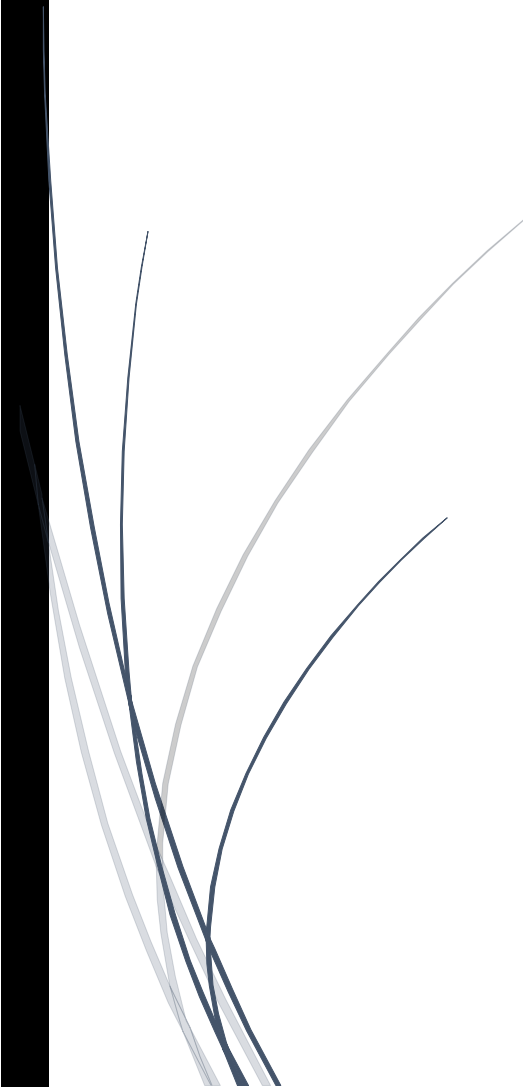4/11/2016

# OPIM 5604 Predictive Modeling

## Final Project – Group 5

## Professor Cruz

NAVDEEP CHEEMA
RAVI KOLLA
JENNIFER NIELSEN
YIFEI TANG

# Table of Contents

# Executive Summary

The dating industry has grown rapidly over the past decade and is comprised of many players. Services such as telephone chat lines, singles bars, speed dating organizations, online dating, mobile dating, and others have shaped the industry into what it is now and faded the stigma associated with using dating services. The U.S. single population is estimated to be close to 100 million people (The Business of Dating). An increasing number of these singles are turning to dating services to find love, and it is working.

Satisfied customers are not pushing their friends to join the newest app in the market, but rather, to the one they had the most luck with. Therefore, reputation is one critical success factor. Dating services are rising their level of concentration to more effectively match customers to potential dates and attract new customers.

Our group project explored the dataset from a speed dating experiment facilitated by Columbia University. Specifically, we observed the perceived importance of 6 attributes and their impact on the subject's decision to go on a date. The purpose of comparison of these attributes can aid in the matching process for any dating service, not just speed dating. Additionally, there is significance in collecting useful data before matching to improve business performance and to provide insight into what men and women are looking for in a partner. For instance, it can be argued that customers who both equally value ambition in a partner are more suited for each other. If this data can be collected in advance, companies can more efficiently create matches. Dating service providers can strengthen their reputation using similar techniques to create successful matches.

Following the SEMMA approach, the speed dating data was first partitioned and variables were defined. We explored the data using various visualization techniques to discover patterns and abnormalities among the variables (e.g. correlations, outliers, missing data). Because the dataset was

already cleaned, no modification steps were performed. Logistical regression, decision tree, and neural network models were then created to predict the subject's likelihood of accepting or rejecting a date with a potential mate from the speed dating event.  Finally, the models were measured and evaluated for usefulness. It was decided that the logistical regression model was the most optimal due to its simplicity and efficiency.

## Objectives

- Design and implement a model that would most accurately predict a subject's decision to go on a date based on his/her attribute preference in a partner

- Examine and report on patterns of gender behavior  differences when dating

- Explore trends to inform  stakeholders  how to market more effectively for client acquisition and retention

- Discover potential additional  business opportunities for dating service providers

- Provide recommendations that could make the speed date process more effective

## Data Introduction

### Speed Dating Event

The experimental speed dating exercise took place for 2 years from 2002 to 2004 as part of a Columbia University study. Students were recruited to be subjects and completed a pre-event survey online. In this survey, they were asked to rate their preference for 6 attributes (intelligence, attractiveness, sincerity, fun, shared interests, and ambition) in a partner.

During the speed dating event, subjects met with a potential mate for 4 minutes and filled out a scorecard afterwards. This scorecard indicated each subject's decision to either accept or reject a date

with who they just met with. This decision became our response variable in the modeling techniques. Subjects were also asked to rate their partner on those six attributes on the scorecard.

## Dataset

The dataset contains 194 variables, which contain the data for the speed dating event. The data has the participant's name, age, address, profession, race, preferences and the attributes they look for in partner. For this project, we've decided to not include all the variables, and work only on what attributes have the most significance for male and female participants. The team has decided to use the following six attributes that were most significant to choose the partner:

- Attractive (Attr)
- Sincere (Sinc)                    *please take note of the abbreviations*
- Intelligent (Intel)               *which are referenced throughout the paper*
- Fun (Fun)
- Ambitious (Amb)
- Shared Interests (Shar)

The decision variable is present in the dataset which tells us whether the participant is willing to go on another date with the partner.

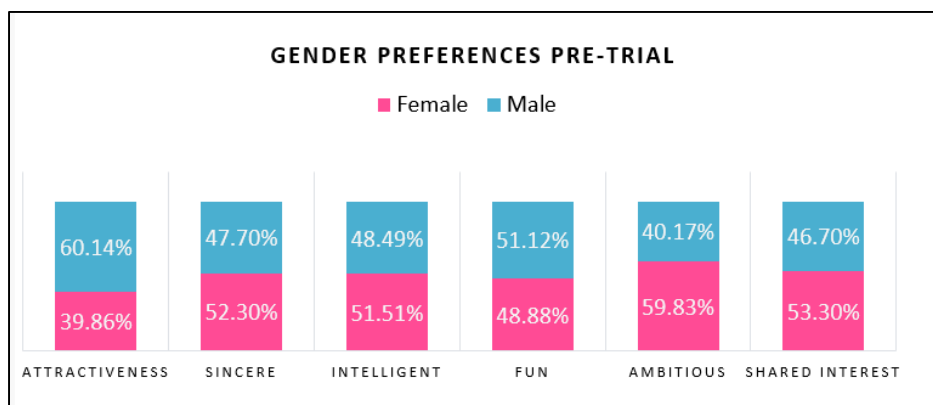| iid | unique subject number, group (wave id gender) |
| --- | --- |
| dec | Participant's decision which tells us whether the participant is willing to go on another date with the current date. 1-Yes, 0-No |
| attr | Attractiveness of the partner. Rated on a scale of 1-10. (1 being the least, and 10 being the highest) |
| sinc | Describes how sincere the partner is. Rated on a scale of 1-10. (1 being the least, and 10 being the highest) |
| intel | Intelligence of the partner. Rated on a scale of 1-10. (1 being the least, and 10 being the highest) |
| fun | Fun factor of the partner. Rated on a scale of 1-10. (1 being the least, and 10 being the highest) |
| amb | ambitiousness of the partner. Rated on a scale of 1-10. (1 being the least, and 10 being the highest) |
| shar | shred interests of the partner. Rated on a scale of 1-10. (1 being the least, and 10 being the highest) |
| gender | Female=0 Male=1 |

# Data Analysis

## Sample

The dataset is a sample of 8378 rows of data from the speed dating event, where each row contains the data for individual decisions. We have selected 6 explanatory variables (attributes mentioned above) and 1 decision variable (see Figure 1). The data has been partitioned into training and validation sets, based on the "iid" column (uniquely identifies each participant in the experiment) and used the "stratified random" technique.
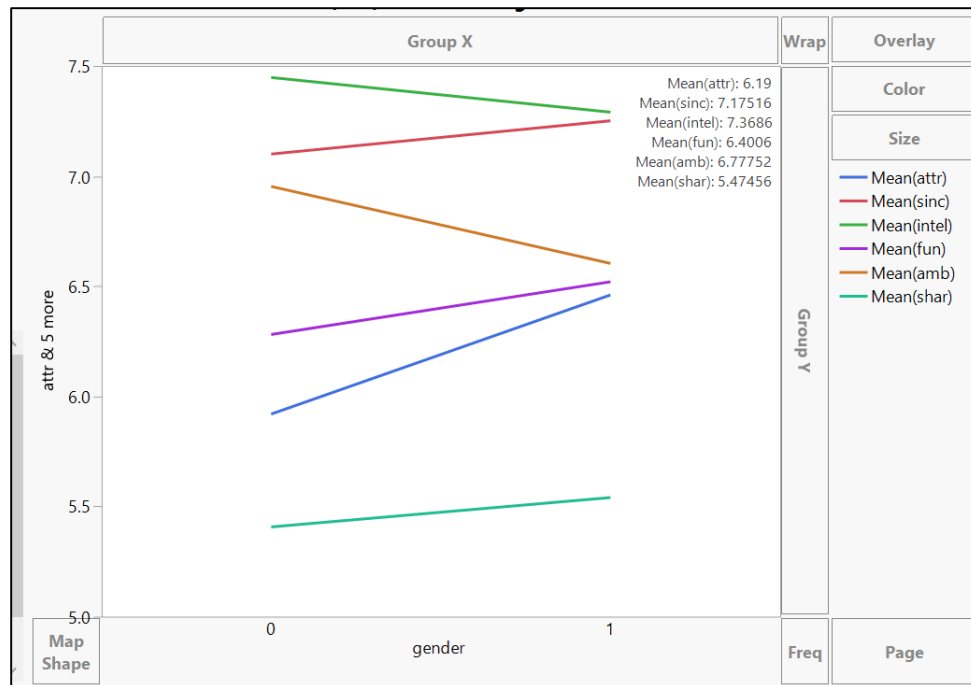
## Explore

In our dataset, the attributes are all continuous variables with the values ranging from 1-10. The decision variable and the gender are ordinal values, with values 1 or 0. The following exploring steps were performed:

- Began exploring the data for missing values (see Figure 2). Dataset has 8378 rows, of which 1338 rows have missing values. The missing data accounts for 15.9% of the data, but the team decided to leave the missing values as they are. These missing values belong to participants and it only makes sense to leave them blank, as one can't simply attribute one person's behavior to another person.

- Observed that there were no inconsistent values. Also, the dataset is fairly balanced, in terms of male and female participation (see Figure 3).

- Analyzed the data for outliers to see if there are any extreme values for the behavioral attributes. The boxplots in Figure 4 suggests that there might be some outliers in the dataset, but in reality, these values are bounded to a defined scale of 1-10.

- Performed "Mahalanobis outlier analysis" to see if the 2 methods tell a different story (see Figure 5).

- Searching for a potential reduction in dimensionality, the team decided it would be interesting to look for correlation between all the variables. When the correlation was performed (see Figure 6), we found that there is no negative correlation among the variables. It makes sense because these are behavioral attributes, and there wouldn't be any scenario where one attribute would cause the other attribute to go down. (ex: A highly intelligent person is not so attractive). These attributes are unique for each person and we can't predict a definite pattern. With not so strong correlation among these variables, performing principal component analysis, doesn't add much value to model and it doesn't call for the reduction of dimensionality.

- A quick review of the data provides some interesting insights. First, we report the valuation of attributes by men and women. Women put greater weight on intelligence than men do, while men place more value on physical appearance. We observe that a man's demand for intelligence and ambition does not extend to women who are more intelligent or ambitious than he is. In fact, a man is significantly less likely to accept a woman who is more ambitious than he. Also, women put more emphasis on their partner's ambitiousness, while men are more concerned about the sincerity and the fun factor when choosing the partner. The common attribute where men and women both put similar value is on the similar interests. The visual on the left and below illustrates these differences.

**GENDER PREFERENCES PRE-TRIAL**

■ Female ■ Male

| ATTRACTIVENESS | SINCERE | INTELLIGENT | FUN | AMBITIOUS | SHARED INTEREST |
|---|---|---|---|---|---|
| 60.14% (Male) | 47.70% (Male) | 48.49% (Male) | 51.12% (Male) | 40.17% (Male) | 46.70% (Male) |
| 39.86% (Female) | 52.30% (Female) | 51.51% (Female) | 48.88% (Female) | 59.83% (Female) | 53.30% (Female) |

Group X | Wrap | Overlay
Color
Size

7.5

Mean(attr): 6.19
Mean(sinc): 7.17516
Mean(intel): 7.3686
Mean(fun): 6.4006
Mean(amb): 6.77752
Mean(shar): 5.47456

Mean(attr)
Mean(sinc)
Mean(intel)
Mean(fun)
Mean(amb)
Mean(shar)

7.0

attr & 5 more

6.5

Group Y

6.0

5.5

5.0

Map
Shape

0    1

gender

Freq | Page

## Modify

There was not an opportunity to modify this data in the project, although for future exploration and recommendations one suggestion is to combine data set with additional age groups to provide a more comprehensive case comparison.

## Model

In the model analysis, we established three models: logistic regression model, decision tree and neural network, and selected the most favorable by comparing simplicity and efficiency of these three models.

## Logistical Regression

Ran the Stepwise for female and male respectively, refer to images below.



For female, all the six attributes are chosen, and for male, intelligence is excluded by Stepwise. We ran the logistic regression model for female to see the contribution of each attribute.



As can be seen from the report, among all the variables intelligence is the least significant one and is insignificant at 0.01 level.

As a result, we decide to choose five attributes to fit the model for both male and female: attractiveness, shared interests, fun, ambition and sincerity.

Both female and male care more about attractive while they decide to accept or decline the date. One noteworthy thing is that males pay much more attention to attractiveness than females (see chart on right).



LogWorth of Attributes

*Decision Tree*

Following the logistical regression, a decision tree model was built. The below table summarizes the key measurements from this model:

| Female | Male |
| --- | --- |
|  |  |

| | RSquare | N | Number of Splits |
| --- | --- | --- | --- |
| Training | 0.223 | 3124 | 7 |
| Validation | 0.212 | 1060 | |

| | RSquare | N | Number of Splits |
| --- | --- | --- | --- |
| Training | 0.309 | 3160 | 18 |
| Validation | 0.259 | 1034 | |

## Column Contributions

| Term | Number of Splits | G^2 | | Portion |
|------|------------------|-----|---|---------|
| attr | 2 | 583.234646 | | 0.6357 |
| fun | 3 | 197.368805 | | 0.2151 |
| shar | 2 | 136.804268 | | 0.1491 |
| sinc | 0 | 0 | | 0.0000 |
| intel | 0 | 0 | | 0.0000 |
| amb | 0 | 0 | | 0.0000 |

## Column Contributions

| Term | Number of Splits | G^2 | | Portion |
|------|------------------|-----|---|---------|
| attr | 5 | 1027.15819 | | 0.7605 |
| fun | 3 | 147.078151 | | 0.1089 |
| shar | 4 | 117.851327 | | 0.0873 |
| intel | 3 | 35.3006654 | | 0.0261 |
| amb | 2 | 12.2866138 | | 0.0091 |
| sinc | 1 | 11.0433387 | | 0.0082 |

*Neural Network*

A neural network model was also built, with the following results:

| Female | Male |
|--------|------|



**Diagram**

### Female

**Confusion Rates**

| Actual dec | Predicted Rate 0 | 1 |
|------------|------------------|---|
| 0 | 0.846 | 0.154 |
| 1 | 0.404 | 0.596 |

**Confusion Rates**

| Actual dec | Predicted Rate 0 | 1 |
|------------|------------------|---|
| 0 | 0.834 | 0.166 |
| 1 | 0.396 | 0.604 |

**Training — dec**

| Measures | Value |
|----------|-------|
| Generalized RSquare | 0.3890205 |
| Entropy RSquare | 0.2544037 |
| RMSE | 0.4037183 |
| Mean Abs Dev | 0.3264346 |
| Misclassification Rate | 0.2464327 |

**Validation — dec**

| Measures | Value |
|----------|-------|
| Generalized RSquare | 0.3683719 |
| Entropy RSquare | 0.2377743 |
| RMSE | 0.4106212 |
| Mean Abs Dev | 0.3339572 |
| Misclassification Rate | 0.2534722 |

### Male

**Confusion Rates**

| Actual dec | Predicted Rate 0 | 1 |
|------------|------------------|---|
| 0 | 0.779 | 0.221 |
| 1 | 0.243 | 0.757 |

**Confusion Rates**

| Actual dec | Predicted Rate 0 | 1 |
|------------|------------------|---|
| 0 | 0.770 | 0.230 |
| 1 | 0.278 | 0.722 |

**Training — dec**

| Measures | Value |
|----------|-------|
| Generalized RSquare | 0.4477922 |
| Entropy RSquare | 0.2952619 |
| RMSE | 0.4015523 |
| Mean Abs Dev | 0.321564 |
| Misclassification Rate | 0.2313433 |

**Validation — dec**

| Measures | Value |
|----------|-------|
| Generalized RSquare | 0.4212878 |
| Entropy RSquare | 0.2739682 |
| RMSE | 0.4094664 |
| Mean Abs Dev | 0.328166 |
| Misclassification Rate | 0.2535991 |

### ⊿ **Receiver Operating Characteristic**

| dec | Area |
|-----|--------|
| — 0 | 0.8236 |
| — 1 | 0.8236 |

### ⊿ **Receiver Operating Characteristic**

| dec | Area |
|-----|--------|
| — 0 | 0.8430 |
| — 1 | 0.8430 |

### ⊿ **Lift Curve**

| dec | |
|-----|---|
| — 0 | |
| — 1 | |

### ⊿ **Lift Curve**

| dec | |
|-----|---|
| — 0 | |
| — 1 | |

## Assess

The 3 models were compared side-by-side. The value of $R^2$ is important because it measures how close the data is to the fitted regression line. This statistic provides some information about the goodness of fit of a model.

| Measures of Fit for dec | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|
| Creator | .2 .4 .6 .8 | Entropy RSquare | Generalized RSquare | Mean -Log p | RMSE | Mean Abs Dev | Misclassification Rate | N |
| Fit Ordinal Logistic | | 0.2640 | 0.4066 | 0.5032 | 0.4083 | 0.3338 | 0.2502 | 7054 |
| Partition | | 0.2678 | 0.4107 | 0.4981 | 0.4069 | 0.3304 | 0.2480 | 8378 |

The confusion matrix table visually described the performance of the model. The results are summarized in the below chart:

| | Logistic Regression | | Decision Tree | | Neural Network[1] | |
|---|---|---|---|---|---|---|
| | 1 | 0 | 1 | 0 | 1 | 0 |
| Dec 1 | 29.32% | 13.82% | 27.75% | 14.24% | 59.59% | 13.54% |
| Dec 0 | 11.20% | 45.66% | 10.56% | 47.45% | 10.85% | 46.02% |
| Misclassification Rates | 0.2502 | | 0.2480 | | 0.2439 | |
| |  | |  | | * | |

All the three models have relative low Entropy RSquare due to ordinal responses. By comparing the misclassification rates, we find that all the models have similar performance. However, the simpler one is

---

[1] Unable to save probability formula for Neural Network model

the better one. Logistic regression model is the simplest model and easier to understand even though its misclassification rate is a little bit higher than those of any other two.

## Model Interpretation

The value in the table below is the mean rating of each attribute by gender. Given the mean rating of a partner, a man is 46.86% likely to choose to go on a date, and woman, 30.75%.

|  | Male | Female |
|---|---|---|
| attr | 6.46 | 5.92 |
| sinc | 7.25 | 7.10 |
| fun | 6.52 | 6.28 |
| amb | 6.60 | 6.95 |
| share | 5.54 | 5.41 |
| linear | -5.27130 | -4.44932 |
| fail | 0.531356 | 0.692544 |
| succeed | 0.468644 | 0.307456 |

If we add 1 to the rating of each attribute, some changes may happen in the probability of a successful date. As can be seen below:

|  | Prob [1] - Male | Change | Prob [1] - Female | Change |
|---|---|---|---|---|
| attr+1 | 0.6347960 | 35.45% | 0.402210 | 30.82% |

| | | | | |
|---|---|---|---|---|
| sinc+1 | 0.4227500 | -9.79% | 0.292927 | -4.73% |
| fun+1 | 0.5402740 | 15.28% | 0.373844 | 21.59% |
| amb+1 | 0.4296930 | -8.31% | 0.281075 | -8.58% |
| share+1 | 0.5319400 | 13.51% | 0.369480 | 20.17% |

1. The probability of a successful date will increase with the rise rating of attractiveness, fun and shared interests for both male and female. However, increased sincerity and ambition will have a negative effect on the probability of a successful date.

2. For both male and female, attractiveness makes the most difference. One more score in attractiveness will bring more than 35% probability for men to say YES, which is 5% higher than for women. That is to say, men care more about attractiveness than women when they choose speed dating partner.

3. The second most important factor is fun both for male and for female. Women pay more attention to fun than men.

*(Note: refer to formulas in Appendix B for further analysis.)*

## Business Insights

Speed dating was trademarked by Rabbi Yaacov Deyo in 1998 for singles specifically looking for a serious relationship and potential life partner. Maximizing the efficiency of identifying a marriage partner did not compromise quality of match-making results or quality of the clientele.

As the idea integrated into the dating industry, organizations tried implementing this process into areas that ended up weakening the results and cheapening the quality. While industry followers were initially successful in cashing in on this trend, it was apparent by 2003 that the market had lost control of this idea. This forum was initially so successful because it was about quality for a group of people looking

for a serious relationship. In order to revive this idea back to its maximum earning potential developers and investors had to redesign their approach.

In this speed dating study specifically we found that planning of the *pre dating* questions is valuable. Framing some of the pre dating questions to get insight on attributes and experiences of *prior* dating partners and experiences would be a recommendation for data analysts to explore. These insights on prior experiences could potentially help offer events and services that would increase match potential more than *desired* traits.

Attributes is an important area to explore when marketing for online dating – specifically speed dating. For example, in pre-trial ratings subjects anticipated intelligence would be something that would be a significant influence on the success rate of their match/preference. However, when actual results post -dating were examined, the model we used eliminated this attribute as a significant factor for males immediately during the match process with relevance to females quickly following in the 2nd round. Women weighted the attributes more evenly than the men did, with intelligence on top and ambition on the bottom.

## Recommendations

### Top Recommendations:

- Speed dating organizations should invest more effort into attractiveness and fun, since these were the top attributes in the model, as opposed to the insignificant attributes (e.g. ambition).
- Redefine the rating system to include most meaningful attributes and goals of target demographic to most effectively match participants to event.

- Offer separate events to keep participants involved most well matched.  For example, speed dating for those looking for a serious relationship and separate events for those looking more for companionship/friendship.

- To increase participation and success in this niche market, focus speed dating events on traditional relationship seeking audience.

- Offer events to financial planning associations to pitch their products and funds to potential investors. (Puko)

## Other Recommendations:

- Offer speed dating events to specific demographics. Studies show 80% of those dating prefer to date someone of their ethnic background.

- To expand customer base customize speed dating forum to offer services to business market. Employers look for low cost venues to help in the pre-screening process when replacing or filling a position. Employment agency fees range from 10%-25% of filled position salary, while use of headhunters for upper management /executive level positions can climb as high as 213% of annual salary. (Employment Agency Cost)

- Expand services to students, specifically foreign students for networking opportunities. It is often difficult to acclimate to a new culture. This type of service could help bridge the gap students and job seekers find when entering the workplace as an international candidate.

# References

*Employment Agency Cost*. Ed. INC CostHelper. n.d. 10 April 2016.
   <http://personalfinance.costhelper.com/employment-agencies.html>.

Merhar, Christina. *Small Business Employee Benefits and HR Blog*. February 4 2016. 10 April 2016.
   <http://www.zanebenefits.com/blog/bid/312123/Employee-Retention-The-Real-Cost-of-Losing-an-Employee>.

Puko, Tim. *The Wall Street Journal*. 5 February 2016. 10 April 2016.
   <http://blogs.wsj.com/moneybeat/2016/02/05/in-search-for-investors-hedge-funds-try-speed-dating/>.

*Random Facts*. 8 January 2010. April 2016. <http://facts.randomhistory.com/dating-and-relationship-facts.html>.

*The Business of Dating*. n.d. WikiDot. 10 April 2016. <http://thebusinessofdating.wikidot.com/the-dating-industry>.

# Appendix A: Exploration Screenshots

*Figure 1*



*Figure 2*

*Figure 3*



*Figure 4*

*Figure 5*



*Figure 6*

# Appendix B: Formulas

Leaner by gender:

$$\text{If}\begin{cases} gender == 1 \Rightarrow & \begin{aligned} & -0.6784391970108 * attr \\ & + 0.1859039621447 * sinc \\ & + -0.2870339366126 * fun \\ & + 0.15751352323395 * amb \\ & + -0.2535233187351 * shar \end{aligned} \\[1em] gender == 0 \Rightarrow & \begin{aligned} & -0.4157748641939 * attr \\ & + 0.06917439780028 * sinc \\ & + -0.2962762226984 * fun \\ & + 0.12709507296394 * amb \\ & + -0.2775904211013 * shar \end{aligned} \\[1em] else \Rightarrow & . \end{cases}$$

Cum [0] By gender:

$$\text{If}\begin{cases} gender == 1 \Rightarrow & \dfrac{1}{1 + Exp\left(-5.3968948396349 - Linear\ By\ gender\right)} \\[1em] gender == 0 \Rightarrow & \dfrac{1}{1 + Exp\left(-5.2613550602782 - Linear\ By\ gender\right)} \\[1em] else \Rightarrow & . \end{cases}$$

Probability [0] By gender:

$$\text{If}\begin{cases} gender == 1 \Rightarrow Cum[0]\ By\ gender \\ gender == 0 \Rightarrow Cum[0]\ By\ gender \\ else \Rightarrow . \end{cases}$$

Probability [1] By gender:

If
| | | | |
|---|---|---|---|
| gender == 1 | ⇒ | 1 - Cum[0] By gender |
| gender == 0 | ⇒ | 1 - Cum[0] By gender |
| else | ⇒ | . |