

Part 2: Finding Associations

Thursday 18th December, 2025



In Part 1, we started looking at some trends in the data. There are fun associations that fans would likely find amusing or obvious. But there are also associations that may suggest deeper cultural norms about how certain categories of people are depicted in fiction. In this section, we look at four dimensions that speak to important demographic categories: `straight_queer`, `young_old`, `masculine_feminine`, and `rich_poor`. It is important to note that, by categorizing characters based on respondents' ratings on these dimensions, we are assessing how *perceptions* of these categories are related to *perceptions* of other dimensions.

1 Imports

```
# importing data libraries
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
import os
import plotly.express as px
import ipywidgets as widgets
from IPython.display import display
import random
from sklearn.preprocessing import StandardScaler
import finaltools as ft
import pingouin as pg
```

2 Reading the Data

```
# read processed data from part 1 notebook
char_score_data = pd.read_csv("data/processed/char_score_data.csv")
char_score_data.head()
```

[illegible]5 rows \times 467 columns

3 Association Visualization

```
# selected bap features/dimensions of interest
target_dimensions = ["straight_queer", "young_old", "masculine_feminine", "rich_poor"]
corr = char_score_data[target_dimensions].corr().round(2) # round to 2 decimals

# set up figure
plt.figure(figsize=(8, 6))

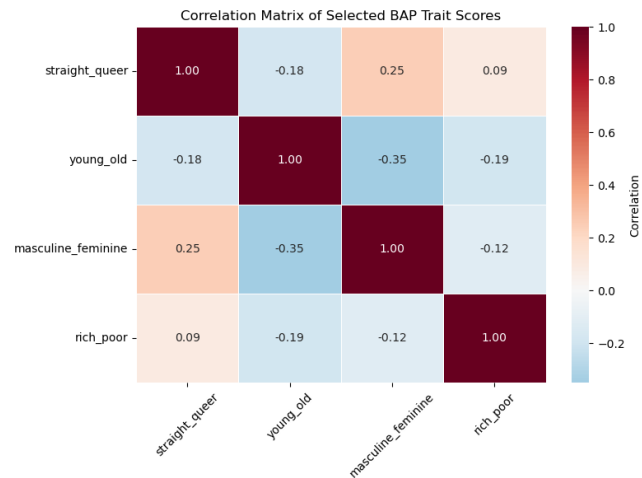
# create heatmap with diverging red -blue colors and annotations
sns.heatmap(
    corr,
    annot=True,          # show correlation values
    fmt=".2f",           # 2 decimal places
    cmap="RdBu_r",       # diverging red -blue
    center=0,            # center at 0 for +/- interpretation
    cbar_kws={'label': 'Correlation'},
    linewidths=0.5
)

# labels and title
plt.xticks(rotation=45)
plt.yticks(rotation=0)
plt.title("Correlation Matrix of Selected BAP Trait Scores")
plt.tight_layout()

# save plot to visualizations folder
plt.savefig("visualizations/selected_dim_correlation_map.png", dpi=300, bbox_inches="tight")

# show plot
```

```
plt.show()
```



We see that the features are not very strongly (pairwise) correlated with each other. This is good as we don't need to reduce dimensionality.

4 Assessing Individual Dimensions

While these dimensions are not strongly (pairwise) correlated with *each other*, they may be correlated with other dimensions in the data set. To start, we can standardize the data (through an array) and then turn it back into a **pandas** dataframe.

```
data_values = char_score_data.iloc[:, 3:]
scaler = StandardScaler()
data_values = scaler.fit_transform(data_values)
char_scores_scaled_df = pd.DataFrame(data_values)
char_scores_scaled_df.columns = char_score_data.iloc[:, 3:467].columns.to_list()
```

Next, we create a 500x500 correlation matrix. While this would be a mess to visualize, we can select our target dimensions and see if they are highly correlated with any other dimensions.

```
target_corrs = ft.target_mode_of_correlations(char_scores_scaled_df, target_dimensions)
target_corrs_df = pd.DataFrame(target_corrs)
target_corrs_df
```

	straight_queer	young_old	masculine_feminine	rich_poor
0	straight_queer (1.0)	young_old (1.0)	masculine_feminine (1.0)	rich_poor (1.0)
1	androgynous_gendered (-0.71)	historical_macho (0.61)	metrosexual_pretentious (0.68)	proletariat_bourgeoisie (-0.84)
2	hipster_basic (-0.47)	old-fashioned_progressive (-0.6)	giggling_chortling (0.56)	blue-collar_ivory-tower (-0.79)
3	macho_metrosexual (0.45)	casual_vintage (0.57)	glamorous_spandoppressed (-0.58)	privileged (-0.72)
4	normal_weird (0.43)	beautiful_ugly (0.54)	tailor_blacksmith (-0.56)	fancy_presidential (-0.68)
5	abstract_concrete (-0.42)	eternal_geriatric (0.52)	militaristic_decadent (0.56)	lowbrow_highbrow (-0.67)
6	freak_normie (-0.42)	attractive_repulsive (0.51)	person_dog-person (-0.53)	extravagant_thrifty (0.66)
7	cat-person_dog-person (-0.42)	giggling_chortling (0.51)	feminist_sexist (-0.53)	scruffy_manicured (-0.64)
8	quirky_predictable (-0.42)	apprentice_master (0.5)	scruffy_manicured (0.49)	frugal_lavish (-0.63)
9	autistic_neurotypical (-0.42)	typical_conservative (0.5)	kingaroos_dolphin (0.49)	polished_unpolished (0.63)
10	classical_avant-explorer_garde (0.39)	builder_refined (0.49)	refined_rugged (-0.47)	refined_rugged (0.61)

```
target_corrs_df.to_csv("data/target_correlations.csv", index=False)
```

By social sciences standards, there are some strong correlations here. For example, **straight_queer** is inversely correlated with **androgynous_gendered**, which suggests that the more queer a character is, the more likely their depiction is androgynous. However, these correlations are not controlling for the influence of other dimensions. We can use the **pingouin** library to calculate the partial correlations, which show us correlations between dimensions while controlling for other dimensions.

```
target_pcorrs = ft.target_mode_of_correlations(char_scores_scaled_df, target_dimensions, model)
target_pcorrs_df = pd.DataFrame(target_pcorrs)
target_pcorrs_df
```

	straight_queer	young_old	masculine_feminine	rich_poor
0	straight_queer (1.0)	young_old (1.0)	masculine_feminine (1.0)	rich_poor (1.0)
1	androgynous_gendered	old_per snapper	perf_saint	sexist_proletariat_bourgeoisie
	(-0.37)	(0.18)	(-0.25)	(-0.3)
2	macho_metrose	gamer non-gamer	macho_metrose	suppressed_privileged
	(0.13)	(0.15)	(0.22)	(-0.25)
3	musical_off-key	celebrity_boy/girl	lingling_chortle	blue-collar_ivory-tower
	(-0.1)	next-door (-0.14)	(-0.18)	(-0.2)
4	open-minded_close-minded	attractive_repulsive	person_dog-person	celebrity_boy/girl-next-door
	(-0.09)	(0.14)	(-0.18)	(0.1)
5	modest_flamboyant	juvenile_mature	chic_cheesy	frugal_lavish
	(0.09)	(0.13)	(-0.13)	(-0.1)
6	cat-person_dog-person	vibrant_geriatric	glamorous_spang	good-cook_bad-cook
	(-0.09)	(0.12)	(-0.13)	(0.09)
7	focused-on-the-present_focused-on-the-future (...)	slow-talking_fast-talking	chivalrous_businesslike	sickly_healthy
		(-0.11)	(0.12)	(0.09)
8	hugs_handshakes	cas_rock	goth_flower-child	unlucky_fortunate
	(0.08)	(0.11)	(0.11)	(-0.09)
9	kinky_vanilla	modern_historical	calpepy_disarming	entrepreneur_employee
	(-0.08)	(0.11)	(-0.11)	(0.09)
10	pronatalist_child-free	apprentice_master	narcissistic_low-self-esteem	open-minded_close-minded
	(0.08)	(0.11)	(0.1)	(0.09)

```
target_pcorrs_df.to_csv("data/target_partial_corr.csv", index=False)
```

These coefficients are far less suggestive of strong relationships. However, given how many redundant dimensions we have in the data, this might simply be an issue of too much noise and too unsophisticated of a method. We can revisit these questions after doing some dimension reduction in Part 3.