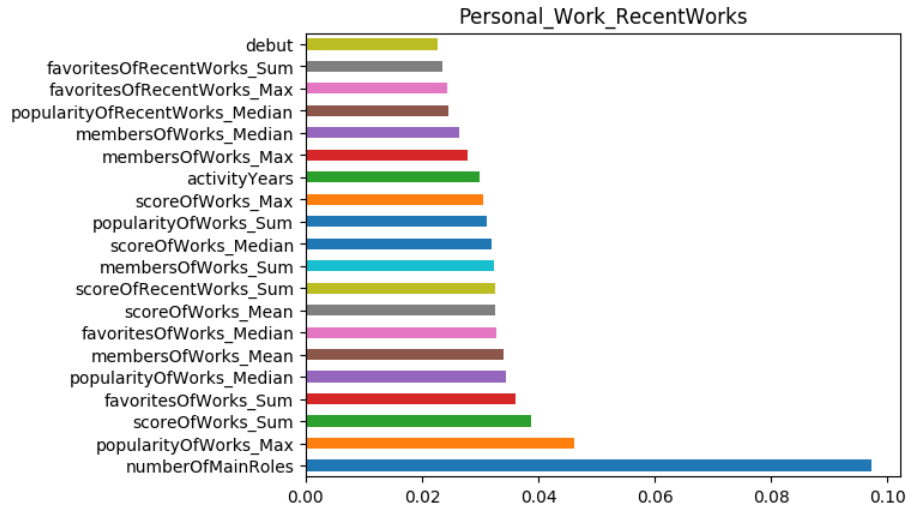


Figure 3.12: Feature information contribution for Personal + Work + Recent works model.

3.2.4 All categories

Table 3.5: All categories R2 score results

	AllFeatures
DecisionTreeClassifier	-0.19
DecisionTreeRegressor	0.49
GaussianNB	0.01
KNeighborsClassifier	0.13
LinearDiscriminantAnalysis	0.53
LinearRegression	0.55
SVM	0.01

Table. 3.5 Decision Tree Classifier has bad performance when using all categories, this can be because it overfits train data. Linear Discriminant Analysis and Linear Regression doesn't seem to do that bad, Decision Tree Regressor also. Of course, once we saw the scatter plots particular predictions were wrong most of the times for small values, showing a very different reality from value of metrics.

As Fig. 3.13 shows, even when using all features "number of main roles" is still the one that separates more values, having almost double importance than the second one (sum of score of works).

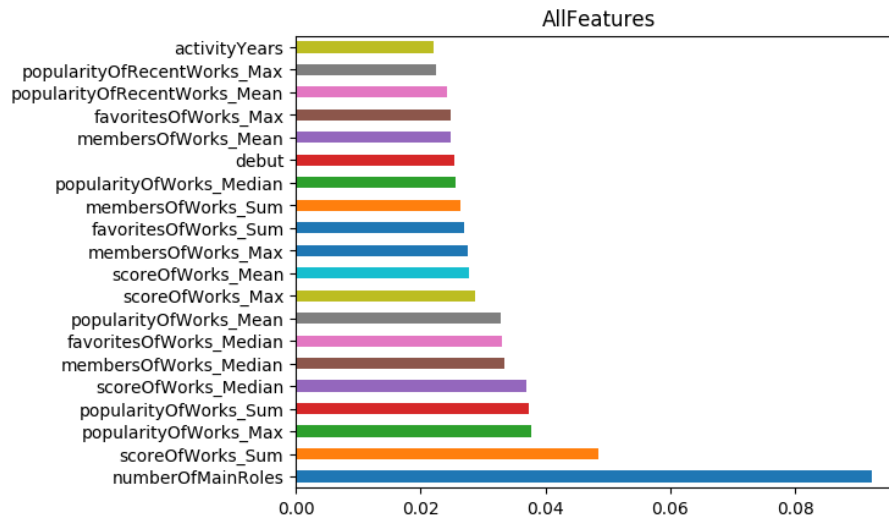


Figure 3.13: Feature information contribution for model using all categories.

3.3 Conclusion

Since particular predictions were usually wrong (mostly for small values) this couldn't be used in a practical way to predict popularity of new seiyuu nor to recover lost values.

TODO, LO BUENO ES QUE PUDIMOS VER IMPORTANCIA DE LOS FEATURES, QUE NUMBER OF MAIN ROLES ES IMPORTANTE, QUE LAS METRICAS SON DIFICILES DE USAR CUANDO HAY MUCHA VARIANZA?, PUDIMOS COMPARAR LOS DISTINTOS ALGORITMOS, VIENDO CUANDO ERAN MAS O MENOS PERFORMANTES...PO...NE...LE

Conclusion

TODO

Bibliography