

## **BNPL Presentation Script**

### **Introduction and team**

**RON**

**[SLIDE 1]**

Good afternoon, we are team 8, researching solutions for BNPL partner merchant rankings. I am Ron, **[SLIDE 2]** and our team consists of Cindy, Anderson, Henry and Yujie.

**[SLIDE 3]**

Our main task is to conclude a top 100 merchants list for a generic BNPL company to onboard. In the first three weeks we focused on fraud detection, missing data imputation and merchant segmentation. Then, we engineered variables and developed our final ranking algorithm.

**[SLIDE 4]**

We each contributed an average of 15-20 hours per week, and on top of regular in-person meetings employed Trello, Github and Zoom for remote collaboration.

### **Introduce business problem**

**ANDERSON**

**[SLIDE 5]**

The main task we had was to assist a Buy Now, Pay Later firm in the selection of 100 merchants. As a firm that has just begun offering a Buy Now, Pay Later feature, we believe its business goal should be to earn profit while maintaining a long-term business establishment. Hence, we believe that the size of cashflows a merchant can bring for the BNPL firm cannot be the sole criterion when picking business partners.

High returns may come with high risks, merchants could make a fortune one month but go low for the rest of the year, which increases financial risk for the BNPL firm. As a newly established firm may not have a strong cash reserve, this may threaten survival.

Also, we believe that the newly established firm may not have a large and strong consumer base that uses its services. Hence, we are interested in selecting merchants with a consumer base that has the potential to provide long term expansion for the firm.

Thus, our solution aims to select merchants that benefit the firm in terms of financial stability and development.

## Fraud

CINDY

### [SLIDE 6]

Fraud is critical when choosing partner merchants since **high fraud rate** leads to **significant losses** for the BNPL firm.

To deal with this problem, we built a fraud rate predictor using provided data, and then eliminated likely fraudulent transactions.

### [SLIDE 7]

Features we engineered for modelling includes simple aggregations 【动作指导：指向屏幕】 , but also the ratio of them divided by their historical averages. For example, the first variable in the table is the dollar value per order spent by a consumer on a suspected day of fraud divided by their own average on other days. This is good for capturing deviance from the consumer's own typical behaviour.

### [SLIDE 8]

We ran our RFR model over all transactions data, and eliminated the **top 0.1%** of around 16,000 transactions using domain knowledge. The effective boundary for fraud labelling is thus 20% using our fraud predictor.

## **Clustering**

**HENRY**

### **[SLIDE 9]**

Clustering is our solution to two problems. First, the dataset was missing take rate information for certain merchants, and we had to impute the take rate to be able to consider them in the final 100 or else must discard potentially valuable merchants for our firm. We decided to use within-cluster mean-imputation to get accurate take rate estimates. Second, we had to define 3 to 5 merchant segments according to the business requirements, which we obtained using clusters.

### **[SLIDE 10]**

We built our clustering models with **monthly revenue, monthly number of orders and historical total distinct customers** as our features, which quantifies the scale of a business. 【动作指导：指向屏幕】

We log-transformed all features due to the right skewness of features values, and ultimately used a Gaussian Mixture Model to form 3 clusters, in which we assumed that our data were sampled from a mixture of Gaussian distributions.

### **[SLIDE 11]**

As we can see from the 3D scatter plot, the transformed data forms a bridge-like shape in 3D space, and the separation between clusters are well defined.

### **[SLIDE 12]**

We can roughly think of the segmentation as three business size levels: small, medium, and large.

## **Finance theory + ranking algorithm introduction**

**RON**

### **[SLIDE 13]**

Now to our ranking algorithm.

### **[SLIDE 14]**

To achieve our business goal of stability, we borrow the financial metric of Sharpe Ratio, which quantifies risk compensation. For our project, we borrow this concept to calculate a metric defined as the BNPL firm's

incoming cashflows divided by its standard deviation over different fortnightly periods.

We also borrow more financial concepts by considering the 100 selected merchants as an ‘investment portfolio’ and further define the market as the portfolio of all our selectable merchants.

Our aim is thus quantified as finding 100 merchants with maximum future portfolio sharpe ratio, or effectively the maximum future benefit score.

### **[SLIDE 15]**

Our algorithm to achieve this goal is quite intuitive: it first scores merchants using a heuristic function and then ranks by their scores, taking the top 100 to form our recommended portfolio.

The optimal parameters are found using a tuning process that interplays two blocks of chronologically split data to ensure our model is future looking, with the feedback mechanism finding the set of parameters which maps to the greatest future optimum. By exploring the mappings using different parameter combinations, the heuristic function learns the characteristics that constitute a good merchant for the BNPL portfolio.

### **Heuristic function features**

#### **CINDY**

### **[SLIDE 16]**

We’ll now discuss our heuristic function features

### **[SLIDE 17]**

The first two features account for the cashflows and cashflow risks of a merchant. However, we scaled the standard deviation feature by dividing it on the first variable to avoid overrepresentation of the mean.

These two variables capture fraud as removed transactions affect cashflows which penalises high fraud rate merchants.

Correlation to the market portfolio was included because it is regularly used by financial professionals for stock market analysis.

### Repeated purchase rate

**ANDESON**

**[SLIDE 18]**

Our fourth feature is the repeated purchase rate of a merchant, defined by this formula. This metric increases if the merchant has many repeated customers, but the BNPL firm should prefer to include merchants where this metric is low into its portfolio because 1. Repeated customers tend to make conservative purchases which is not beneficial to cashflow growth; and 2. More distinct customers bring implicit benefit of unpaid advertisement.

### Persona

**YUJIE**

**[SLIDE 19]**

We also explored the customer profiles and derived a persona score for each candidate merchant to capture whether such merchant's consumers can provide us with long-term development opportunities.

We created these variables from aggregating Australian Census Data by postcode, and then proportionated them for each merchant based on the postcode frequency of their consumers.

**[SLIDE 20]**

Components of this score includes age and students, which the firm would desire to be young and highly educated because first, young people tend to spend more impulsively, second, they are more open-minded to new technology Third, highly educated students are likely to be the next generation of high income consumers.

### **[SLIDE 21]**

The score also considers income, which is desired to be high because high-income consumers are more capable to repay regularly, reducing business risk.

### **[SLIDE 22]**

We also include the postcode population as a scaling factor

Hence, our BNPL firm should prefer merchants with high persona scores in its portfolio.

## **Result**

### **[SLIDE 23]**

Before we present our recommended merchants, we would like to first analyse their characteristics.

### **[SLIDE 24]**

We first see that the two fundamental characteristics which relate to our aim of stable growth: fortnightly mean cashflow and its scaled sd, shows significant differences between portfolio and non-portfolio merchants. Specifically, the within portfolio median of fortnightly mean cashflow at around 260k is much higher than that of the others at approximately 21k.

In contrast, the scaled sd of our portfolio merchants have a very narrow distribution with median of 0.21, whilst the remaining merchants' are much higher at 0.46.

These cross sectional glimpses at statistics of the final portfolio indicate that the algorithm was successful in achieving the stability goal of high cashflow at low sd which would give our portfolio a high future benefit score.

## **[SLIDE 25]**

We also analysed the effectiveness of the last two features, by comparing our current top portfolio to a top portfolio ranked using a heuristic function of just the first three cashflow based features.

We notice that, compared to the 3 featured portfolio, the 5 featured portfolio has more merchants with lower per order transaction amount, but also included some selected high per order amount merchants. This is the typical characteristic of a financial portfolio, with most merchants stable and some risky ones where risk is compensated by fair reward. This demonstrates that including repeated customer rate and persona score has enhanced our model.

Overall, it can be seen that a portfolio most beneficial to stability and growth should contain a majority of merchants with small order quantity and low per order transaction amount, while having around 20% of high per order transaction amount merchants.

## **[SLIDE 26]**

### **ANDERSON**

We now confidently recommend - using this optimal heuristic function - to onboard this list of 100 merchants. We are especially pleased that 11 firms which were kept in consideration thanks to our segmented mean-imputation made it into our final portfolio, many being highly ranked. Please refer to the appendix for the rest of our portfolio.

## **[SLIDE 27, 28, 29]**

This is the merchant ranking by segment and also the proportions of each segment in the final portfolio. As we can see from the pie chart, medium sized merchants made up the majority of the portfolio. The ranking by segment can also be found in the appendix.

### **Limitation**

## **[SLIDE 30]**

The greatest limitation of the model is the potential for overfit as we had repeated use of data for tuning and recommendation - though manipulated differently. However, we believe the current process is the best solution to the problem as further splitting the 1.5 years of data we currently have to avoid overfit will cause inaccurate estimates of feature values.