**Part 1 (Using only internal dataset):**

Determine the groups of products based on what users with similar web pages browsed. This will be done based on the set of factors below:

a. Given:

   i. Dataset 3 (Internal data gathered by company) distilled down to a list of products purchased, web pages browsed, and what each person clicked on per page. Each row corresponds to one person's data.

b. Use:

   i. Collaborative Filtering.[1]

c. To:

   i. Obtain a response matrix per product that will detail 1 if the product is frequently viewed with the product that corresponds to the row label, and 0 if not.

**Rationale:**

The rationale behind this approach is that by using collaborative filtering, we can mathematically determine correlation between each product sold/viewed and the most common items viewed alongside it. I have determined the collaborative filtering model to be the best model to use because of its simplicity to explain to the retail company, and its relevance to the problem at hand.[1]

**Part 2 (Combining multiple datasets):**

Using datasets 1 and 2 you could join them (SQL or Python) by the similar features and create a master dataset that can be sold to their alma mater from dataset 1. Based on the history of payments in dataset 2 and the financial net worth from dataset 1, we can determine a ranking of who would be most likely to donate to the university.

- a. Given:
    - ii. Master dataset as described in the previous paragraph.
- d. Use:
    - i. A clustering algorithm such as DBSCAN or OPTICS[2]
- e. To:
    - i. Obtain calculated clusters where you can determine which groups are most likely to donate based on income and past payment history.

**Rationale:** This approach is useful because we want to rank items based on criteria that involve multiple factors, and where simply sorting by a single metric would not capture the full picture. The reason I decided to choose DBSCAN or OPTICS as opposed to another clustering algorithm is because they rank based on similarities, and they do not take a predetermined amount of clusters, meaning that they will determine the clusters themselves.[2] After visualizing the data we can determine the clusters with the highest set of values as the ones most likely to donate to the university. This is because we assume they have the most disposable income. Using this generated data we can sell this to universities in the list to have them invest in pursuing these individuals.

References:
1. https://developers.google.com/machine-learning/recommendation/collaborative/basics
2. https://scikit-learn.org/1.5/modules/clustering.html