



# Speed Dating

## Brief analysis of a speed dating event

Alessandro Cogollo, 10571078

Systems and Methods For Big and Unstructured Data - Politecnico di Milano

# 1. Introduction

This paper describes a brief analysis conducted over a dataset about a set of speed dating experiments; the analysis has been conducted for educational purposes, as the reported queries represent only an overview of the potential investigations that can be conducted on this dataset. However, the dataset is of interest, since, as will be shown in the subsequent sections of this paper, it is possible to analyze numerous aspects related to the perception, interests, and expectations of the participants, also based on aspects studied by social sciences. With a more in-depth analysis, more numerous datasets, and techniques that go beyond the topics of this course it would be possible and interesting to deepen the research using predictive and inference algorithms to predict possible matches, as well as identify patterns regarding the preferences of any further participants.

The chosen dataset collects 1,012,321 entries, and has been downloaded via Kaggle, at the [following link](#). More specifically, it collects experimental data relating to speed dating events held between 2002 and 2004. The duration of the single date was set at four minutes, during which each participant had the opportunity to get to know and look for a possible "match" with (mostly) strangers of the opposite sex. The dataset contains different types of information, from information regarding the participant's lifestyle to interests and hobbies; more precise information on the structure of the dataset will be explained in the following sections.

The technology chosen is [elasticsearch](#), mainly for two reasons; the first concerns the possibility of using it as a "suggester" of results, thanks to the possibility of "boosting" the results, obtaining answers more suited to specific preferences, even though, as explained later, the dataset didn't contain lot of text to be inspected with an analyser. The second concerns the possibility of exploiting [Kibana](#) to display interesting insights relating to the queries performed. Last but not least, Elasticsearch has been chosen due to a personal interest in exploring this technology.

## 2. Data Wrangling

The chosen dataset already had very high quality and usability, measured in terms of completeness, credibility, and compatibility (both usage license and format), therefore the data wrangling process required was minimal.

A first inspection of the dataset was performed using the Pandas and NumPy libraries, to evaluate the possible need to normalize the data. The inspection result immediately showed the presence of “b” characters at the beginning of some fields (has\_null, gender, match, etc); although the dataset can also be used in these conditions, for reasons of consistency and quality of the dataset the interested columns have been cleaned with the python function shown below.

```
1 usage
1 def substr_b(x):
2     return x[2:(-1)]

1 usage
1 def syntaxcorr(df):
2     for col in df:
3         if isinstance(df[col][0], str) and (df[col][0])[0:2] == "b'" and (df[col][0])[-1] == "'":
4             df[col] = df[col].apply(lambda x: substr_b(x))
```

Further analysis involved exploring the statistics of the dataset, from which it was clear that the fields: attractive\_o and fun\_o must have been normalized in the range 0 to 10, as well as the fields gaming, and reading.

```
1 usage
1 def norm_met(df):
2     df['met'] = df['met'].apply(lambda x: set_one_zero(x))

1 usage
1 def mm_scaler(df, field, ran):
2     avg = df[field].mean()
3     df[field].fillna(value=avg)
4     scaler = MinMaxScaler()
5     scaler.fit(np.array(df[field]).reshape(-1, 1))
6     df[field] = scaler.transform(np.array(df[field]).reshape(-1, 1))
7     df[field] = df[field] * ran
8     df[field] = round(df[field], 0).astype('float32')
```

The “met” field should have been 0 or 1 (Boolean), but was between 0 and 8, thus it has been transformed using the following functions.

```

1 usage
✓ def set_one_zero(x):
    y = 0
    if x > 0:
        y = 1
    return y

1 usage
✓ def norm_met(df):
    df['met'] = df['met'].apply(lambda x: set_one_zero(x))

```

The first letter of each word in the “field” column has been transformed into all caps for consistency, even though term queries are case insensitive by default.

```

1 usage
def to_title(st):
    return st.title()

1 usage
def uniform_uppercase(df, field):
    df[field] = df[field].apply(lambda x: to_title(x))

```

Finally, derived values have been recomputed to be coherent with new values. The final dataset used has been attached to the delivery folder (dump subfolder) and named “speeddating-v1”; python code used for the data wrangling process has been attached as well.

### 3. Dataset

As mentioned in previous sections, the dataset contains 1,012,321 entries (excluding NULL values), in particular, the original CSV file is composed of 123 columns and 8378 rows.

The dataset gathers different information about the experiment, in particular, each entry contains information about a couple in a specific “wave” (round). For each entry, **demographic information** is collected, such as age, gender, “race”, and field of study of both participants. Additional information covers six attributes: **attractiveness**, **sincerity**, **intelligence**, **fun**, **ambition**, and **shared interests**. The dataset also includes other questionnaire data gathered from participants at different points in the process. These fields include **dating habits**, **self-perception** across key attributes, **beliefs on what others find valuable in a mate** and **lifestyle information**. Finally, the dataset contains information regarding the result of the match, the interest in dating, etc. For all relevant fields, derived information was calculated, mainly related to clusters to which certain values can be mapped to, for example about interests etc.

The dataset has been distributed on two shards, since its dimensions are quite small, and for this analysis performed there wasn’t a strong requirement on speed.

Regarding the implementation, the mapping is quite standard, but modified with respect to the version suggested by the elasticsearch ingestion tool, as most of the numerical fields have been casted to integer values; the only values that have been left as double (float64) are preferences, and correlation fields. String values have been set as “keywords”, as most of them are about clusters (d\_funny\_important, d\_intelligence\_partner, etc), therefore need to be searchable exactly for their specific value. The only exception is made by the attribute field which needs to be searchable for its content, since entries have not been homogenized. Some attributes, such as met, or match, could have been boolean, but would have required further data wrangling and would have brought low advantages.

Below can be found a very small portion of the mapping, the rest will be attached in the delivery folder under the name of “mapping.txt”.

```
{
  "mappings": {
    "_meta": {
      "created_by": "file-data-visualizer"
    },
    "properties": {
      "age": {
        "type": "integer"
      },
      "age_o": {
        "type": "integer"
      },
      "ambition": {
        "type": "integer"
      },
      "ambition_partner": {
        "type": "integer"
      }
    }
  }
}
```

## 4. Queries

The chosen queries represent a set of interesting requests both from an educational point of view and from social sciences point of view, with the aim of obtaining an insight into the preferences, demographics, and tendencies of the participants in the experiment. The queries were chosen with different degrees of complexity, and are shown below. All the results of the queries have been attached in the delivery folder.

### a. Highest difference in age of a match

The query below returns first filter aggregates all entries for which the field “match” is set to 1 (TRUE) and returns the maximum value (diff\_age) of these aggregates. The result is returned with the last line of the query and shown below.

```
1 GET /speeddating-v4/_search
2 {
3   "size" : "1",
4   "aggs": {
5     "more-diff-match" : {
6       "filter" : {
7         "term": {"match": "1"}
8       },
9       "aggs": {
10        "diff_age": { "max": { "field": "d_age" } }
11      }
12    },
13    "_source": ["diff_age"]
14  }
15 }

1 {
2   "took": 1,
3   "timed_out": false,
4   "_shards": {
5     "total": 1,
6     "successful": 1,
7     "skipped": 0,
8     "failed": 0
9   },
10  "hits": {
11    "total": {
12      "value": 8378,
13      "relation": "eq"
14    },
15    "max_score": 1,
16    "hits": [
17      {
18        "_index": "speeddating-v4",
19        "_id": "bFq81IwBkPnrhWx6FobH",
20        "_score": 1,
21        "_source": {}
22      }
23    ]
24  },
25  "aggregations": {
26    "more-diff-match": {
27      "doc_count": 1380,
28      "diff_age": {
29        "value": 34
30      }
31    }
32  }
33 }
```

The query could have been created also by using the range operator, as follows.

```

1 GET /speeddating-v4/_search
2 {
3   "size": 1,
4   "aggs": {
5     "more-diff-match": {
6       "filter": {
7         "range": {"match": {"gt": "0"}}
8       },
9       "aggs": {
10        "diff_age": { "max": { "field": "d_age" } }
11      }
12    }
13  },
14  "_source": ["diff_age"]
15 }
16

```

## b. Oldest participant to the speed dating experiment

This simple query returns the single element with the highest age in the dataset.

```

1 GET /speeddating-v4/_search
2 {
3   "size": 1,
4   "sort": [ {
5     "age": { "order": "desc" }
6   } ]
7 }
8

```

```

1 {
2   "took": 2,
3   "timed_out": false,
4   "_shards": {
5     "total": 1,
6     "successful": 1,
7     "skipped": 0,
8     "failed": 0
9   },
10  "hits": {
11    "total": {
12      "value": 8378,
13      "relation": "eq"
14    },
15    "max_score": null,
16    "hits": [
17      {
18        "_index": "speeddating-v4",
19        "_id": "Clq81Iw8kPnrhWx6N5bf",
20        "_score": null,
21        "_source": {
22          "d_funny_important": "[21-100]",
23          "d_reading": "[6-8]",
24          "ambtition_important": 0,
25          "d_age": 22,
26          "intelligence_important": 25,
27          "movies": 3,
28          "music": 8,
29          "d_tv": "[0-5]",
30          "d_expected_num_interested_in_me": "[0-3]",
31          "samrace": 0,
32          "shopping": 1,
33          "d_tvsports": "[0-5]",
34          "d_attractive": "[6-8]",
35          "d_ambition": "[0-5]",
36          "dining": 8,
37          "d_shared_interests_o": "[6-8]",
38          "intelligence": 7,
39          "d_funny": "[6-8]",
40          "d_ambtition_important": "[0-15]",
41          "d_sincere": "[6-8]",
42          "field": "Soa -- Writing",
43          "d_sincere_important": "[0-15]"

```

### c. Correlation between races and matches?

This query has been formulated to try to understand if there's any correlation between races and matches. First an aggregation on different races is performed, then the number of matches per aggregate is counted.

```
1 GET /speeddating-v4/_search
2 {
3   "size" : 0,
4   "aggs" : {
5     "race-count" : {
6       "terms" : {"field" : "race"},
7       "aggs" : {
8         "num-matches" : {
9           "terms" : {"field": "match"}
10        }
11      }
12    }
13  }
14 }
```

```
1 {
2   "took": 0,
3   "timed_out": false,
4   "_shards": {
5     "total": 1,
6     "successful": 1,
7     "skipped": 0,
8     "failed": 0
9   },
10  "hits": {
11    "total": {
12      "value": 8378,
13      "relation": "eq"
14    },
15    "max_score": null,
16    "hits": []
17  },
18  "aggregations": {
19    "race-count": {
20      "doc_count_error_upper_bound": 0,
21      "sum_other_doc_count": 0,
22      "buckets": [
23        {
24          "key": "European/Caucasian-American",
25          "doc_count": 4727,
26          "num-matches": {
27            "doc_count_error_upper_bound": 0,
28            "sum_other_doc_count": 0,
29            "buckets": [
30              {
31                "key": 0,
32                "doc_count": 3939
33              },
34              {
35                "key": 1,
36                "doc_count": 788
37              }
38            ]
39          }
40        },
41        {
42          "key": "Asian/Pacific Islander/Asian-American",
43          "doc_count": 1982,
44          "num-matches": {
45            "doc_count_error_upper_bound": 0,
46            "sum_other_doc_count": 0,
47            "buckets": [
48              {
```



Results shows the following numbers:

Race	Total	No Match	Match	Percentage
European/Caucasian-American	4727	3939	788	16.67%
Asian/Pacific Islander/Asian-American	1982	1715	267	13.47%
Latino/Hispanic American	664	541	123	18.52%
Black/African American	420	335	85	20.24%
Unknown	63	49	14	22.22%
Other	522	419	103	19.73%

#### d. Matches expectancies between men?

This query tries to understand what the expectancy in term of number of matches for the men participating in the experiment:

```
1 GET /speeddating-v4/_search
2 {
3   "aggs": {
4     "male_perception": {
5       "filter": { "term": { "gender": "male" } },
6       "aggs": {
7         "avg_expected_matches": { "avg": { "field": "expected_num_matches" } }
8       }
9     }
10  },
11  "size" : "0"
12 }
13
```

The average is: 3.39, which can be rounded to 3. Executing the same query for the women, the result is slightly lower, and is 3.02, which can still be rounded to 3.

#### e. How much attractive people consider attractiveness important when looking for an attractive partner?

Find out the importance that attractive people are giving to finding an attractive partner.

```

1 GET /speeddating-v4/_search
2 {
3   "query": {
4     "range": {
5       "attractive_o" : {"gt": 8}
6     }
7   },
8   "aggs": {
9     "avg-attr": {
10      "avg": {
11        "field": "attractive_important"
12      }
13    }
14  },
15  "size" : 0
16 }

```

```

2   "took": 1,
3   "timed_out": false,
4   "_shards": {
5     "total": 1,
6     "successful": 1,
7     "skipped": 0,
8     "failed": 0
9   },
10  "hits": {
11    "total": {
12      "value": 869,
13      "relation": "eq"
14    },
15    "max_score": null,
16    "hits": []
17  },
18  "aggregations": {
19    "avg-attr": {
20      "value": 23.490058479532163
21    }
22  }

```

**f. Return the IDs and field of entries about men, better if in the field of any science and with an age between 25 and 35**

Find out the importance that attractive people are giving to finding an attractive partner.

```

1 GET /speeddating-v4/_search
2 {
3   "query": {
4     "bool": {
5       "must": [
6         { "match" : { "gender" : "male" } } ],
7       "should": [
8         { "match" : { "field" : "Science" } },
9         { "range": { "age" : { "gt": 18, "lt" : 25 } } } ] ]
10  }, "_source" : [ "_id", "field" ]
11 }

```

```

1 {
2   "took": 1,
3   "timed_out": false,
4   "_shards": {
5     "total": 1,
6     "successful": 1,
7     "skipped": 0,
8     "failed": 0
9   },
10  "hits": {
11    "total": {
12      "value": 4194,
13      "relation": "eq"
14    },
15    "max_score": 5.1814966,
16    "hits": [
17      {
18        "_index": "speeddating-v4",
19        "_id": "mFq81IwBkPnrhWx6I41T",
20        "score": 5.1814966,
21        "_source": {
22          "field": "Computer Science"
23        }
24      },
25      {
26        "_index": "speeddating-v4",
27        "_id": "mVq81IwBkPnrhWx6I41T",
28        "score": 5.1814966,
29        "_source": {
30          "field": "Computer Science"
31        }
32      },
33      {
34        "_index": "speeddating-v4",
35        "_id": "mlq81IwBkPnrhWx6I41T",
36        "score": 5.1814966,
37        "_source": {
38          "field": "Computer Science"
39        }
40      },
41      {
42        "_index": "speeddating-v4",
43        "_id": "m1q81IwBkPnrhWx6I41T",

```

**g. Return the average value of all the possible interests, chosen from those participants that matched**

This query is interesting as it permits to determine which is the interest leading to a higher number of matches. It could have been performed less “mechanically” by iterating over an array in which interests were specified.

```

1 GET /dating-final/_search
2 {
3   "query":
4   { "match":
5     { "match": "b'1'" }
6   },
7   "aggs": {
8     "avg_sports": { "avg": { "field": "sports" } },
9     "avg_tv_sports": { "avg": { "field": "tvsports" } },
10    "avg_hiking": { "avg": { "field": "hiking" } },
11    "avg_exercise": { "avg": { "field": "exercise" } },
12    "avg_dining": { "avg": { "field": "dining" } },
13    "avg_museums": { "avg": { "field": "museums" } },
14    "avg_art": { "avg": { "field": "art" } },
15    "avg_gaming": { "avg": { "field": "gaming" } },
16    "avg_clubbing": { "avg": { "field": "clubbing" } },
17    "avg_reading": { "avg": { "field": "reading" } },
18    "avg_tv": { "avg": { "field": "tv" } },
19    "avg_theater": { "avg": { "field": "theater" } },
20    "avg_movies": { "avg": { "field": "movies" } },
21    "avg_concerts": { "avg": { "field": "concerts" } },
22    "avg_music": { "avg": { "field": "shopping" } },
23    "avg_yoga": { "avg": { "field": "yoga" } }
24  },
25  "size": 0
26 }
27

```

```

1 {
2   "took": 0,
3   "timed_out": false,
4   "shards": {
5     "total": 1,
6     "successful": 1,
7     "skipped": 0,
8     "failed": 0
9   },
10  "hits": {
11    "total": {
12      "value": 1380,
13      "relation": "eq"
14    },
15    "max_score": null,
16    "hits": []
17  },
18  "aggregations": {
19    "avg_reading": {
20      "value": 7.76996336996337
21    },
22    "avg_dining": {
23      "value": 7.917948717948718
24    },
25    "avg_theater": {
26      "value": 6.7765567765567765
27    },
28    "avg_art": {
29      "value": 6.875457875457876
30    },
31    "avg_gaming": {
32      "value": 3.9611721611721613
33    },
34    "avg_tv_sports": {
35      "value": 4.547252747252747
36    },
37    "avg_movies": {
38      "value": 7.831501831501831
39    },
40    "avg_sports": {
41      "value": 6.553113553113553
42    },
43    "avg_concerts": {
44      "value": 6.956776556776557
45    }
46  }
47 }

```

- h. Who are the three “top rated” people in the dataset? (terms of attractiveness, intelligence and funny)

```
1 GET /dating-final/_search
2 {
3   "size" : 3,
4   "query": {
5     "match_all": {}
6   },
7   "sort": [
8     { "attractive_o": { "order": "desc" } },
9     { "funny_o": { "order": "desc" } },
10    { "intelligence_o": { "order": "desc" } } ]
11 }
```

```
1 {
2   "took": 2,
3   "timed_out": false,
4   "_shards": {
5     "total": 1,
6     "successful": 1,
7     "skipped": 0,
8     "failed": 0
9   },
10  "hits": {
11    "total": {
12      "value": 8378,
13      "relation": "eq"
14    },
15    "max_score": null,
16    "hits": [
17      {
18        "_index": "dating-final",
19        "_id": "WALau4wBl3DzbDc0Ibyk",
20        "_score": null,
21        "_source": {
22          "d_funny_important": "b'[21-100]'",
23          "d_reading": "b'[6-8]'",
24          "ambtition_important": 3,
25          "d_age": 4,
26          "intellience_important": 25,
27          "movies": 8,
28          "music": 8,
29          "d_tv": "b'[0-5]'",
30          "d_expected_num_interested_in_me": "b'[0-3]'",
31          "samerace": "b'1'",
32          "shopping": 6,
33          "d_tvsports": "b'[0-5]'",
34          "d_attractive": "b'[6-8]'",
35          "d_ambition": "b'[6-8]'",
36          "dining": 10,
37          "d_shared_interests_o": "b'[0-5]'",
38          "intelligence": 7,
39          "d_funny": "b'[6-8]'",
40          "d_ambtition_important": "b'[0-15]'",
41          "d_sincere": "b'[9-10]'",
42          "field": "b'Financial Engineering'",
43          "guess_prob_liked": 9,
44          "sincere": 9,
```

- i. Find a person using certain values, and compute its number of expected matches, and the actual number of its matches

Values used in this query identify 10 entries regarding a single person, since no ID is available to identify a single person; with further data wrangling, thus with the addition of an ID, it would be possible (and easier) to perform the same query the entries about a specific person.

```

1 GET /speeddating-v4/_search
2 {
3   "query": {
4     "bool": {
5       "filter": [
6         { "term": { "gender": "female" } },
7         { "term": { "age": 24 } },
8         { "term": { "wave": 1 } }
9       ]
10    }
11  },
12  "aggs": {
13    "count_matches": {
14      "sum": {
15        "field": "match"
16      }
17    }
18  },
19  "size": 1,
20  "_source": ["id", "expected_num_matches", "count_matches"]
21 }

```

```

1 {
2   "took": 0,
3   "timed_out": false,
4   "shards": {
5     "total": 1,
6     "successful": 1,
7     "skipped": 0,
8     "failed": 0
9   },
10  "hits": {
11    "total": {
12      "value": 10,
13      "relation": "eq"
14    },
15    "max_score": 0,
16    "hits": [
17      {
18        "_index": "speeddating-v4",
19        "_id": "dlq81IwBkPnrhWx6FobH",
20        "_score": 0,
21        "_source": {
22          "expected_num_matches": 3
23        }
24      }
25    ]
26  },
27  "aggregations": {
28    "count_matches": {
29      "value": 2
30    }
31  }
32 }

```

j. Determine the number of matches in participants with a high interest correlation

```

1 GET /speeddating-v4/_search
2 {
3   "query": {
4     "range": {
5       "interests_correlate": { "gt": 0.75 }
6     }
7   },
8   "aggs": {
9     "match-count": {
10      "terms": {
11        "field": "match"
12      }
13    }
14  },
15  "size": 0
16 }

```

In the case of interest correlation greater than 0.75, returns a percentage of more or less: 18.18%, meanwhile for values between 0 and 0.75 it is: 16.73%, when between -0.75 and 0, it is: 15.35%.

```
1 {
2   "took": 0,
3   "timed_out": false,
4   "_shards": {
5     "total": 1,
6     "successful": 1,
7     "skipped": 0,
8     "failed": 0
9   },
10  "hits": {
11    "total": {
12      "value": 110,
13      "relation": "eq"
14    },
15    "max_score": null,
16    "hits": []
17  },
18  "aggregations": {
19    "match-count": {
20      "doc_count_error_upper_bound": 0,
21      "sum_other_doc_count": 0,
22      "buckets": [
23        {
24          "key": 0,
25          "doc_count": 90
26        },
27        {
28          "key": 1,
29          "doc_count": 20
30        }
31      ]
32    }
33  }
34 }
```

## 5. Extra

As extra work, to further deepen the analysis of the dataset, a dashboard with Kibana has been chosen. In particular, an overview interface has been developed, which presents a general overview of the characteristics and statistics of the dataset. A preview is presented below.



Among the views developed there are:

- The **median of the age** of participants, of which the **number** is simply shown
- **Distribution of participants**, showing the percentage of entries involving males and females. Developed with a **donut diagram**, to show the balance of the slices.
- **Distribution of participants by "race"**. We chose to represent it with a **waffle graph** to transfer the percentages to a sample of 100 participants and make the result clearer.
- The **median of preferences**, represented with a **bar histogram**, to visualize the preferences in the character of potential partners.
- The **median of interests**, also represented with a **vertical bar graph**, to visualize the most "mainstream" interests, and those of which people are less "fond". "gaming" appears to be the least interesting for the participants, while "dining", "movies", "dinner", and "music" appear to be the most interesting.
- **Matches per wave**, graph with **percentage area**, to show the total trend of match percentages (green) on the total participants per speed dating round.
- The **percentage of people who met previously**, represented with a **pie chart**, shows us that only just under 5% have already met before the experiment.



The developed dashboard presents an overview of various statistics, which can be clicked on to filter the results by aggregate categories.

For example, it is possible to filter the results by displaying only those relating to a male or female audience, from which, for example, it is possible to view a difference in the satisfaction of interests such as yoga (higher for females, second image), and gaming (higher for males, first image).

