

Introduction to Machine Learning

Fall 2022

University of Science and Technology of China

Project Report	PB20061254, Yifei Zuo, 33%
Lecturer: Jie Wang	PB20061210, Yongfan Lai, 33%
Due: Jan. 11, 2023	PB20061372, Yunjin Zhen, 33%

1 INTRODUCTION

Recommending clothes of suitable sizes to customers based on the information of clothes and users are very important for E-commerce platforms. In this project, we implement several classifiers to predict customers' fit feedback based on a dataset collected from RentTheRunWay.

2 PREPARATION

The training dataset contains 87766 samples, each of which has 14 features and 1 label. We summarize these records as follows:

(1) Item attributes, i.e. `item_name`, `size` and `price`; (2) User attributes. The first one is `user_name`, which is excluded in the test set. The others describe the body characteristics of each user: `age`, `height`, `weight`, `body_type` and `bust_size`; (3) Transaction attributes, i.e. `rented_for`, `usually_wear`; (4) Feedback, i.e. `fit`, `review_summary` and `rating`, among which `fit` is the target variable we want to predict, and the other three are supposed to be inaccessible on the test set.

By observing that most of the inputs has missing values and inconsistent formats, we need to design a thorough data cleansing (and transforming) pipeline that converts the raw data into either categorical or numerical variables we can leverage for training. We also need to deal with the data imbalance issue, since the number of `True` to `Size` samples, consisting of 70% of the whole dataset, is much larger than the other two classes. These two challenges will be discussed in detail in the following sections.

2.1 Data Cleansing

We load the provided json file into memory as tabular data, with each rows being a training sample and each column representing features or labels. Before diving into each separate input column, we first drop the 190 samples corrupted by the byte order mark `\udfff`. Moreover, as the labels are encoded differently in the training and test sets, we unify them to integers 1, 2, 3 for `Small`, `True` to `Size` and `Large`, respectively. This encoding approach captures the ordinal nature of the labels, which is important for the design of our model.

2.1.1 Numerical Features

There are 5 input columns that can be considered purely numerical: `price`, `age`, `height`, `weight`, `usually_wear` and `rating`. We remove the prepended dollar sign from `price`, convert the unit of `height` from feet and inches to centimeters, convert the unit of `weight` from pounds to kilograms, and treat them as float point numbers along with the other three columns.

We also notice that there are some anomalous values, i.e. `age` ≤ 5 or ≥ 100 , `height` > 200 cm and `weight` < 30 kg or > 150 kg, which are more likely to be caused by typos or data corruption than real values. To avoid the harmful effect of these outliers, we set them as `NaN` and replace them with the corresponding median values in later steps.

The challenge is to handle the `bust_size`. We observe that a valid `bust_size` always contains two part: (1) number part in inches; (2) letter part, implying the `cup_size`. Since these two parts are different measurements of women's busts, it is necessary to split them into 2 features. The `cup_size`, ordered as `AA < A < B < C < D < D+ < DD < DDD/E < F < G < H < I < J`, are ordinarily encoded as 0 ~ 12. We now have `bust_size` in 2 numerical values.

After enclosing the 4 numerical features that measures user's body characteristics approximately follow a normal distribution, as is shown in the figure below. We will normalize them to have zero mean and unit variance in the later steps.

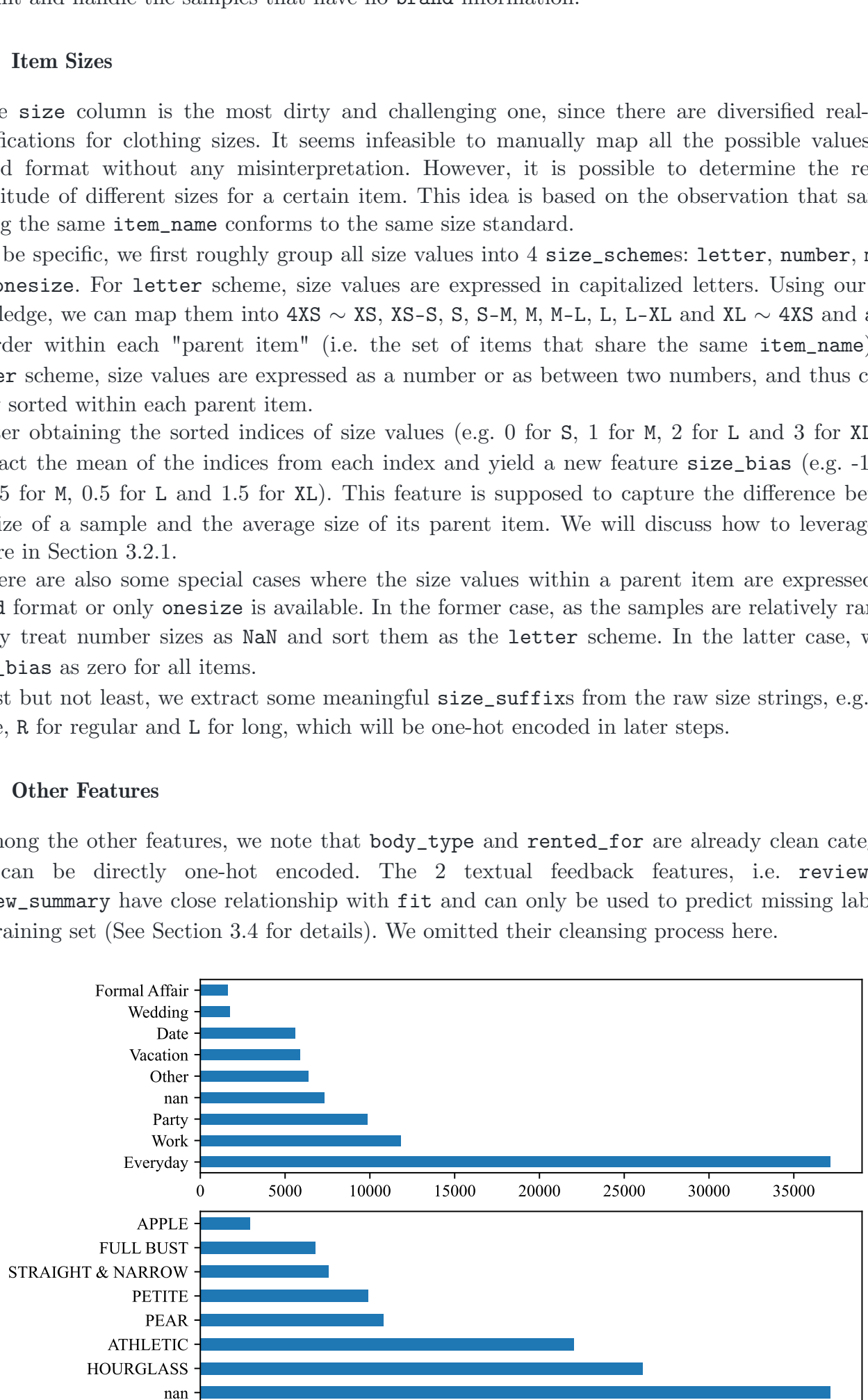


Figure 1: Violin plots of `weight`, `height`, `bust_size` and `cup_size`

2.1.2 Item Names

Next, we take a look at the categorical features. For `item_name`, we notice that it demonstrates some common patterns and can be utilized to extract extra information: that is, the first line of `item_name` typically contains the brand name, while the second line contains the product name with the last word indicating its category, e.g. `Jumpsuit`, `Romper`, `Skirt`, `Top`, `Blouse`. Therefore, we parse each `item_name` string and yield two new features: `brand` and `category`. We also take into account and handle the samples that have no `brand` information.

2.1.3 Item Sizes

The `size` column is the most dirty and challenging one, since there are diversified real-world specifications for clothing sizes. It seems infeasible to manually map all the possible values to a unified format without any misinterpretation. However, it is possible to determine the relative magnitude of different sizes for a certain item. This idea is based on the observation that samples having the same `item_name` conforms to the same size standard.

To be specific, we first roughly group all size values into 4 `size_schemes`: `letter`, `number`, `mixed` and `onesize`. For `letter` scheme, size values are expressed in capitalized letters. Using our prior knowledge, we can map them into `XS < XL`, `XS < S`, `S < M`, `M < L`, `L < XL` and `L < XS` and assign an order within each "parent item" (i.e. the set of items that share the same `item_name`). For `number` scheme, size values are expressed as a number or as between two numbers, and can be easily sorted within each parent item.

After obtaining the sorted indices of size values (e.g. 0 for S, 1 for M, 2 for L and 3 for XL), we subtract the mean of the indices from each index and yield a new feature `size_bias` (e.g. -1.5 for S, -0.5 for M, 0.5 for L and 1.5 for XL). This feature is supposed to capture the difference between the size of a sample and the average size of its parent item. We will discuss how to leverage this feature in Section 3.2.1.

There are also some special cases where the size values within a parent item are expressed in a `mixed` format or only `onesize` is available. In the former case, as the samples are relatively rare, we simply treat number sizes as `NaN` and sort them as the `letter` scheme. In the latter case, we set `size_bias` as zero for all items.

Last but not least, we extract some meaningful `size_sufffix` from the raw size strings, e.g. `P` for `petite`, `R` for `regular` and `L` for `long`, which will be one-hot encoded in later steps.

2.1.4 Other Features

Among the other features, we note that `body_type` and `rented_for` are already clean categories and can be directly one-hot encoded. The 2 textual feedback features, i.e. `review` and `review_summary` have close relationship with `fit` and can only be used to predict missing labels in the training set (See Section 3.4 for details). We omitted their cleansing process here.

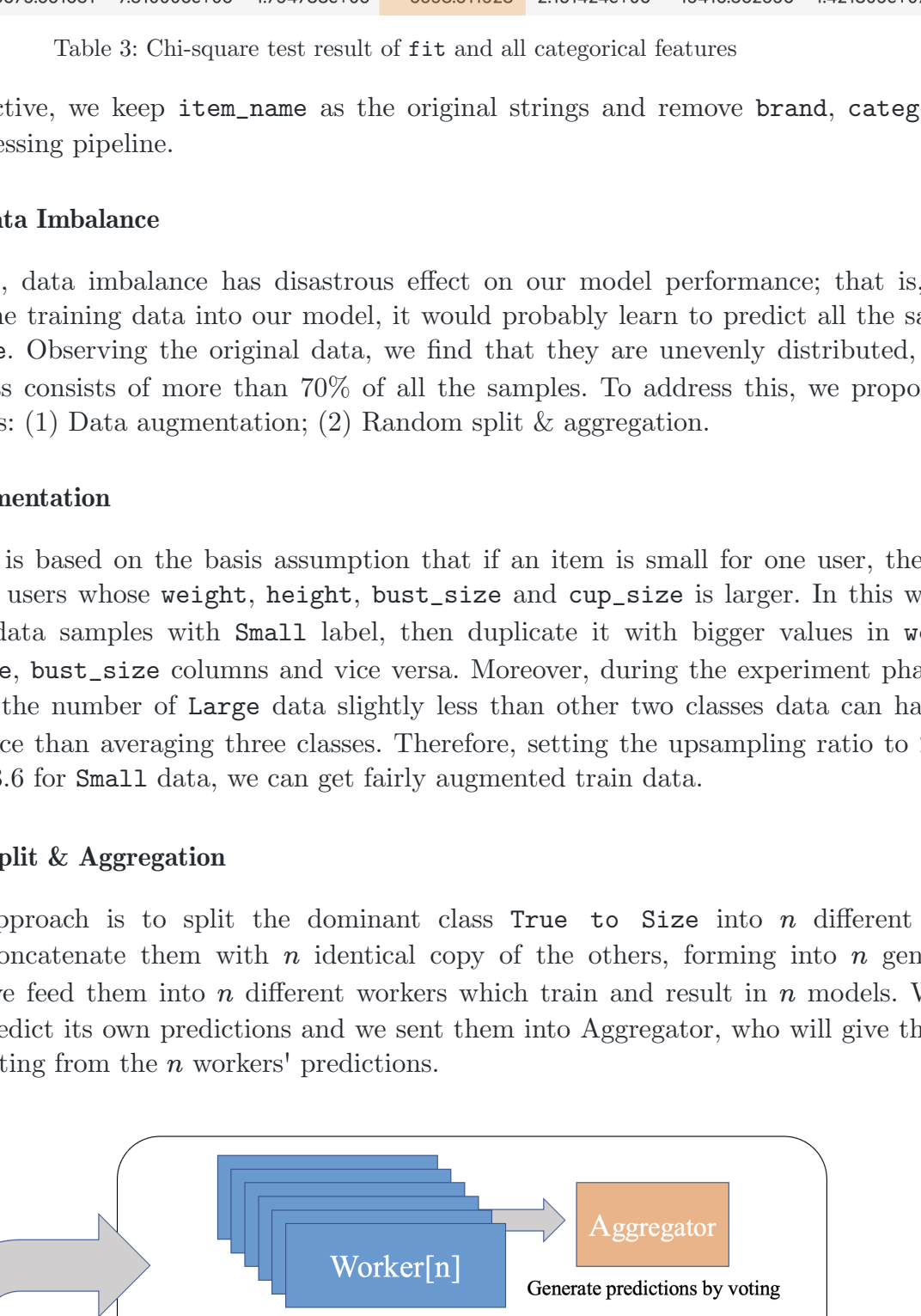


Figure 2: Bar plots of `rented_for`, `body_type` and `size_main` (size without suffix)

2.2 Exploratory Data Analysis

2.2.1 Numerical Features

Based on the current understanding of all data format, we measure the relationship among all features to further explore the dataset. First of all, we focus on numerical features. According to the literal meaning, features like `price`, `age` are less likely to impact the `fit` feature, and by calculating their correlation index, we found that actually none of these numerical features are highly related to `fit`. We also analyze the possible combination of features, e.g. `bm1` in the first figure below.

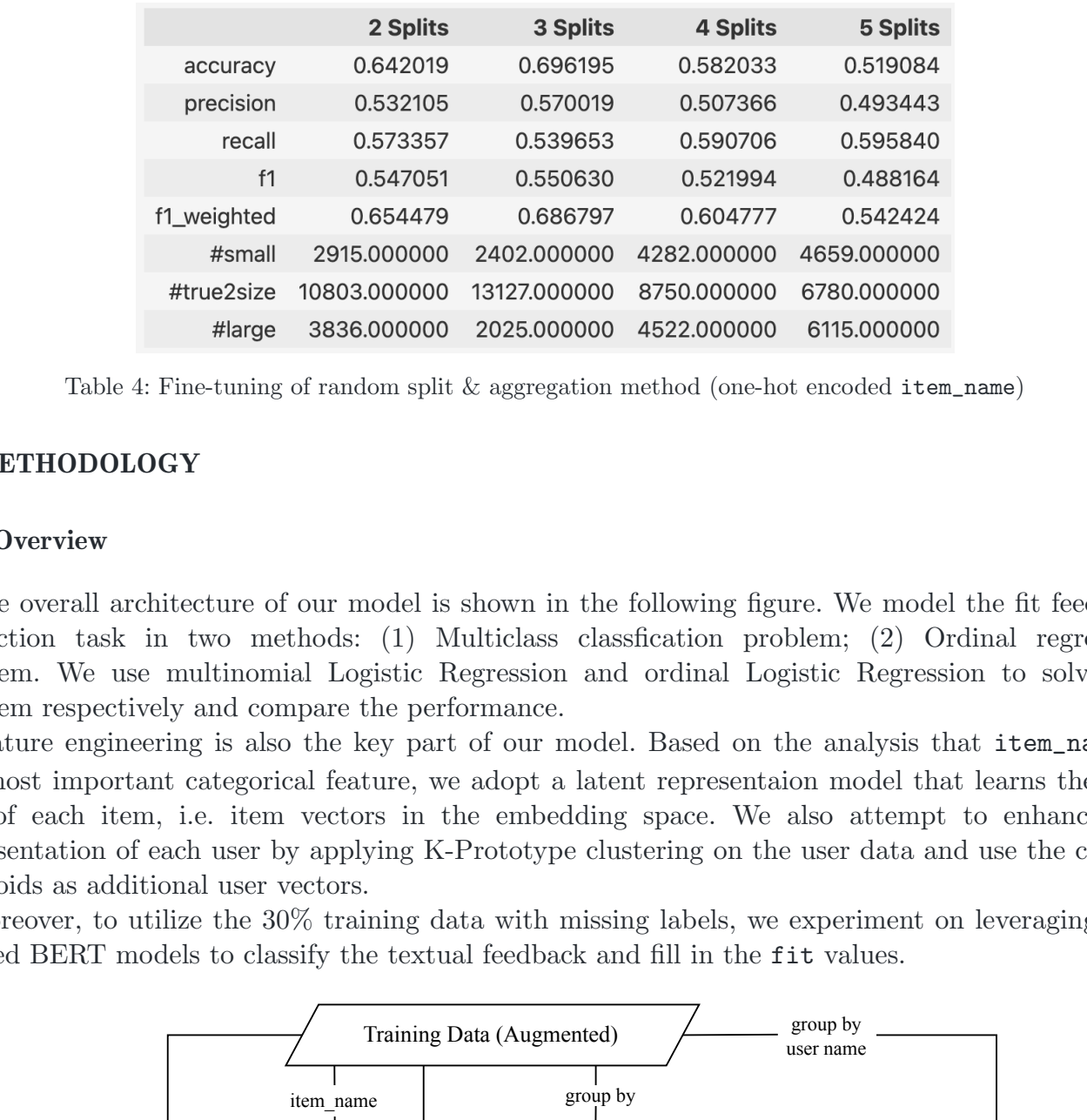


Figure 3: Heatmap of correlation between numerical features

Trivial these numerical columns are, we still try to exploit latent relationship from it, because just discarding them seems brute. The desire to gain further insight of the data results in the Detecting User Prototypes phase in our methodology, which will be explained in section 3.2.2.

2.2.2 Categorical Features

Then, we measure the categorical features. Note that we still keep the original `item_name` despite we have split it into `brand` and `category`. This is because we are not sure yet the performance of splitting. We first analysis the distribution of `fit` value inside each of these categories, e.g. in `rented_for` and `body_type`. Meanwhile, we here keeps the empty values because sometimes empty values could be seen as a category and resulting more robust performance than imputing them.

	Work	Vacation	Date	Everyday	Party	Other	Formal Affair
Small	0.14108	0.122323	0.169825	0.110612	0.124058	0.158237	0.180494
True to Size	0.701001	0.701400	0.719878	0.698057	0.728147	0.677273	0.637841
Large	0.148792	0.157347	0.110297	0.191331	0.110725	0.119315	0.161607

	APPLE	FULL BUST	STRAIGHT & NARROW	PETITE	PEAR	APPLE
Small	0.140480	0.135484	0.135347	0.125928	0.145990	0.150924
True to Size	0.686929	0.695950	0.719878	0.707659	0.684625	0.702852
Large	0.162287	0.168990	0.147651	0.156658	0.180251	0.165456

Table 1: Proportion of fit in each `rented_for` category

	HOURLGLASS	STRAIGHT & NARROW	ATHLETIC	PETITE	FULL BUST	PEAR	APPLE
Small	0.140480	0.135484	0.135347	0.125928	0.145990	0.150924	0.150905
True to Size	0.686929	0.695950	0.719878	0.707659	0.684625	0.702852	0.686929
Large	0.162287	0.168990	0.147651	0.156658	0.180251	0.165456	0.178365

Table 2: Proportion of fit in each `body_type` category

We notice that except for `item_name` and its derivative category, the distribution inside each one do not have manifest difference, which implies that these features are still not the dominant one in terms of predicting `fit`. While `item_name` is too diverse and dirty for us to analyze in this way. So we measure the direct relationship of `fit` and all of these features by Chi-square test. The result matches our expectation from previous observation, `item_name` is the key point and it is even more important than the features we extracted from it.

col2	body_type	brand	category	fit	item_name	rented_for	size
col1	body_type	brand	category	fit	item_name	rented_for	size
body_type	614362.000000	5.188996e+03	7.800388e+02	104.291123	3.221431e+04	957.951907	1.567955e+04
brand	5198.889594	5.439463e+07	1.130522e+06	9545.844958	4.247737e+07	4.2839.006624	7.310006e+05
category	7.800388e+02	5.439463e+07	6.480624e+04	4.2136.858114	4.247737e+07	4.2839.006624	7.310006e+05
fit	104.291123	3.221431e+04	957.951907	1.567955e+04	957.951907	1.567955e+04	1.567955e+04
item_name	3.221431e+04	4.247737e+07	6.480624e+04	4.2136.858114	4.247737e+07	4.2839.006624	7.310006e+05
rented_for	957.951907	4.2839.006624	7.310006e+05	1.567955e+04	1.567955e+04	1.567955e+04	1.567955e+04
size	1.567955e+04	7.310006e+05	1.567955e+04	1.567955e+04	1.567955e+04	1.567955e+04	1.567955e+04

Table 3: Chi-square test result of fit and all categorical features

In this perspective, we keep `item_name` as the original strings and remove `brand`, `category` in our final preprocessing pipeline.

2.3 Handling Data Imbalance

In this project, data imbalance has disastrous effect on our model performance; that is, if we simply feed all the training data into our model, it would probably learn to predict all the samples as `True` to `Size`. Observing the original data, we find that they are unevenly distributed, where the majority class consists of more than 70% of all the samples. To address this, we propose two alternative tactics: (1) Data augmentation; (2) Random split & aggregation.

2.3.1 Data Augmentation

This approach is based on the basis assumption that if an item is small for one user, then it is also small for all users whose `weight`, `height`, `bust_size` and `cup_size` is larger. In this way, we randomly fetch data samples with `Small` label, then duplicate it with bigger values in `weight`, `height`, `cup_size`, `bust_size` columns and vice versa. Moreover, during the experiment phase, we find that letting the number of `Large` data slightly less than the other two classes data can have the better performance than averaging three classes. Therefore, setting the upsampling ratio to 2.7 for `Large` data and 3.6 for `Small` data, we can get fairly augmented train data.

2.3.2 Random Split & Aggregation

The second approach is to split the dominant class `True` to `Size` into n different group randomly, not considering them with an identical copy of the others, forming into n generated datasets. Then we feed them into n different workers which train and result in n models. We use each model to predict its own predictions and we sent them into Aggregator, who will give the final predictions by voting from the n workers' predictions.

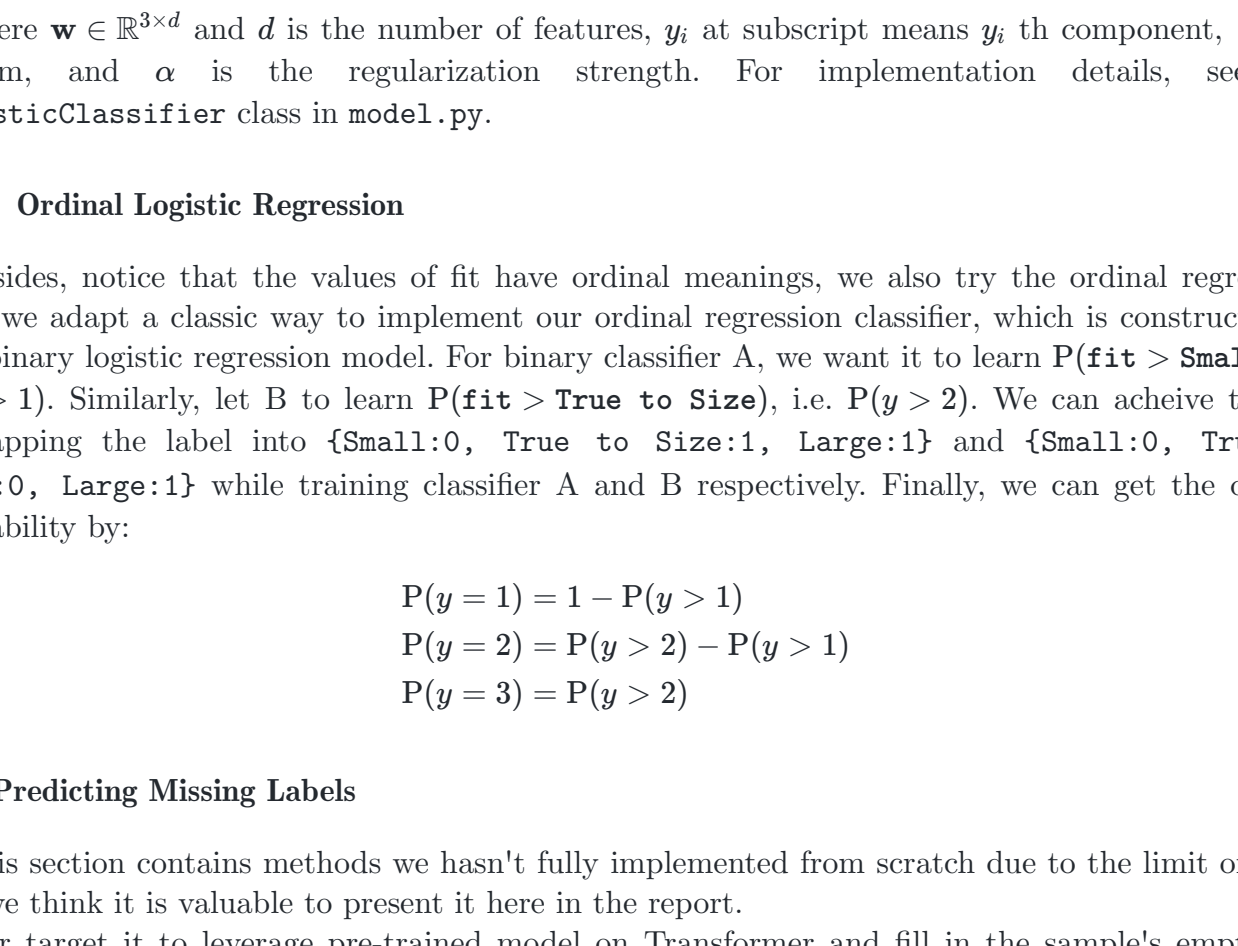


Figure 4: Random split & aggregation method

Noting that we don't necessarily have to split the dominant class into n samples so that the number of samples in each group is approximately identical to the others. Due to the fact that this task is relatively hard and our model won't give a tremendously improving representation, we do need to fine tune the parameter n . A reasonable assumption is that by relatively giving more samples in dominant class, we result will be better because there will be more samples fall into this class. Here's some result of parameter search trained on `LogisticRegression` model: (Noting that we could not reach such performance by simply do an hot encoding due to the fact that we may encounter new value of `item_name`, this is just a demonstration of how fine tune this parameter.)

	2 Splits	3 Splits	4 Splits	5 Splits
accuracy	0.642019	0.696195	0.582033	0.619194
precision	0.523205	0.570709	0.507566	0.634443
recall	0.573357	0.539653	0.590706	0.595840
f1	0.547051	0.550630	0.521994	0.608164
f1_weighted	0.654479	0.686797	0.604777	0.642424
#small	2915.000000	2402.000000	4282.000000	4680.000000
#true2size	10803.000000	13127.000000	8750.000000	6769.000000
#large	3836.000000	2025.000000	4522.000000	6115.000000

Table 4: Fine-tuning of random split & aggregation method (one-hot encoded `item_name`)

3 METHODOLOGY

3.1 Overview

The overall architecture of our model is shown in the following figure. We model the fit feedback prediction task in two methods: (1) Multiclass classification problem; (2) Ordinal regression problem. We use multinomial Logistic Regression and ordinal Logistic Regression to solve the problem respectively and compare the performance.

Feature engineering is also the key part of our model. Based on the analysis that `item_name` is the most important categorical feature, we adopt a latent representation model that learns the true size of each item, i.e. item vectors in the embedding space. We also attempt to enhance the representation of each user by applying K-Prototype clustering on the user data and use the cluster centroids as additional user vectors.

Moreover, to utilize the 30% training data with missing labels, we experiment on leveraging pre-trained BERT models to classify the textual feedback and fill in the fit values.

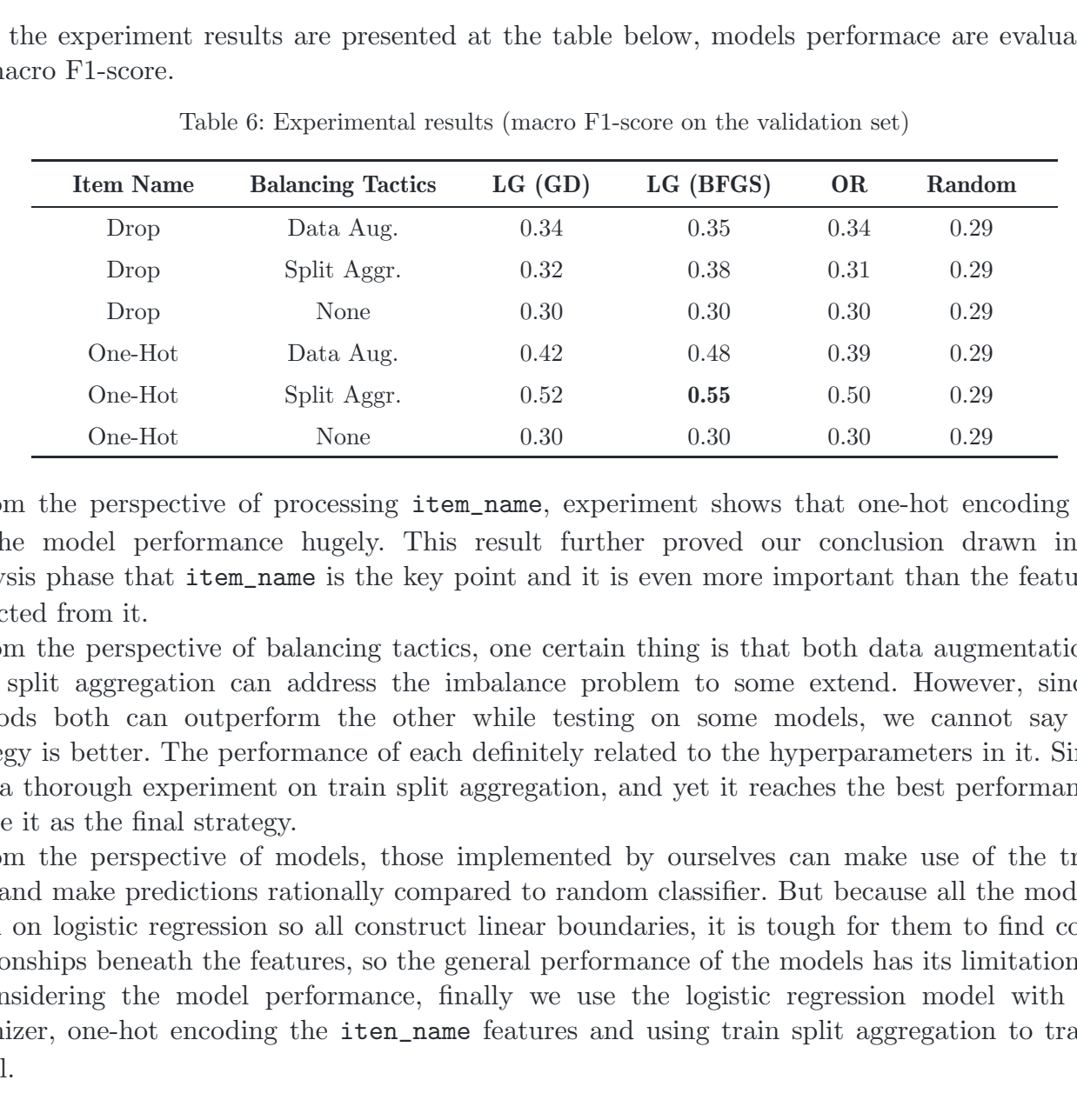


Figure 5: The architecture of our proposed model.

Our model is demonstrated to improve the performance when experimented on each separate component. Unfortunately, due to the limit of time, energy and resource, we are not able to fully implement the algorithms and integrate them all with only whitelist libraries.

3.2 Feature Learning

3.2.1 Leveraging Item Sizes

In the previous sections, we have transformed the item sizes into `size_bias` feature which captures the difference between the item size and the average size of its parent item. The true size of an item can be viewed as a linear function of its size bias. Denote the item, the user, and the parent item of a transaction t by i , u , and p , respectively. We model the true size of the item i as

$$v_i = v_p + \epsilon_i d_p,$$

where v_i is the true size of its parent item, ϵ_i is the bias of item size relative to the parent item, and d_p describes the heterogeneous influence of size bias. We want to learn v_p and d_p for each parent item p . A simple approach for estimating v_p is to use the mean of the true sizes of all users that have rent the parent item with $\epsilon_i = 0$. However, the item can either be too small or too large for some users, leading to noisy results. In order to learn the latent representation of item size more precisely, we drop all irrelevant features and solve the following ordinal regression problem:

$$\min_{\mathbf{w}, b_1, b_2} \sum_{i \in T} \ell(y_i, f(\mathbf{w}, \mathbf{v}_i), b_1, b_2),$$
$$\text{s.t. } \mathbf{w} \geq \mathbf{0}, d_p \geq \mathbf{0},$$

where we define $f(\mathbf{w}, \mathbf{v}_i) = \mathbf{w}^T (\mathbf{v}_i - \mathbf{v}_p)$ as the fitness score between the item and the user, and $\ell(y_i, f(\mathbf{w}, \mathbf{v}_i), b_1, b_2)$ is the total loss of 2 binary logistic classifiers sharing the same weights \mathbf{w} :

$$\ell(y_i, f(\mathbf{w}, \mathbf{v}_i), b_1, b_2) = \begin{cases} \log(1 + e^{f(\mathbf{w}, \mathbf{v}_i) - b_1}), & \text{if } y_i = 1; \\ \log(1 + e^{-f(\mathbf{w}, \mathbf{v}_i) - b_2}), & \text{if } y_i = 2; \\ \log(1 + e^{-f(\mathbf{w}, \mathbf{v}_i) - b_3}), & \text{if } y_i = 3. \end{cases}$$

Note that, in Prob. (1), we add the constraint $d_p \geq \mathbf{0}$ to ensure the monotonicity of the true size v_i with respect to the size bias ϵ_i , and also $\mathbf{w} \geq \mathbf{0}$ to capture the monotonicity of the predicted fitness score with respect to the size bias ϵ_i ; that is, if a size is larger for some user, then its predicted fitness score should be larger. These two constraints are trivial for the Projected Gradient Descent (PGD) algorithm, since the projection operator onto the non-negative orthant is a simple element-wise $\max\{0, \cdot\}$.

Therefore, we implement PGD algorithm to optimize the objective in rounds: (1) Initialize \mathbf{w} , b_1 , b_2 randomly; (2) In each round, we fix the parameters \mathbf{w} , b_1 , b_2 and optimize v_p and d_p for each parent item p ; (3) In each round, we fix v_p and d_p and optimize \mathbf{w} , b_1 , b_2 . We repeat the process for 10 rounds and obtain the final parameters \mathbf{w}^* , d_p^* . See the `ItemWeightOptimizer` class in `preprocess.py` for more details.

3.2.2 Detecting User Prototypes

In this section, we will discuss the technique we use to enhance the representation of each user. Although the training data already contains the true size of each user, i.e. `weight`, `height`, `bust_size` and `cup_size`, these features are unable to express the heterogeneity among different users (e.g. personal preference), let alone considerable amount of missing values.

Following the same logic as we learn the latent variables for item sizes, it is reasonable to group the transactions by each user and learn user-specific features. However, unlike the `item_name`, the column `user_name` cannot uniquely identify a user. Hence, instead of learning the exact latent representation of each user, we attempt to learn several user prototypes by performing clustering algorithms on the user data, and use the cluster centroid as the user vector.

Besides numerical body measurements, our clustering algorithm should be aware of categorical features, i.e. `user_name`, `body_type` and `rented_for`, so as to capture user preferences. We propose using K-Prototype algorithm, which is a hybrid model that combines K-Means for numerical and K-Mode for categorical data. Basically, the algorithm calculates the distance between a data point and a cluster centroid by summing the Euclidean distance between the numerical features and the Hamming distance between the categorical features. By experimenting with different values of cluster number n , we find that: as n increases, the resulting user vectors slightly improve our model performance on both the training and validation set. The classification scores are shown in the following table:

	accuracy	precision	recall	f1	f1_weighted
train (n=500)	60.10%	59.75%	60.10%	59.82%	59.82%
train (n=300)	59.99%	59.62%	59.62%	59.74%	59.74%
train (n=200)	60.02%	59.62%	59.99%	59.70%	59.70%
train (n=400)	59.94%	59.59%	59.94%	59.67%	59.67%
train (n=100)	59.90%	59.59%	59.90%	59.63%	59.63%
train (user vector)	59.83%	59.47%	59.83%	59.54%	59.54%
test (n=500)	49.16%	48.04%	53.30%	45.24%	51.88%
test (n=200)	48.87%	46.03%	53.34%	45.11%	51.59%
test (n=400)	49.02%	45.86%	53.06%	45.09%	51.73%
test (n=300)	48.84%	45.82%	53.03%	44.98%	51.62%
test (user vector)	48.61%	45.85%	53.13%	44.89%	51.36%
test (n=100)	48.66%	45.66%	52.73%	44.74%	51.42%

Table 5: Fine-tuning of K-Prototype clustering (one-hot encoded `item_name`)

These results are satisfactory if we ignore the computational cost of the algorithm. Unfortunately, it cost hours to cluster on the entire dataset using the existing `kmodes` library even after enabling parallel computing. Since the improved F1-score on the validation set is not significant, we decide not to implement this technique in our final model.

3.3 Fit Feedback Classification

In preceding part, we already have a penetrating insight of our tasks. As for the final models, considering the model performance and the implementing difficulties, we choose the Logistic Regression model to implement.

3.3.1 Multinomial Logistic Regression

At first, we implemented the classifier with gradient descent algorithm. However, its performance not so good as the model from `scikit-learn` library. So we imitate the sklearn implementation, using the BFGS algorithm in `scipy.optimize` to minimize the multiclass loss function, and get another form of the classifier. The loss function is formulated as:

$$\ell(\mathbf{w}, \alpha) = -\frac{1}{n} \sum_{i=1}^n \text{softmax}(z_i, \mathbf{w}) \cdot \mathbf{y}_i + \frac{1}{2} \alpha \|\mathbf{w}\|^2$$

where $\mathbf{w} \in$