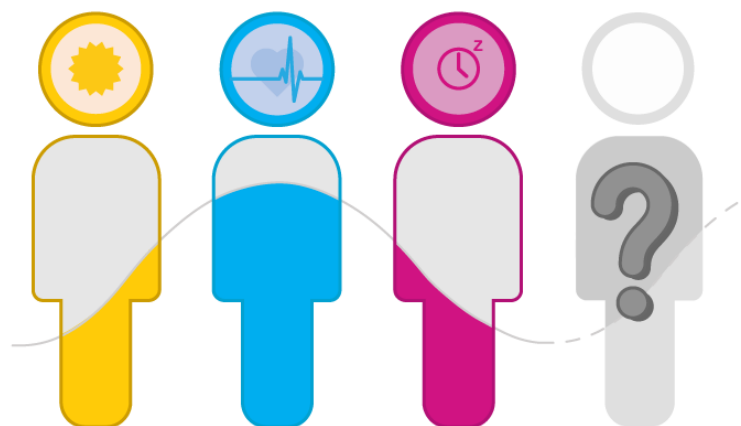




User Cancellation Modelling: *on Clustering of Customer Behaviours*



Victor Wang

Supervisors: Andrew Mellor, Junaid Mubeen

Agenda



- Motivation and Goal
- Background – Business Settings at Whizz
- Model
 - ❑ Customer Journey
 - ❑ Methodology
- Result
 - ❑ Clusters
 - ❑ Feature Analysis
- Conclusion

Motivation and Goal



Motivation



Motivation and Goal



Motivation



Goal – **Churn** (Cancellation) modelling

Who?

- ✓ Identify customers most prone to cancel subscription
- ✓ Assess likelihood

Why?

- ✓ Analyse critical reason triggering cancellation
- ✓ Make bespoke retention policy

Business Settings at Whizz



Cancellation Mechanism



Contractual

- Pupils subscribe to access Whizz products on a 1-month or 1-year contract



Voluntary

- Subscribers make the choice to leave the service at the end of the subscription; otherwise auto-rolled

Business Settings at Whizz



Cancellation Mechanism



Contractual

- Pupils subscribe to access Whizz products on a 1-month or 1-year contract

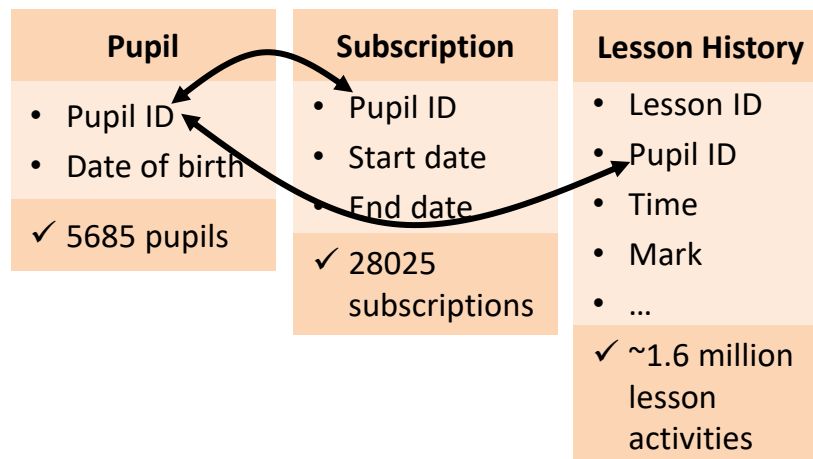


Voluntary

- Subscribers make the choice to leave the service at the end of the subscription; otherwise auto-rolled

Data Records

Time period: 2014-Jan-01 – 2018-Apr-20



Business Settings at Whizz



Cancellation Mechanism



Contractual

- Pupils subscribe to access Whizz products on a 1-month or 1-year contract

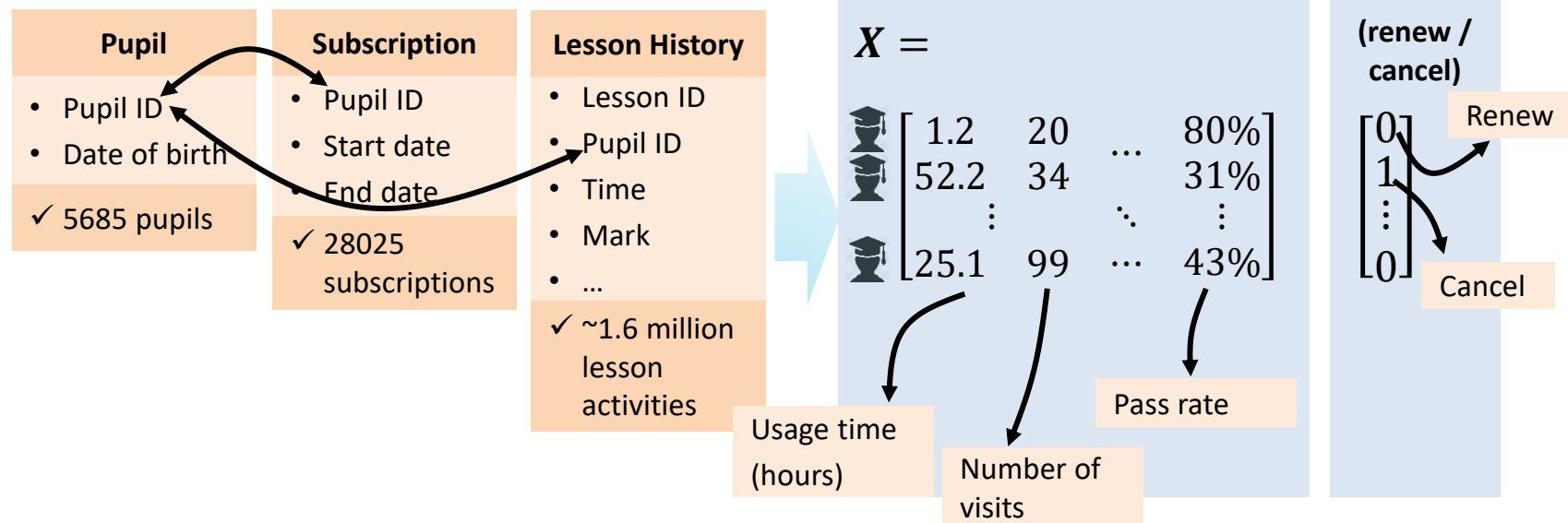


Voluntary

- Subscribers make the choice to leave the service at the end of the subscription; otherwise auto-rolled

Data Records

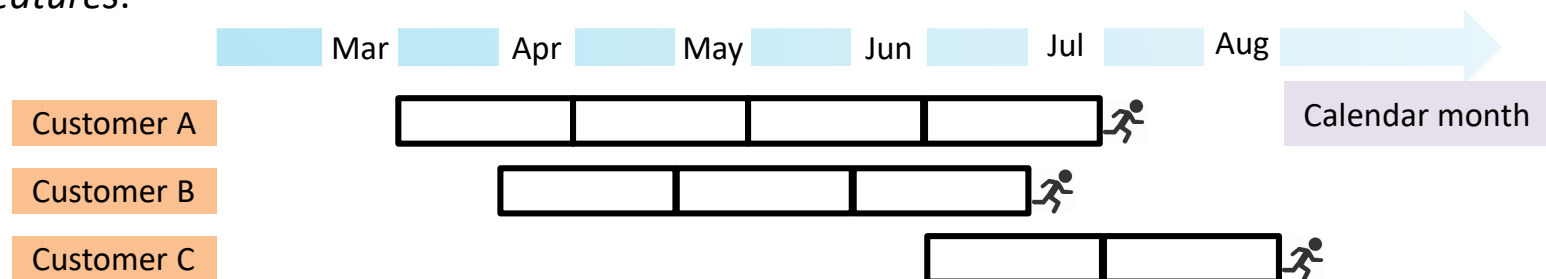
Time period: 2014-Jan-01 – 2018-Apr-20



Customer Journey



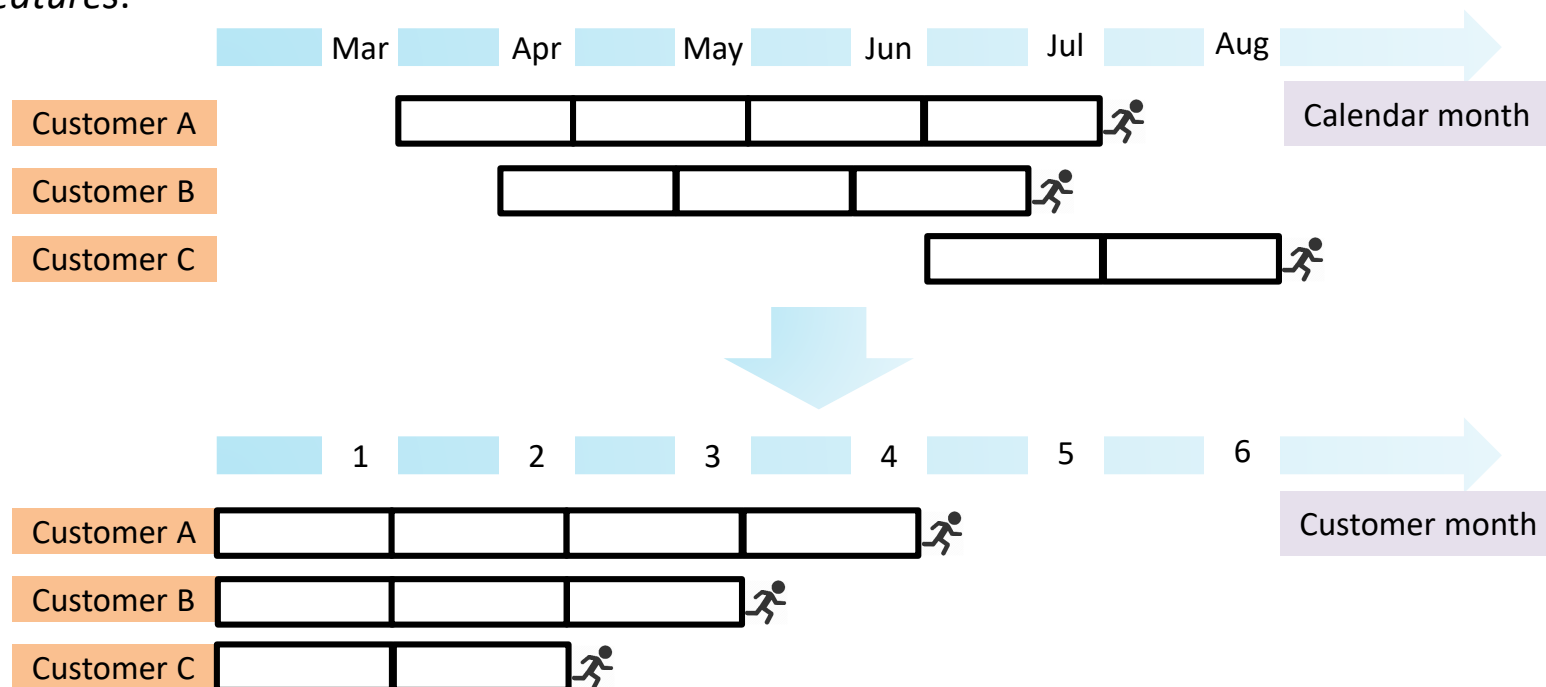
Pupils' activities are split into monthly time periods with each **customer-month** represented by *features*.



Customer Journey



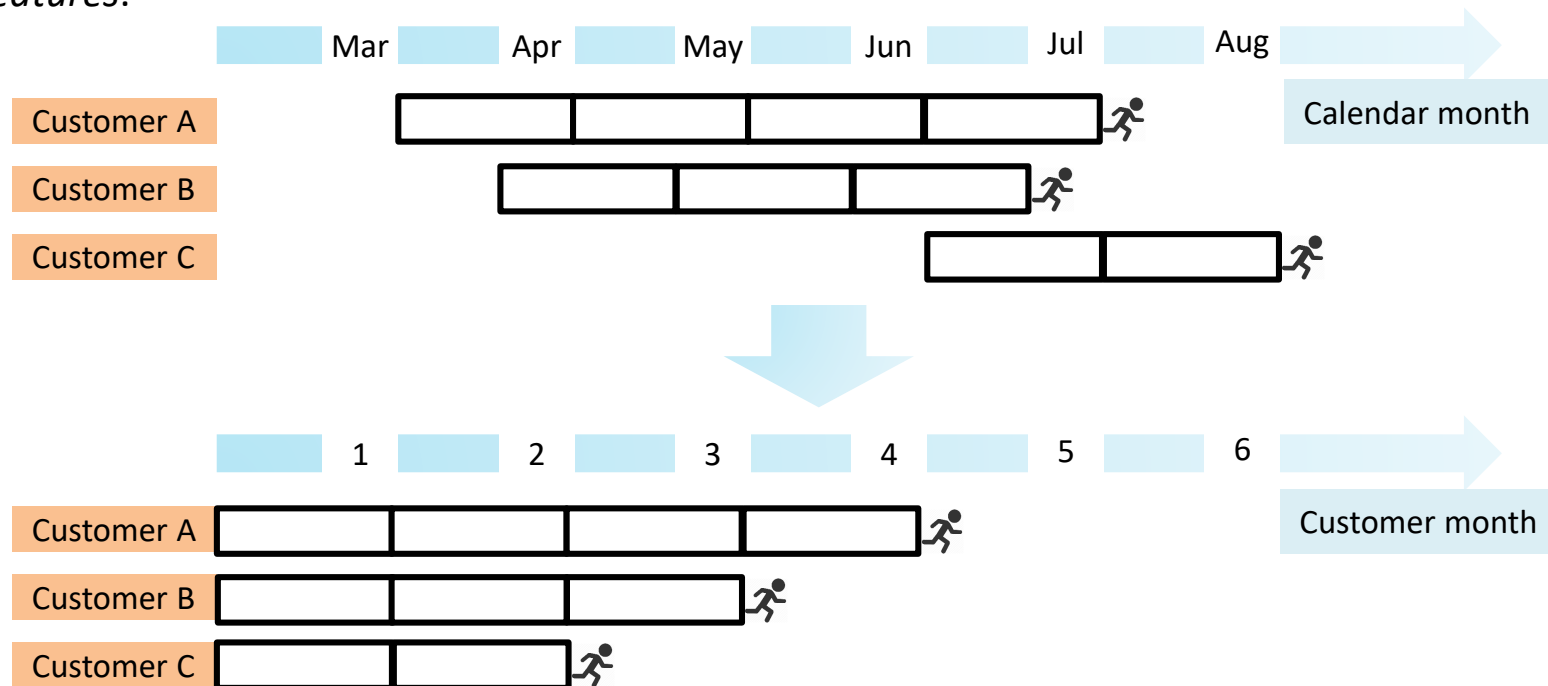
Pupils' activities are split into monthly time periods with each **customer-month** represented by *features*.



Customer Journey



Pupils' activities are split into monthly time periods with each **customer-month** represented by *features*.



We assume the cancellation to be **ONLY** dependent on the current month.

- ✓ The assumption holds since we observe statistically constant churn rate over customer month.
- ✓ The assumption enables us to treat activities in different customer months indifferently.

Methodology



Key Assumptions

- ✓ Churners and non-churners exhibit different behaviours generated by some **states**

States

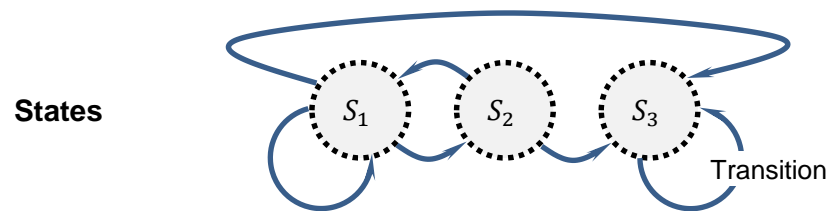


Methodology



Key Assumptions

- ✓ Churners and non-churners exhibit different behaviours generated by some **states**

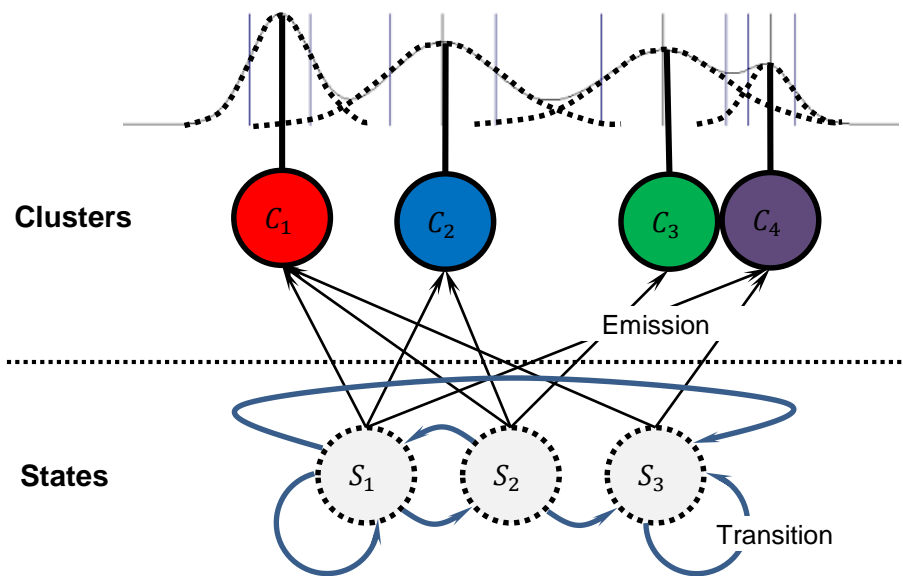


Methodology



Key Assumptions

- ✓ Churners and non-churners exhibit different behaviours generated by some **states**
- ✓ Generative process: State \rightarrow Cluster

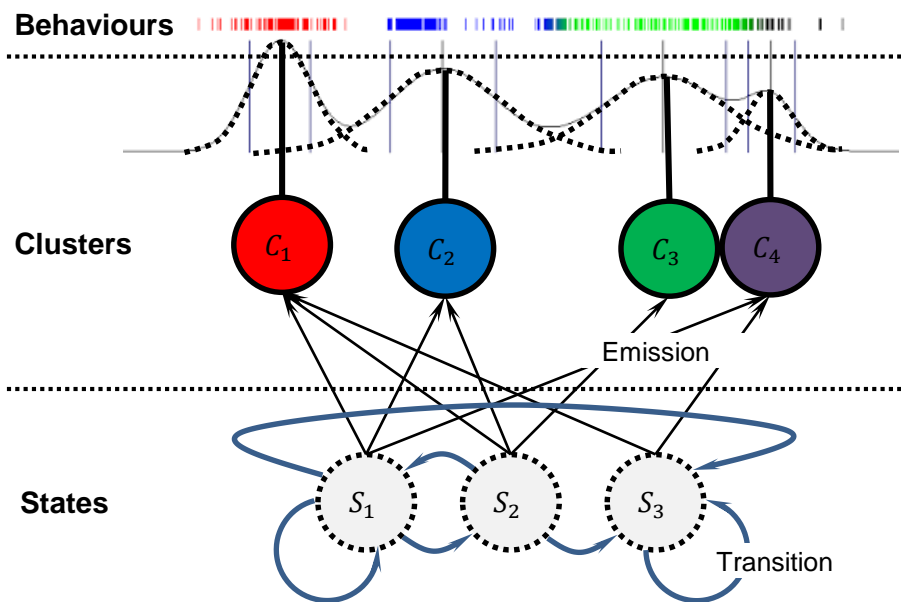


Methodology



Key Assumptions

- ✓ Churners and non-churners exhibit different behaviours generated by some **states**
- ✓ Generative process: State \rightarrow Cluster \rightarrow Behaviour observations

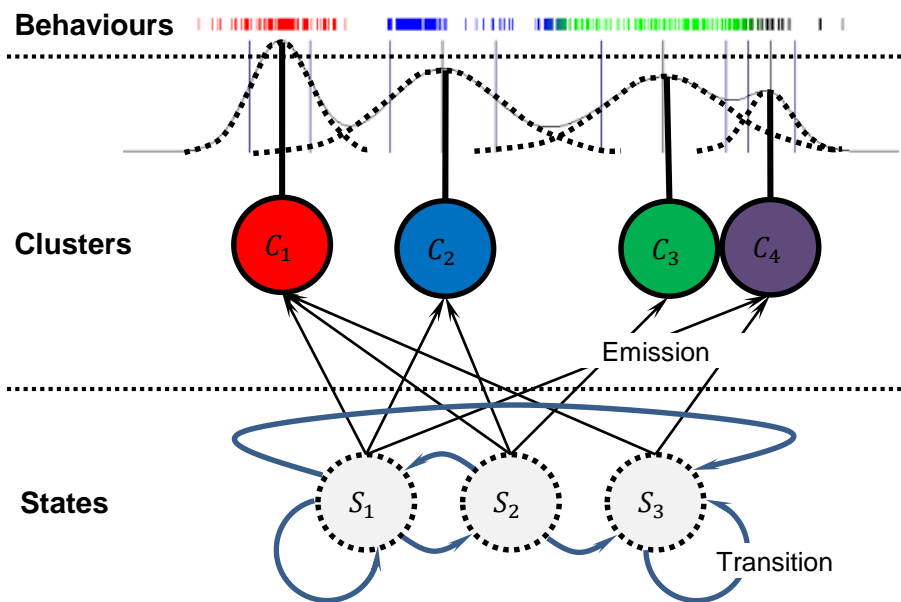


Methodology



Key Assumptions

- ✓ Churners and non-churners exhibit different behaviours generated by some **states**
- ✓ Generative process: State \rightarrow Cluster \rightarrow Behaviour observations



Mixture Model

$$f(\mathbf{X}) = \sum_{k=1}^K \pi_k f_k(\mathbf{X} | \boldsymbol{\theta}_k)$$

Markov State Transition probability

$$A = (a_{ij})$$

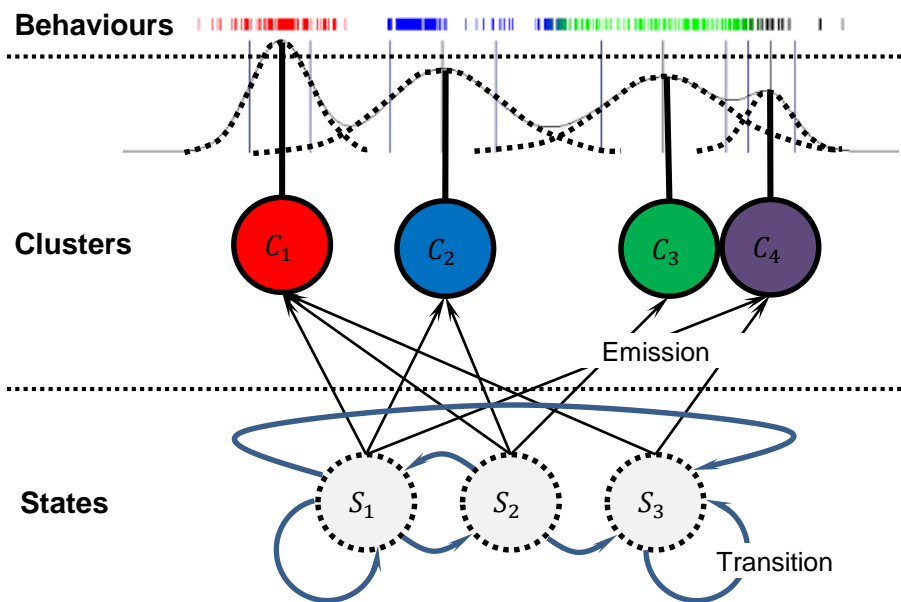
$$a_{ij} = P(s_{t+1} = S_i | s_t = S_j)$$

Methodology



Key Assumptions

- ✓ Churners and non-churners exhibit different behaviours generated by some **states**
- ✓ Generative process: State \rightarrow Cluster \rightarrow Behaviour observations
- ✓ Customers' behaviours are independent from others'



Mixture Model

$$f(\mathbf{X}) = \sum_{k=1}^K \pi_k f_k(\mathbf{X} | \boldsymbol{\theta}_k)$$

Markov State Transition probability

$$A = (a_{ij})$$

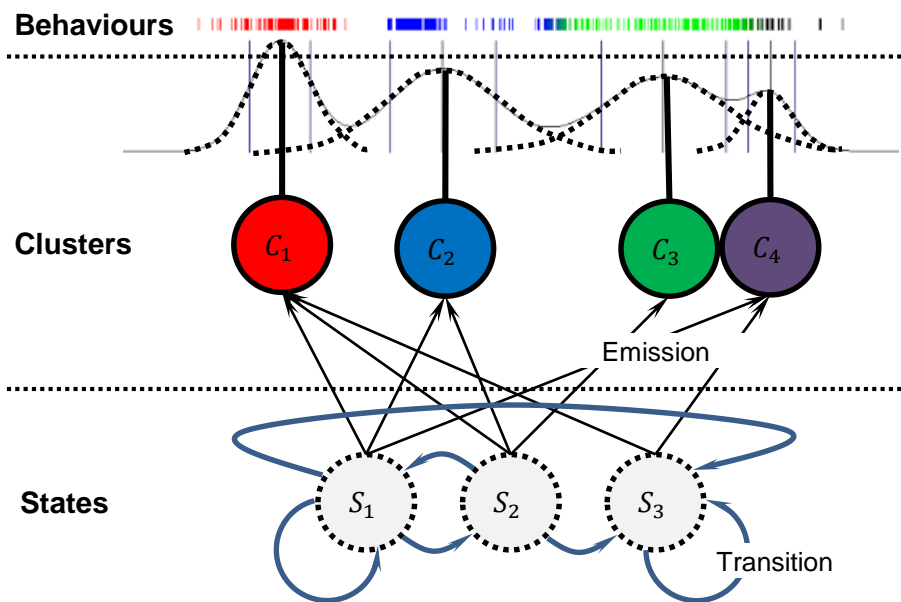
$$a_{ij} = P(s_{t+1} = S_i | s_t = S_j)$$

Methodology



Key Assumptions

- ✓ Churners and non-churners exhibit different behaviours generated by some **states**
- ✓ Generative process: State \rightarrow Cluster \rightarrow Behaviour observations
- ✓ Customers' behaviours are independent from others'



Make inference!

Mixture Model

$$f(\mathbf{X}) = \sum_{k=1}^K \pi_k f_k(\mathbf{X} | \boldsymbol{\theta}_k)$$

Markov State Transition probability

$$A = (a_{ij})$$

$$a_{ij} = P(s_{t+1} = S_i | s_t = S_j)$$

Methodology



Key Assumptions

- ✓ Churners and non-churners exhibit different behaviours generated by some **states**
- ✓ Generative process: State \rightarrow Cluster \rightarrow Behaviour observations
- ✓ Customers' behaviours are independent from others'

Pipeline

Methodology



Key Assumptions

- ✓ Churners and non-churners exhibit different behaviours generated by some **states**
- ✓ Generative process: State \rightarrow Cluster \rightarrow Behaviour observations
- ✓ Customers' behaviours are independent from others'

Pipeline



Feature Extraction

Extract features to represent customers within a specific customer month.



Feature Distributional Modelling

Assess distributions, correlations, etc.



Clustering

Fit mixture model;
Define states



Analytics and Prediction

Interpret clusters

Methodology



Key Assumptions

- ✓ Churners and non-churners exhibit different behaviours generated by some **states**
- ✓ Generative process: State \rightarrow Cluster \rightarrow Behaviour observations
- ✓ Customers' behaviours are independent from others'

Pipeline



Feature Extraction

Extract features to represent customers within a specific customer month.



Feature Distributional Modelling

Assess distributions, correlations, etc.



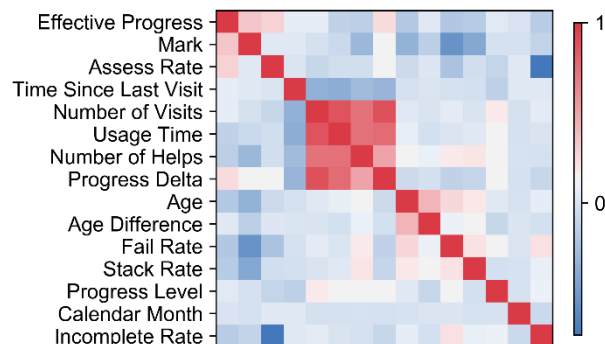
Clustering

Fit mixture model;
Define states



Analytics and Prediction

Interpret clusters



Methodology



Key Assumptions

- ✓ Churners and non-churners exhibit different behaviours generated by some **states**
- ✓ Generative process: State \rightarrow Cluster \rightarrow Behaviour observations
- ✓ Customers' behaviours are independent from others'

Pipeline



Feature Extraction

Extract features to represent customers within a specific customer month.



Feature Distributional Modelling

Assess distributions, correlations, etc.



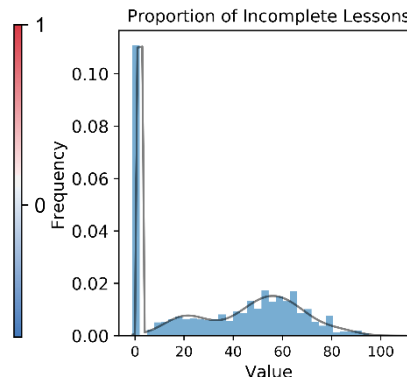
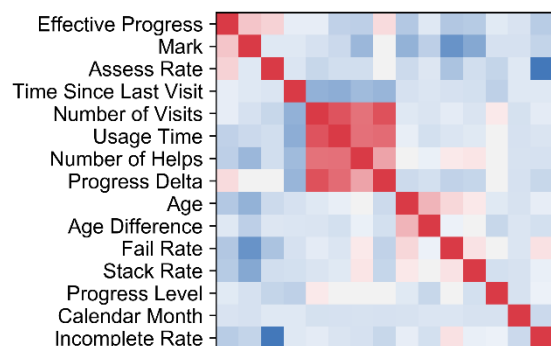
Clustering

Fit mixture model;
Define states



Analytics and Prediction

Interpret clusters



Methodology



Key Assumptions

- ✓ Churners and non-churners exhibit different behaviours generated by some **states**
- ✓ Generative process: State \rightarrow Cluster \rightarrow Behaviour observations
- ✓ Customers' behaviours are independent from others'

Pipeline



Feature Extraction

Extract features to represent customers within a specific customer month.



Feature Distributional Modelling

Assess distributions, correlations, etc.



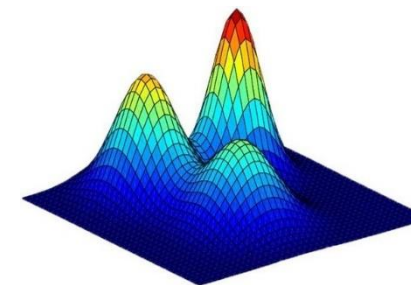
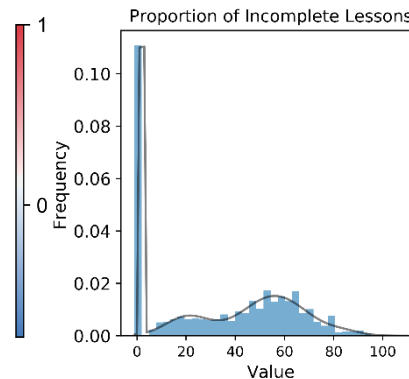
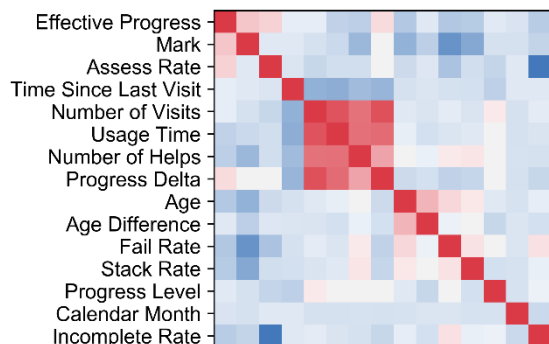
Clustering

Fit mixture model;
Define states



Analytics and Prediction

Interpret clusters



Methodology



Key Assumptions

- ✓ Churners and non-churners exhibit different behaviours generated by some **states**
- ✓ Generative process: State \rightarrow Cluster \rightarrow Behaviour observations
- ✓ Customers' behaviours are independent from others'

Pipeline



Feature Extraction

Extract features to represent customers within a specific customer month.



Feature Distributional Modelling

Assess distributions, correlations, etc.



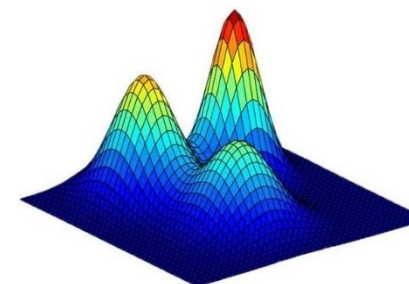
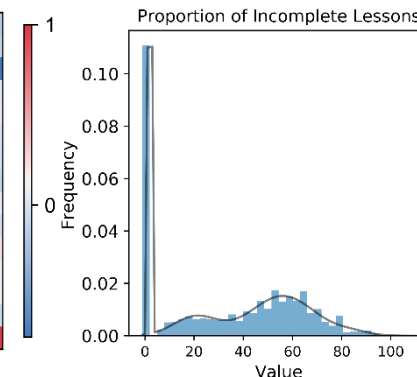
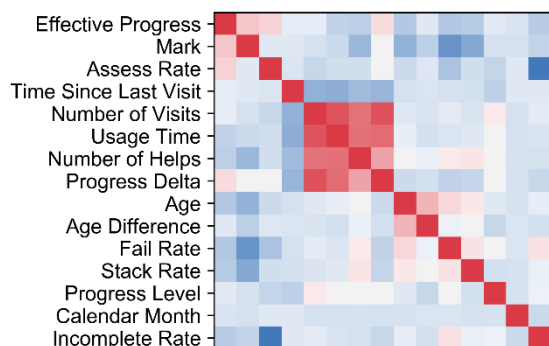
Clustering

Fit mixture model;
Define states



Analytics and Prediction

Interpret clusters



A Closer Look at Clustering...



Different users have different **data available**. For example, inactive users will have no records in features like marks, pass rates, etc.

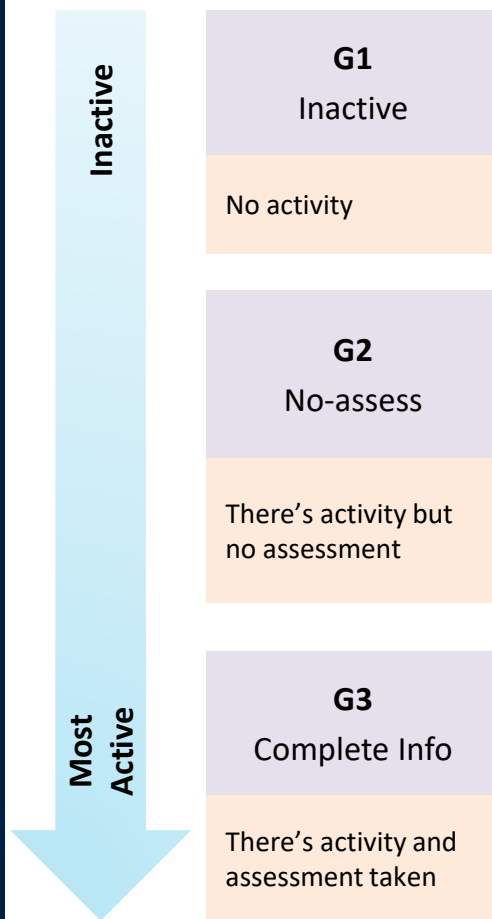
Rather than interpolating missing information, we divide customers by activity level.

A Closer Look at Clustering...



Different users have different **data available**. For example, inactive users will have no records in features like marks, pass rates, etc.

Rather than interpolating missing information, we divide customers by activity level.

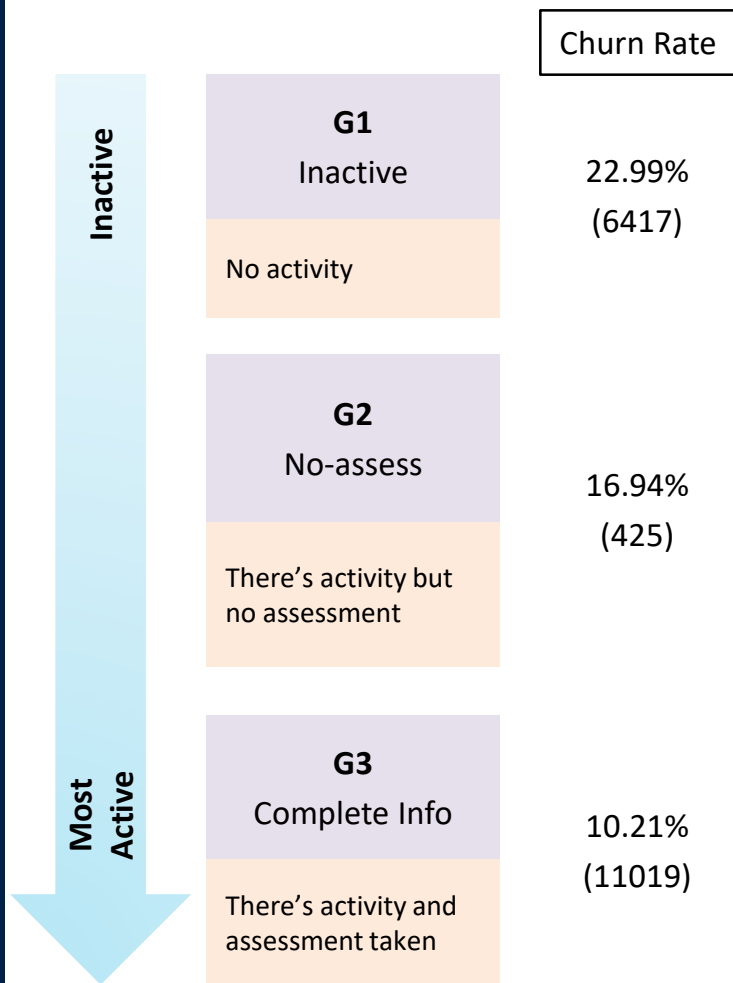


A Closer Look at Clustering...



Different users have different **data available**. For example, inactive users will have no records in features like marks, pass rates, etc.

Rather than interpolating missing information, we divide customers by activity level.

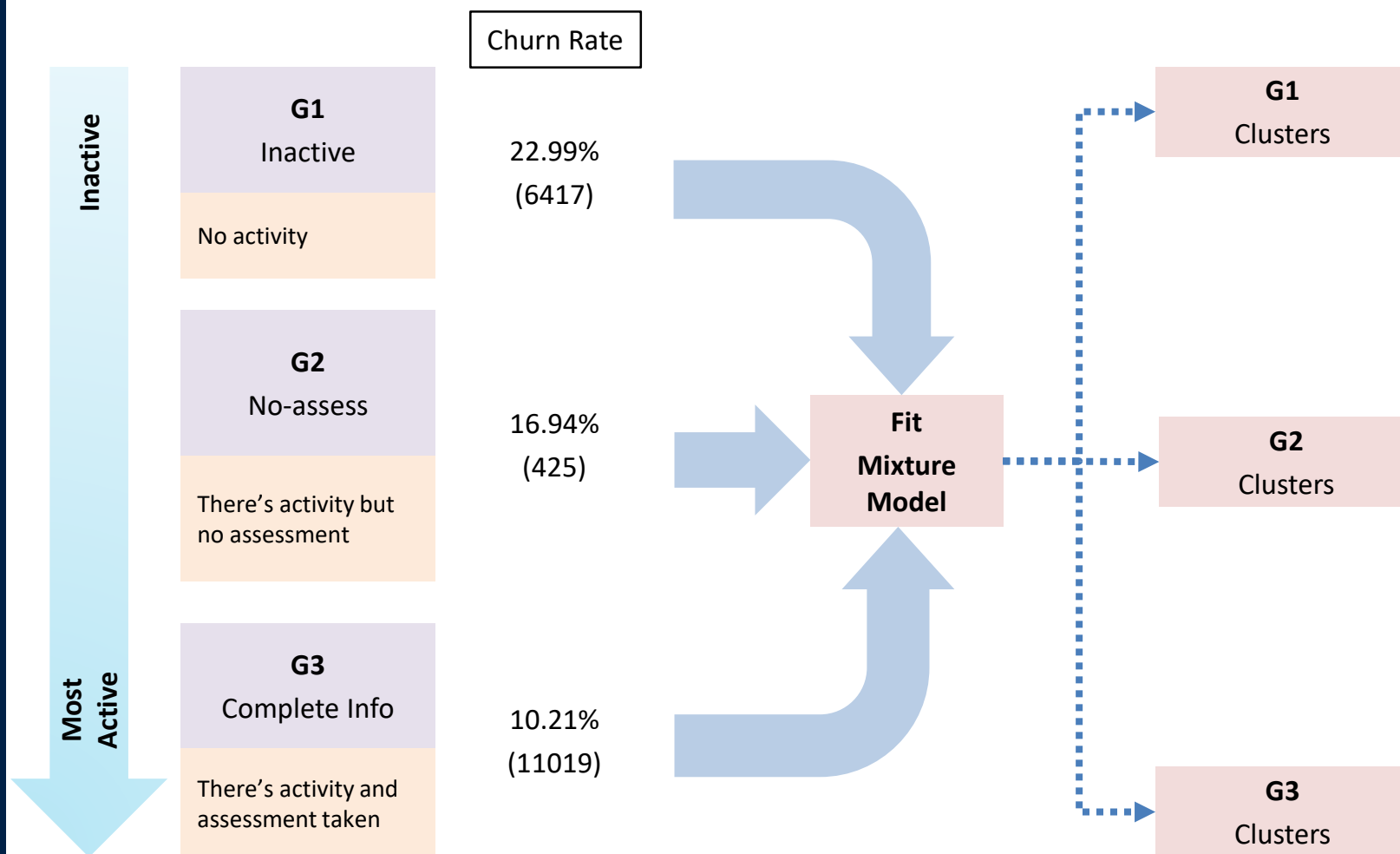


A Closer Look at Clustering...



Different users have different **data available**. For example, inactive users will have no records in features like marks, pass rates, etc.

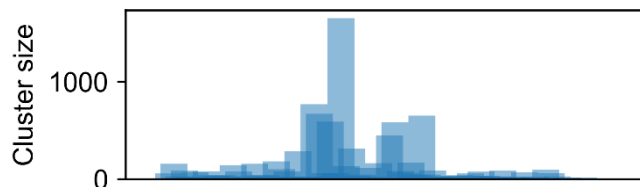
Rather than interpolating missing information, we divide customers by activity level.



Cluster Churn Rates (Who)



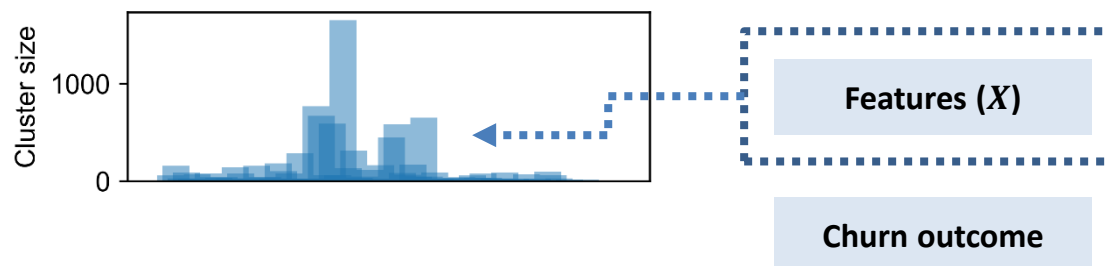
Identified Clusters: Size



Cluster Churn Rates (Who)



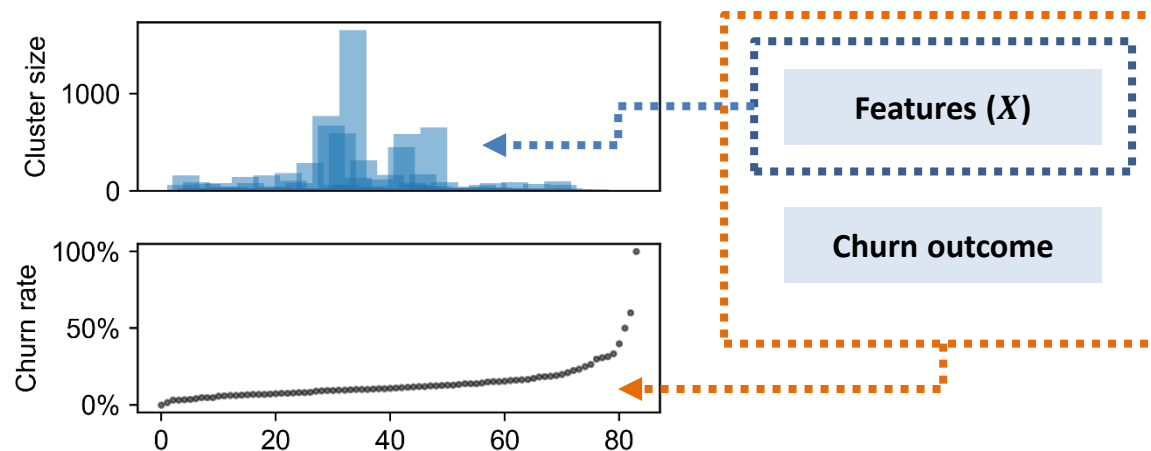
Identified Clusters: Size



Cluster Churn Rates (Who)



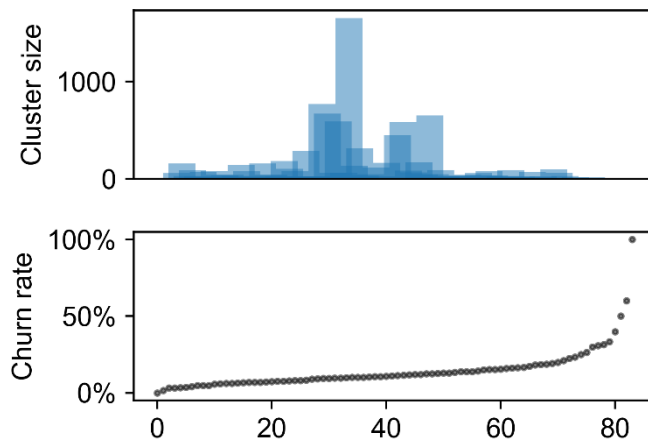
Identified Clusters: Size and Churn Rate



Cluster Churn Rates (Who)



Identified Clusters: Size and Churn Rate



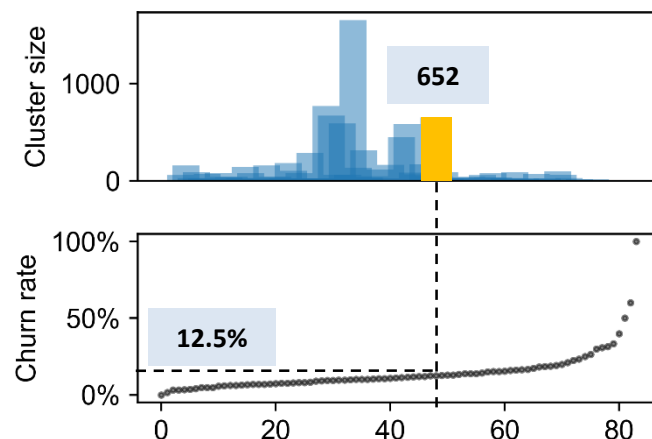
Observations:

- ✓ The purely behaviour based model infer many clusters:
 - ❑ the corresponding churn rate ranges from 0% to 100%.
- ✓ It is more challenging to identify users of extreme high/low churn probability:
 - ❑ the cluster size tends to be much smaller for those extreme high/low churn rate.

Cluster Churn Rates (Who)



Identified Clusters: Size and Churn Rate



Observations:

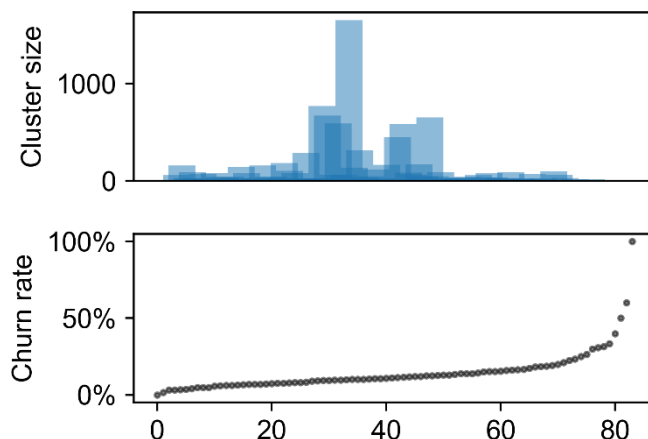
- ✓ The purely behaviour based model infer many clusters:
 - ❑ the corresponding churn rate ranges from 0% to 100%.
- ✓ It is more challenging to identify users of extreme high/low churn probability:
 - ❑ the cluster size tends to be much smaller for those extreme high/low churn rate.

From Cluster Churn Rate to Churn Probability of Individual Pupil

Cluster Churn Rates (Who)



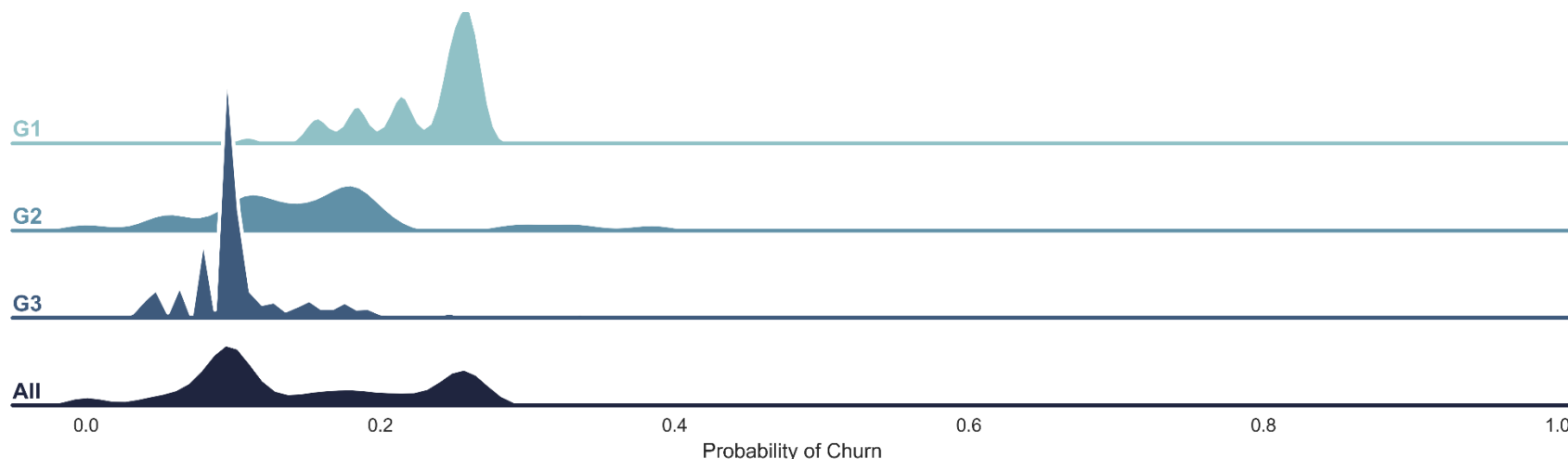
Identified Clusters: Size and Churn Rate



Observations:

- ✓ The purely behaviour based model infer many clusters:
 - ❑ the corresponding churn rate ranges from 0% to 100%.
- ✓ It is more challenging to identify users of extreme high/low churn probability:
 - ❑ the cluster size tends to be much smaller for those extreme high/low churn rate.

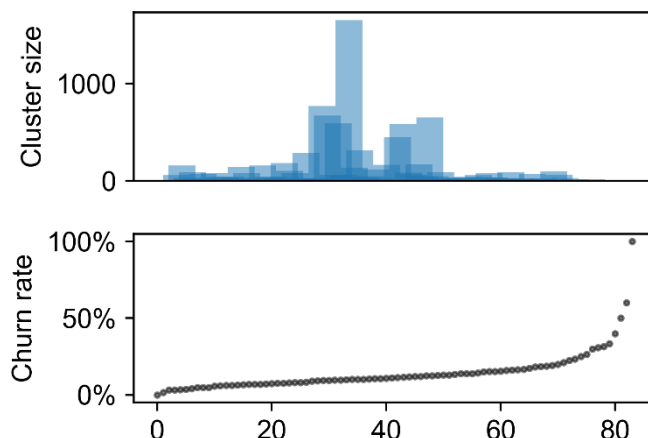
From Cluster Churn Rate to Churn Probability of Individual Pupil



Cluster Churn Rates (Who)



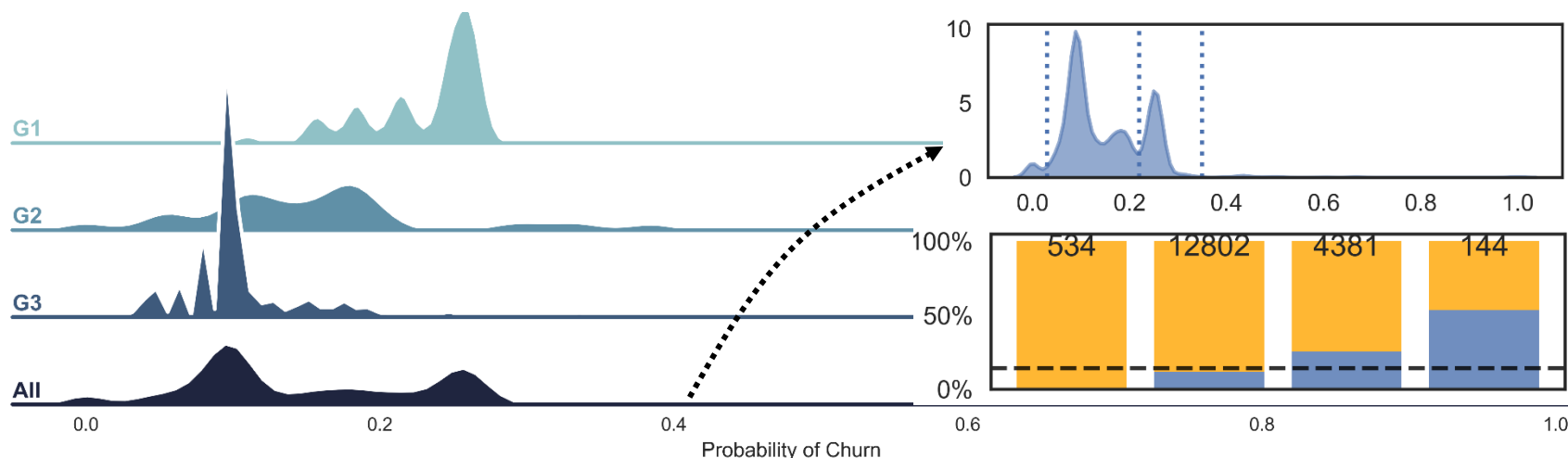
Identified Clusters: Size and Churn Rate



Observations:

- ✓ The purely behaviour based model infer many clusters:
 - ❑ the corresponding churn rate ranges from 0% to 100%.
- ✓ It is more challenging to identify users of extreme high/low churn probability:
 - ❑ the cluster size tends to be much smaller for those extreme high/low churn rate.

From Cluster Churn Rate to Churn Probability of Individual Pupil

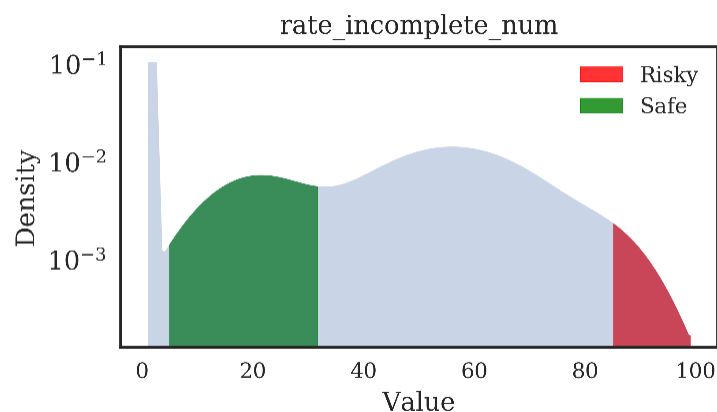


Feature Analysis (Why)



The mixture model methodology, by its nature, is designed for identifying clusters, but not explicitly for analysing feature importance.

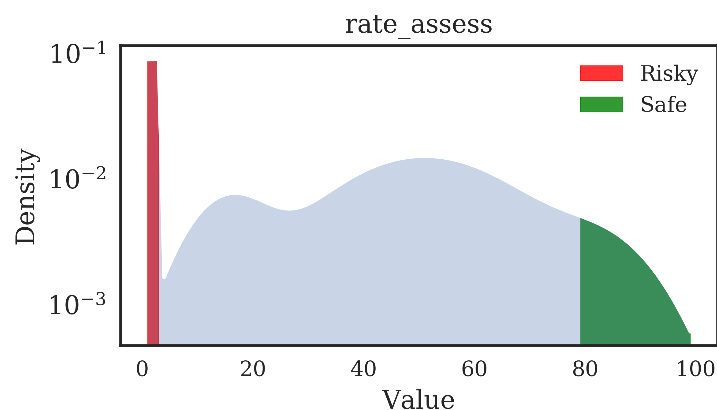
However, by counting frequencies of components of **risky** and **safe** clusters, we can see how feature impacts churn outcome.



Risky cluster: churn rate $> 50\%$

Safe cluster: churn rate $< 5\%$

We can look at the component each cluster falls in for features.



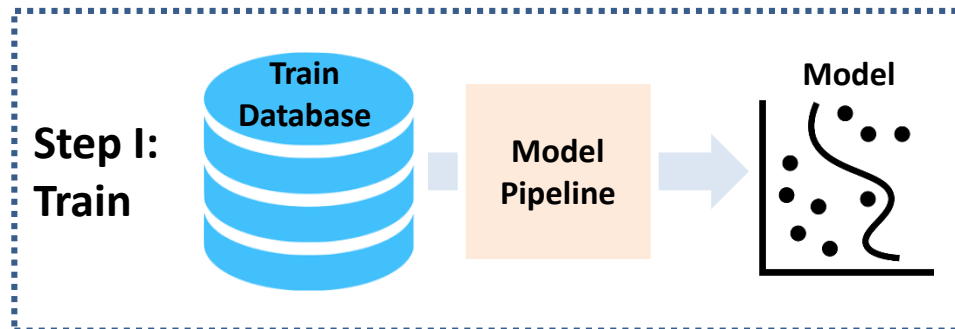
Observations:

- ✓ Subscribers having many incomplete lesson records are most likely to churn.
- ✓ Subscribers taking few assessment are most likely to churn

Prediction Workflow



How does Whizz use this model to predict the risk of churn for a set of subscribers?



Prediction Workflow



How does Whizz use this model to predict the risk of churn for a set of subscribers?

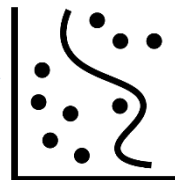
**Step I:
Train**



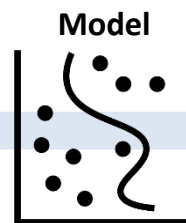
Model
Pipeline



Model



**Step II:
Predict**



2%

11%

19%

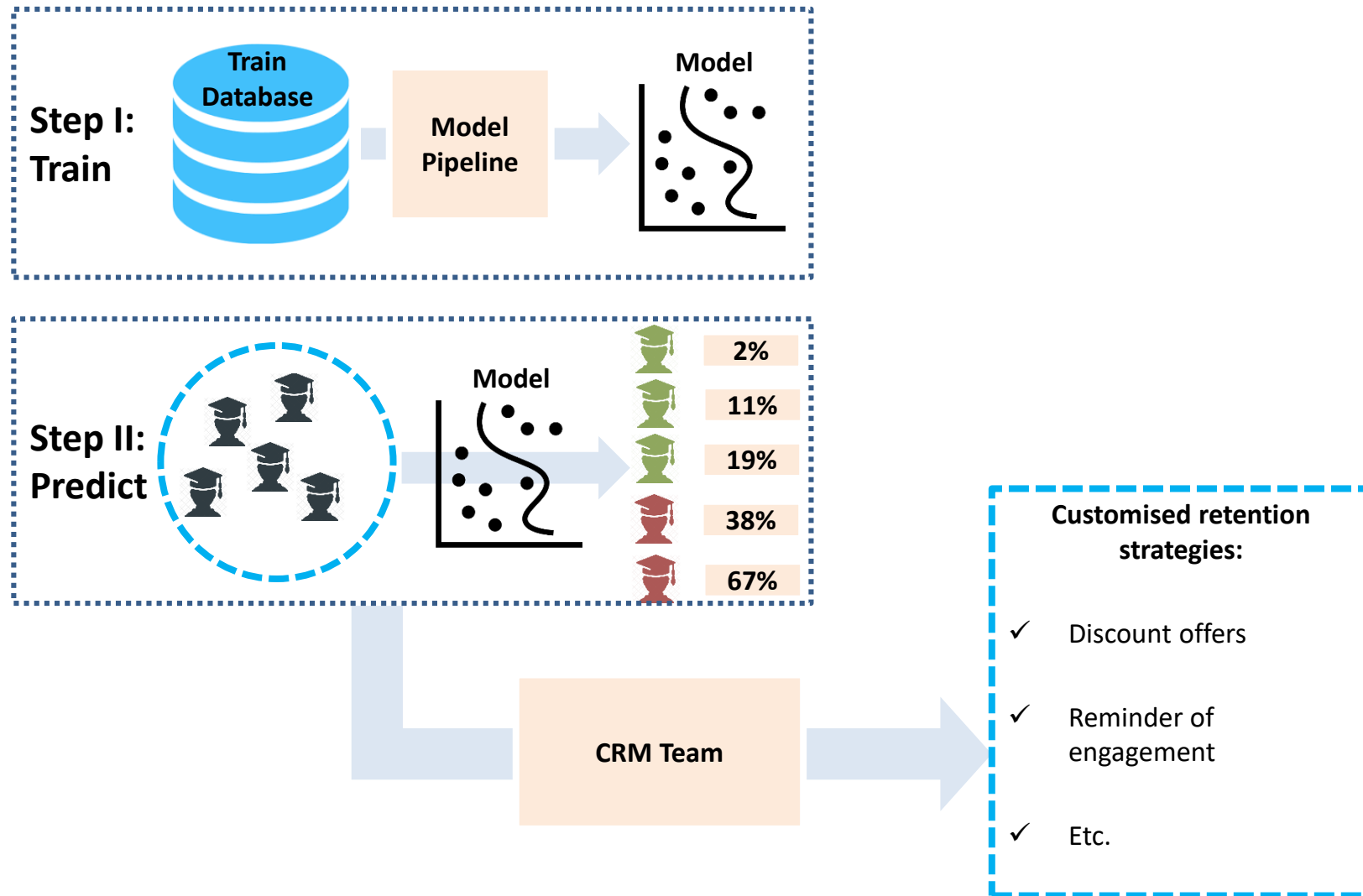
38%

67%

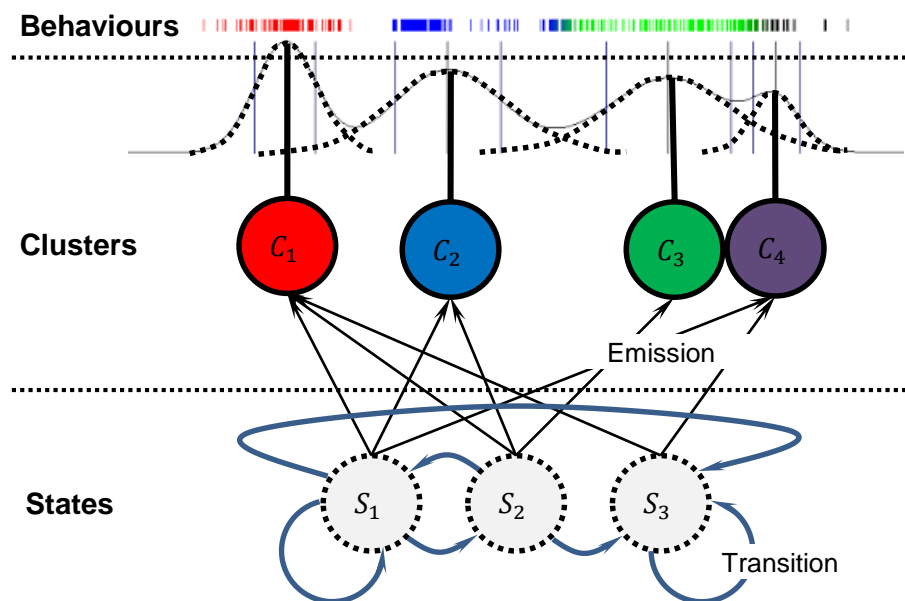
Prediction Workflow



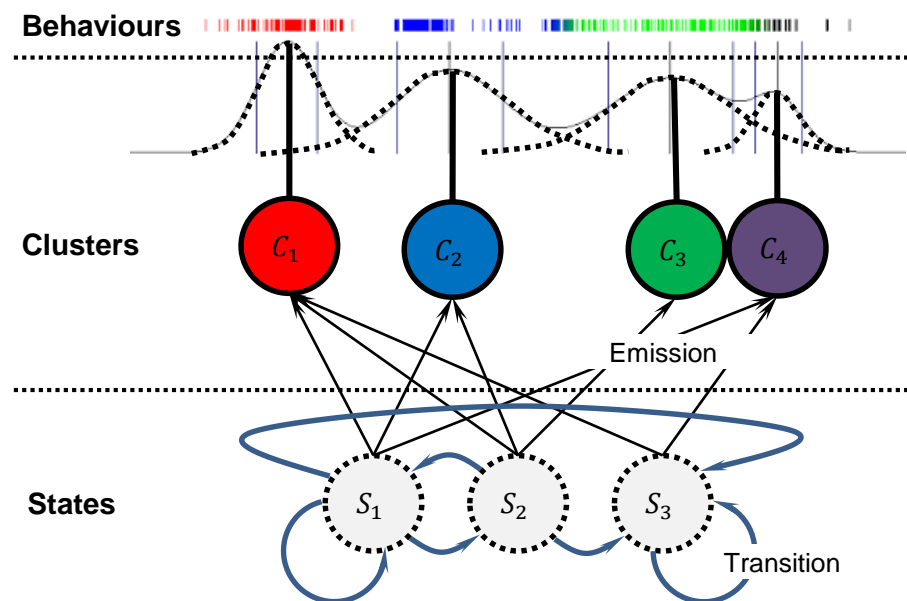
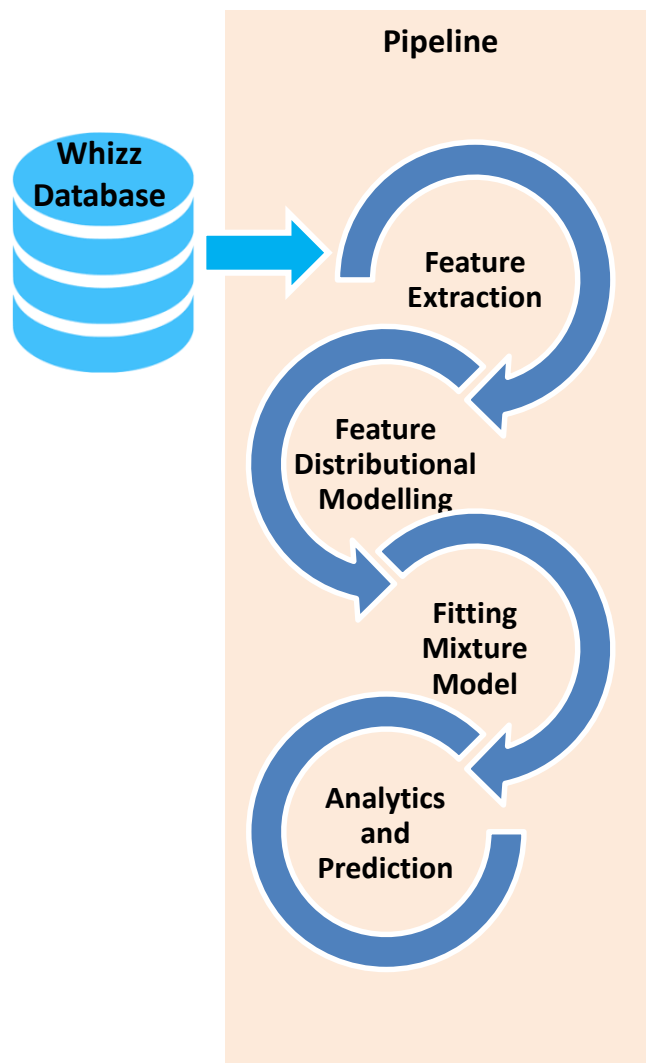
How does Whizz use this model to predict the risk of churn for a set of subscribers?



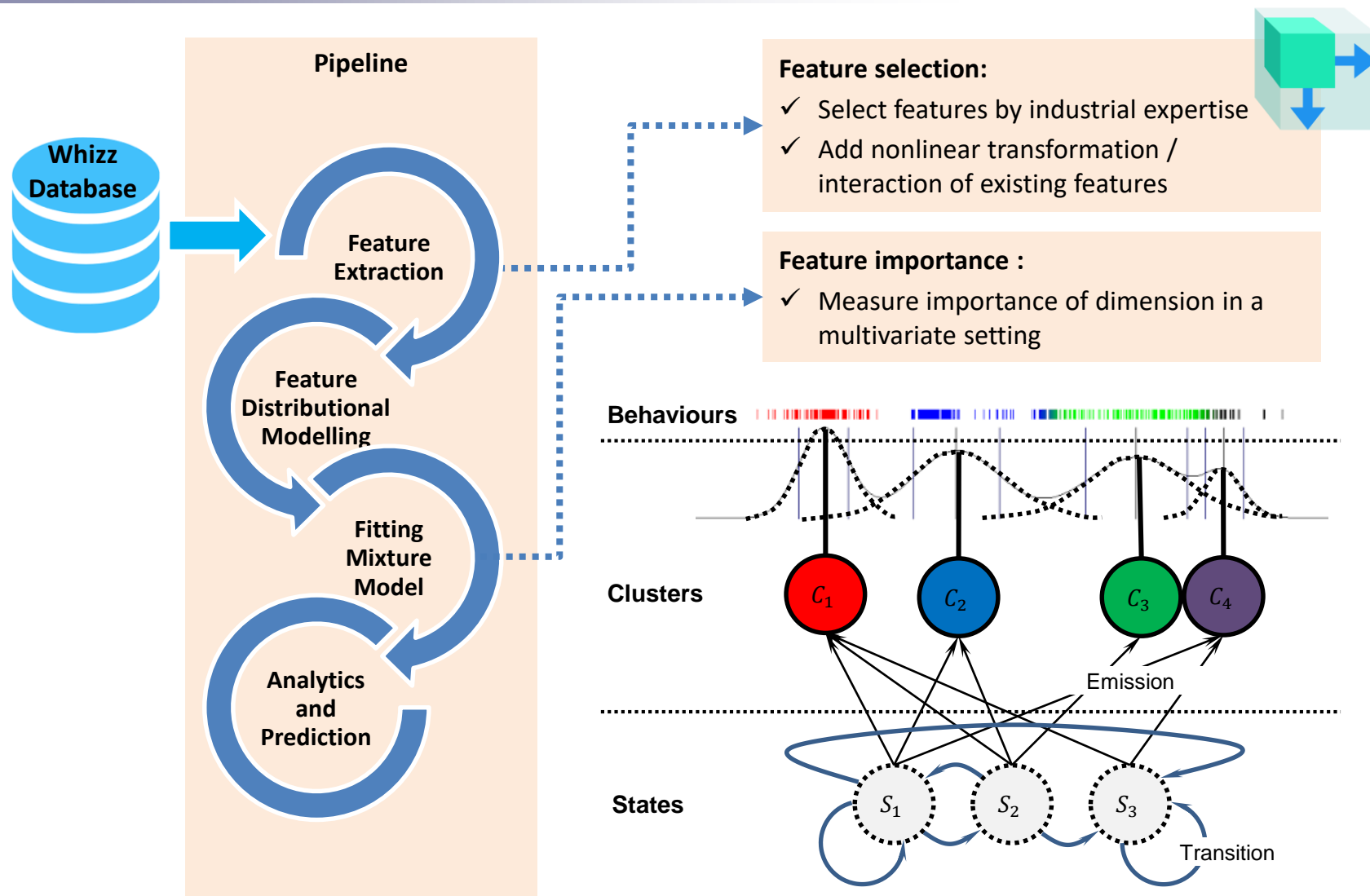
Conclusion



Conclusion



Conclusion

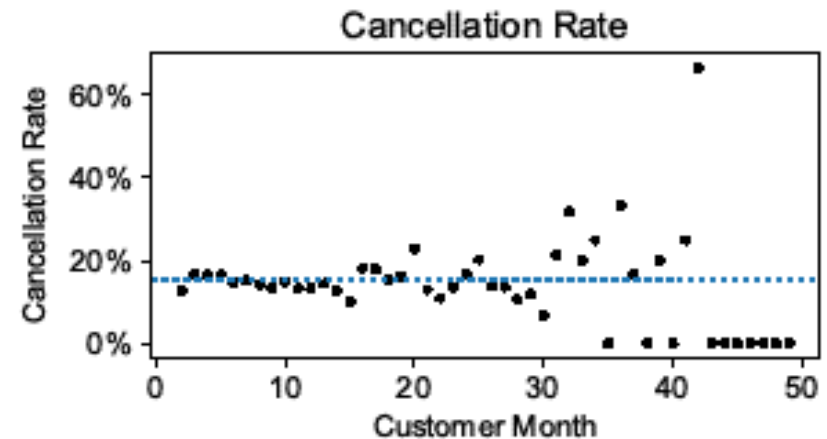
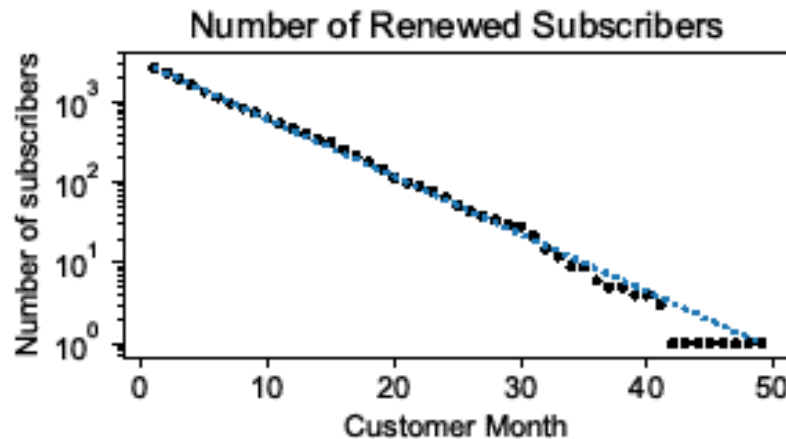


Appendix

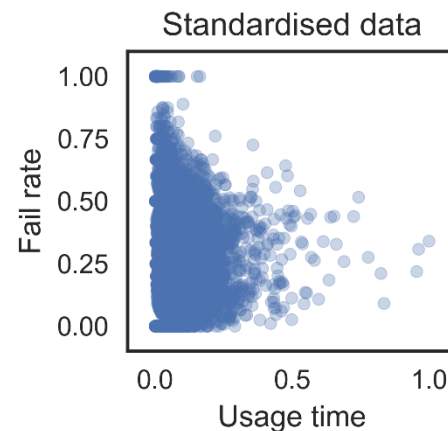
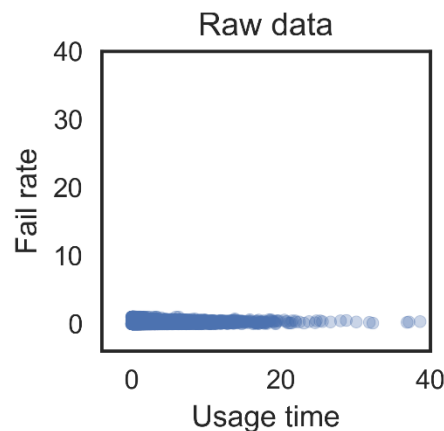
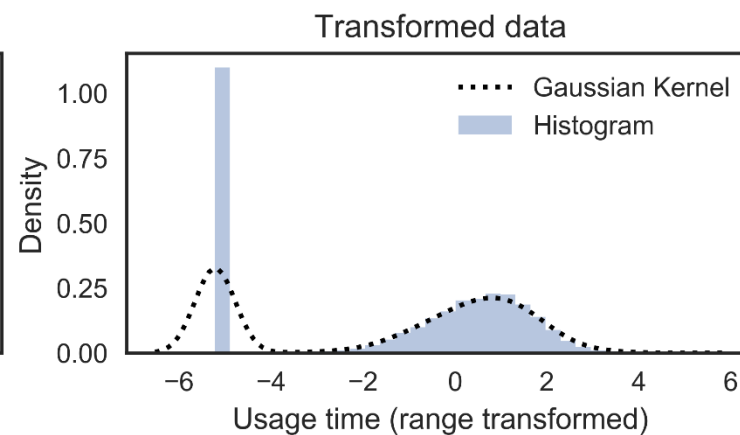
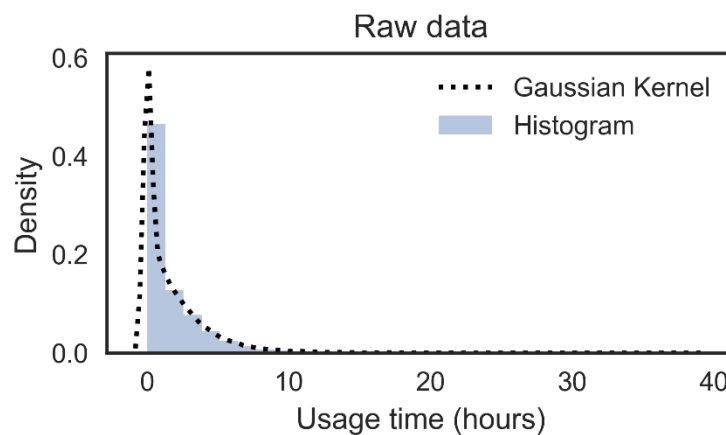
Customer Month Independence

We assume the cancellation to be **ONLY** dependent on the current month.

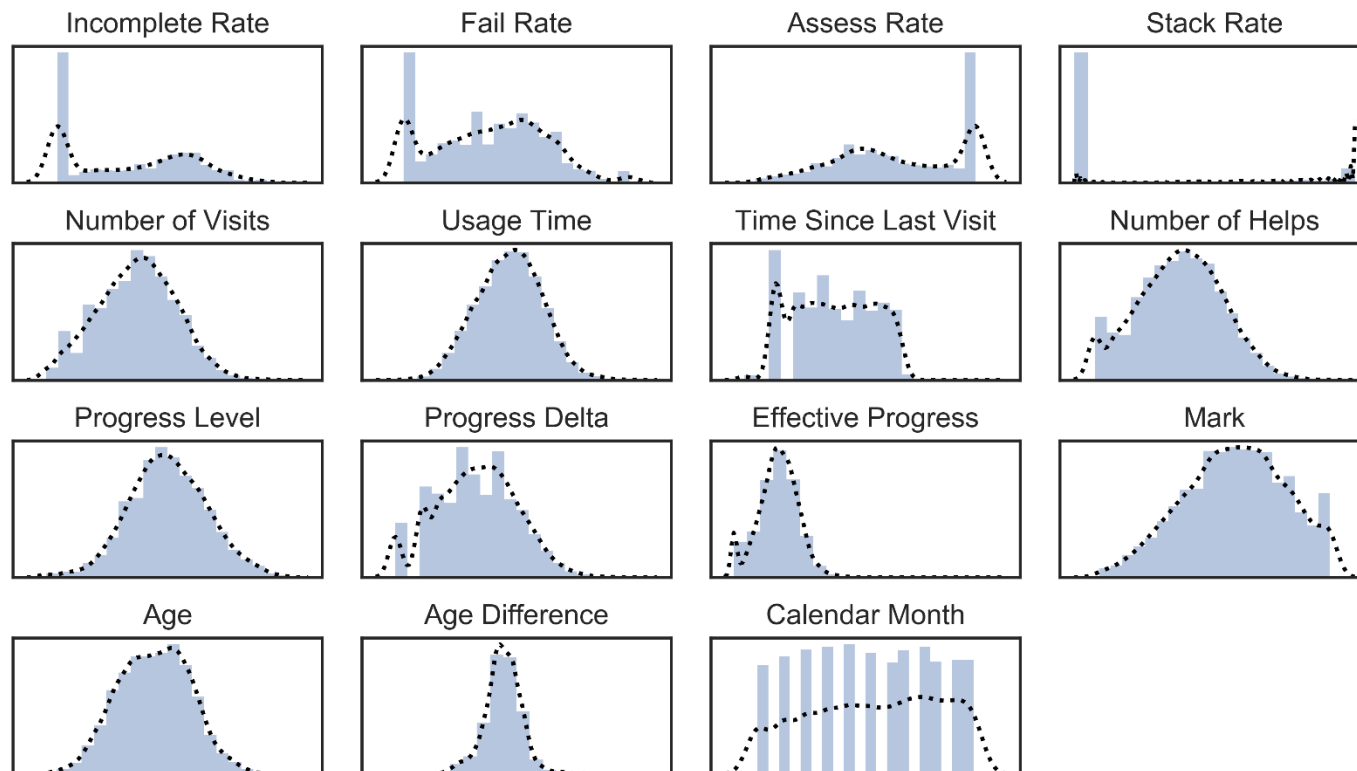
- ✓ The assumption holds since we observe statistically constant churn rate over customer month.
- ✓ The assumption enables us to treat activities in different customer months indifferently.



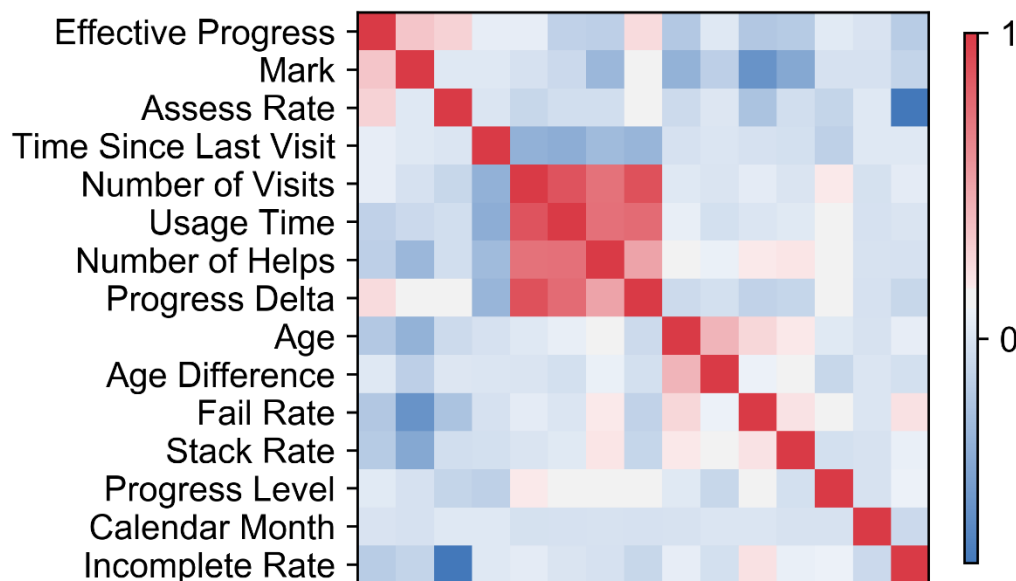
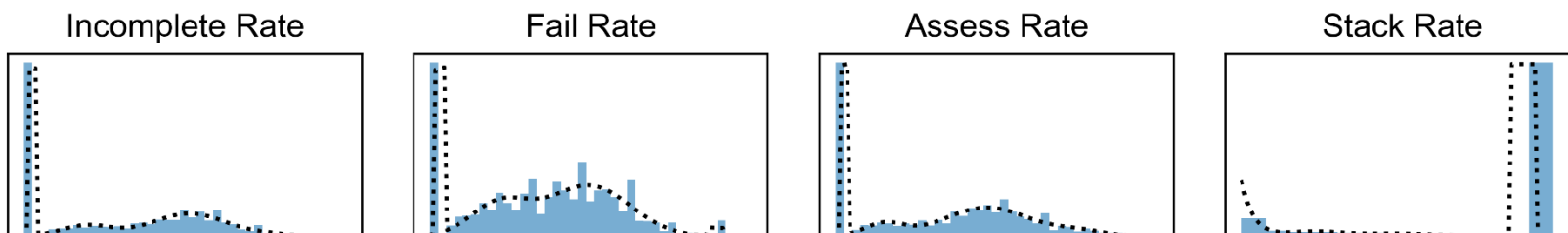
Data Transformation



Distributional Modelling – I



Distributional Modelling – II



Dirichlet Process Mixture – I



DP mixture model

Use DP as prior for θ

$$x|\theta \sim f(\cdot|\theta), \quad \theta|G \sim G(\cdot), \quad G \sim DP(\cdot|G_0, \alpha)$$

Definitions A Dirichlet Process (DP) is a distribution of a random measure. Let G_0 be a base distribution (measure) for our cluster density parameter $\theta \in \Theta$, a measurable space, and let α be a positive, real-valued scalar. A random measure G is then distributed according to *Dirichlet Process* with scaling parameter α and base measure G_0 :

$$G \sim DP(\cdot|G_0, \alpha), \quad (16a)$$

if for all $K \in \mathbb{N}$, and all $\{\Theta_1, \dots, \Theta_K\}$ finite partitions of Θ :

$$(G(\theta_1), \dots, G(\theta_K)) \sim \text{Dir}(\alpha G_0(\theta_1), \dots, \alpha G_0(\theta_K)), \quad (16b)$$

where $\text{Dir}(\cdot)$ denotes the *Dirichlet distribution*. The Dirichlet distribution is a distribution of the standard $K - 1$ simplex. Let $\boldsymbol{\pi} = \{\pi_k\}_{k=1}^K$ with $\sum_{k=1}^K \pi_k = 1$ and $\forall k : \pi_k \geq 0$, and let $\boldsymbol{\alpha} = (\alpha_1, \dots, \alpha_K)$ with $\alpha_1, \dots, \alpha_K \geq 0$. Then

$$\mathbb{P}(\boldsymbol{\pi}|\boldsymbol{\alpha}) = \text{Dir}(\alpha_1, \dots, \alpha_K) = \frac{1}{\text{Beta}(\boldsymbol{\alpha})} \prod_{k=1}^K \pi_k^{\alpha_k-1} = \frac{\Gamma(\sum_k \alpha_k)}{\prod_k \Gamma(\alpha_k)} \prod_{k=1}^K \pi_k^{\alpha_k-1}, \quad (16c)$$

where $\text{Beta}(\cdot)$ is the beta function, $\Gamma(\cdot)$ is the gamma function.

Dirichlet Process Mixture – II



Clustering Effect We use DP as a prior to distribution of cluster parameter θ :

$$\theta|G \sim G(\cdot) \quad \text{and} \quad G \sim \text{DP}(\cdot|G_0, \alpha). \quad (17)$$

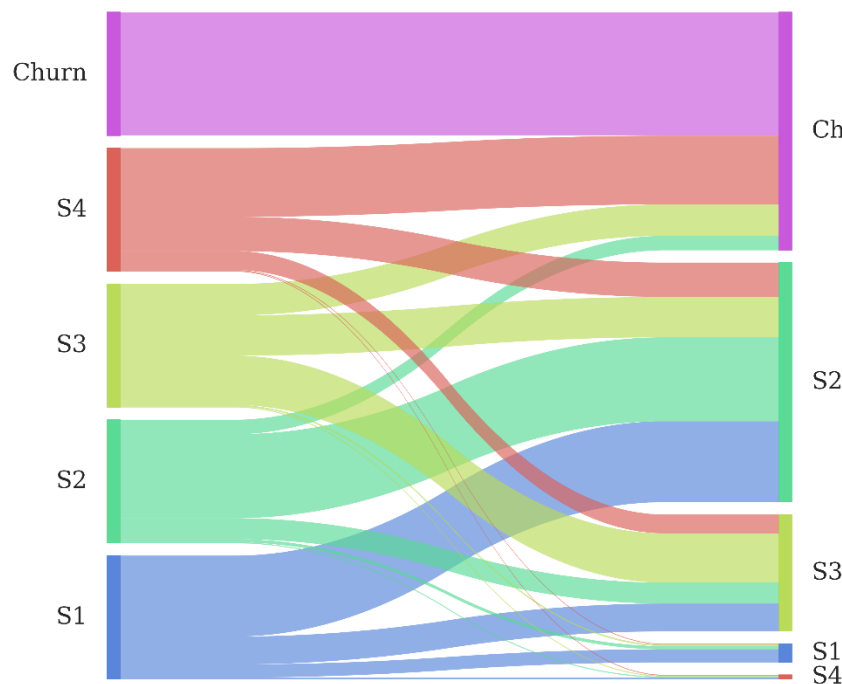
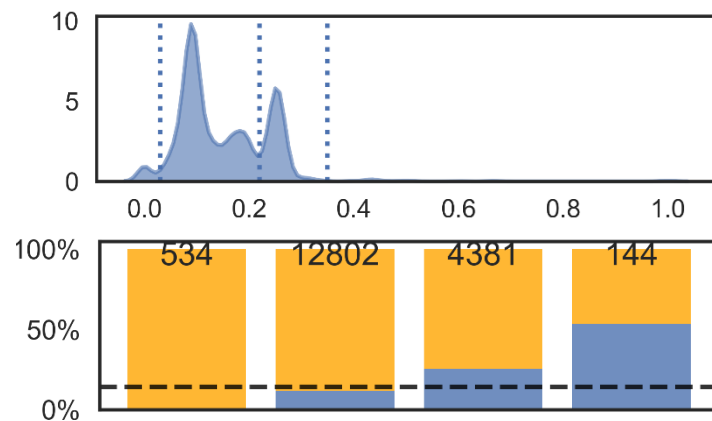
This model exhibits a “clustering effect” which enables us to infer number of clusters from data rather than pre-defining it. Suppose we independently draw n random values $\theta^{(j)}$ from G under the model (17), then Blackwell and MacQueen’s urn representation theorem [1] states that, marginalising out the random measure G , the joint distribution of the collection of variables $\{\theta^{(1)}, \dots, \theta^{(n)}\}$ exhibits a clustering effect:

$$\mathbb{P}(\theta^{(j)}|\theta^{(1)}, \dots, \theta^{(j-1)}) \propto \alpha G_0(\theta^{(j)}) + \sum_{l=1}^{j-1} \delta_{\theta^{(l)}}(\theta^{(j)}), \quad (18)$$

where $\delta_{\theta^{(l)}}(\cdot)$ is a Dirac delta at $\theta^{(l)}$. Thus the variables $\{\theta^{(1)}, \dots, \theta^{(n)}\}$ are randomly partitioned according to which variables are equal to the same value. Moreover, let $\{\theta_1, \dots, \theta_K\}$ denote the distinct values of the drawn samples $\{\theta^{(1)}, \dots, \theta^{(j-1)}\}$, let $\{\kappa_1, \dots, \kappa_{j-1}\}$ be the assignment variables such that $\theta^{(l)} = \theta_{\kappa_l}$. Then,

$$\mathbb{P}(\theta^{(j)}|\theta^{(1)}, \dots, \theta^{(j-1)}) \propto \frac{\alpha}{j-1+\alpha} G_0(\theta^{(j)}) + \sum_{k=1}^K \frac{|\{l : \kappa_l = k\}|}{j-1+\alpha} \delta_{\theta^{(l)}}(\theta^{(j)}). \quad (19)$$

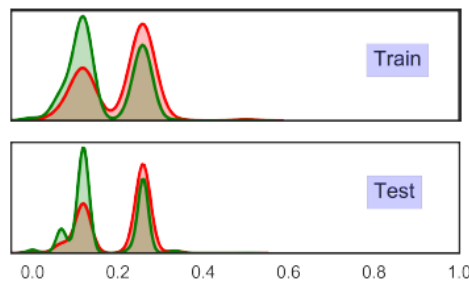
Temporal Transition of States



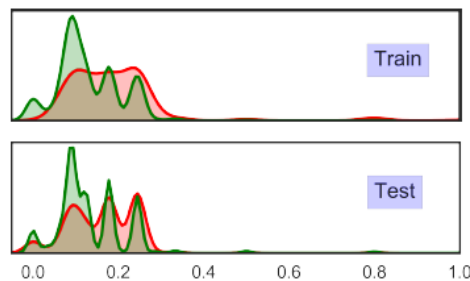
✓ S2 is the largest destination state, which makes sense as S2 represents a “normal” state in which the churn probability is approximately the same as population average.

✓ Barely pupils transit from other states to S4 (the riskiest state). This seems imply that the reason for strong intention of churn might be purely external and irrelevant to customer experience at Whizz, since there is little chance of transiting to S4.

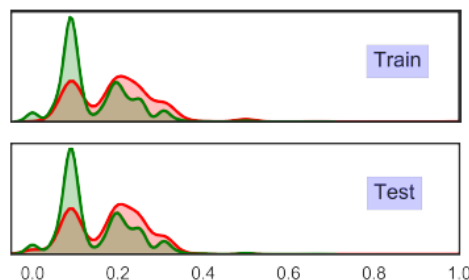
Overfitting



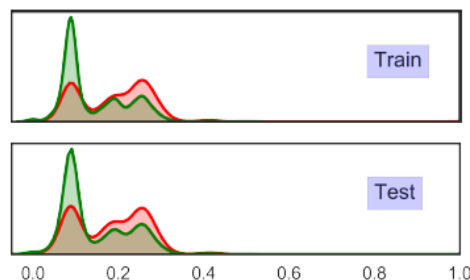
(a) Train Pct = 5%



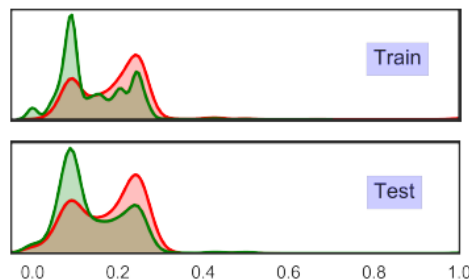
(b) Train Pct = 10%



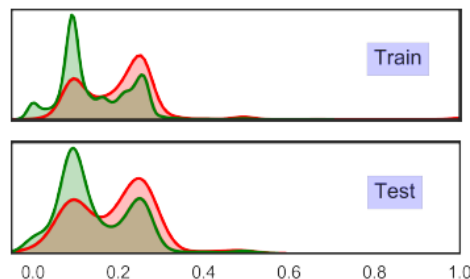
(c) Train Pct = 50%



(d) Train Pct = 70%



(e) Train Pct = 90%



(f) Train Pct = 95%