

# DSAA 5002 Final Project

2025.10.27

Feiyu Huang

Han Linghu

# General Requirements

1. **Topic:** Related to the course material, e.g., data preprocessing, classification, clustering, anomaly detection.
2. **Format:** KDD, double-column, **maximum 6 pages** including reference and appendix.  
Refer to <https://kdd2025.kdd.org/research-track-call-for-papers/>
3. **Sections:**  
**At least contain:** introduction, related work, method, experiments.  
**Encouraged:** theoretical derivation.
4. **Writing:** Use concise language and avoid complex sentences.
5. **Novelty:** Show how you differentiate from related work.
6. **Code:** Required.

**Grading: 50% (Presentation 25% + Report 25%)**

**Report Due: 2025.12.19 (End of Day)**

**Presentation Dates: Dec 8<sup>th</sup>, 15<sup>th</sup>**

Students with **even** ID numbers will present on Dec 8th,  
and those with **odd** ID numbers will present on Dec 15th.

**The report and presentation file need to be submitted to Canvas**

# General Requirements

---

- **Individual Project.**
- No previously published papers are allowed. Submissions based on published work will receive **0 points**.
- Plagiarism in any form is strictly prohibited. Violations will result in **0 points**.
- At the beginning of **both** presentation and report, clearly state:
  1. your **project title**; and
  2. the **specific course topic** your project relates to**Failure to do so will incur a 20% penalty on the total project score.**

# About Report Writing (25%)

---

**Clear contribution (differentiation from previous existing works)**

**Well-structured format**

**Well-organized narration logic**

**Make it complete.**

# About Presentation (25%)

**General:** 5 min presentation for each person.

**Format:** Due to the large enrollment this semester, project presentations will be conducted in two formats:

**In-person:** A random half of each class section will present in class during scheduled class time. The presentation schedule and order will be announced in advance.

**Online:** Within three days of the presentation date, remaining students must upload a 5-minute recorded presentation video to Baidu Cloud (or another file-storage service) and submit only the video link on Canvas (include the extraction/access code if required).

Students with **even** ID numbers must upload a video by Dec 11<sup>th</sup>.

Students with **odd** ID numbers must upload a video by Dec 18<sup>th</sup>.

**All presentations, live or recorded, must adhere to the same time limit and content expectations.**

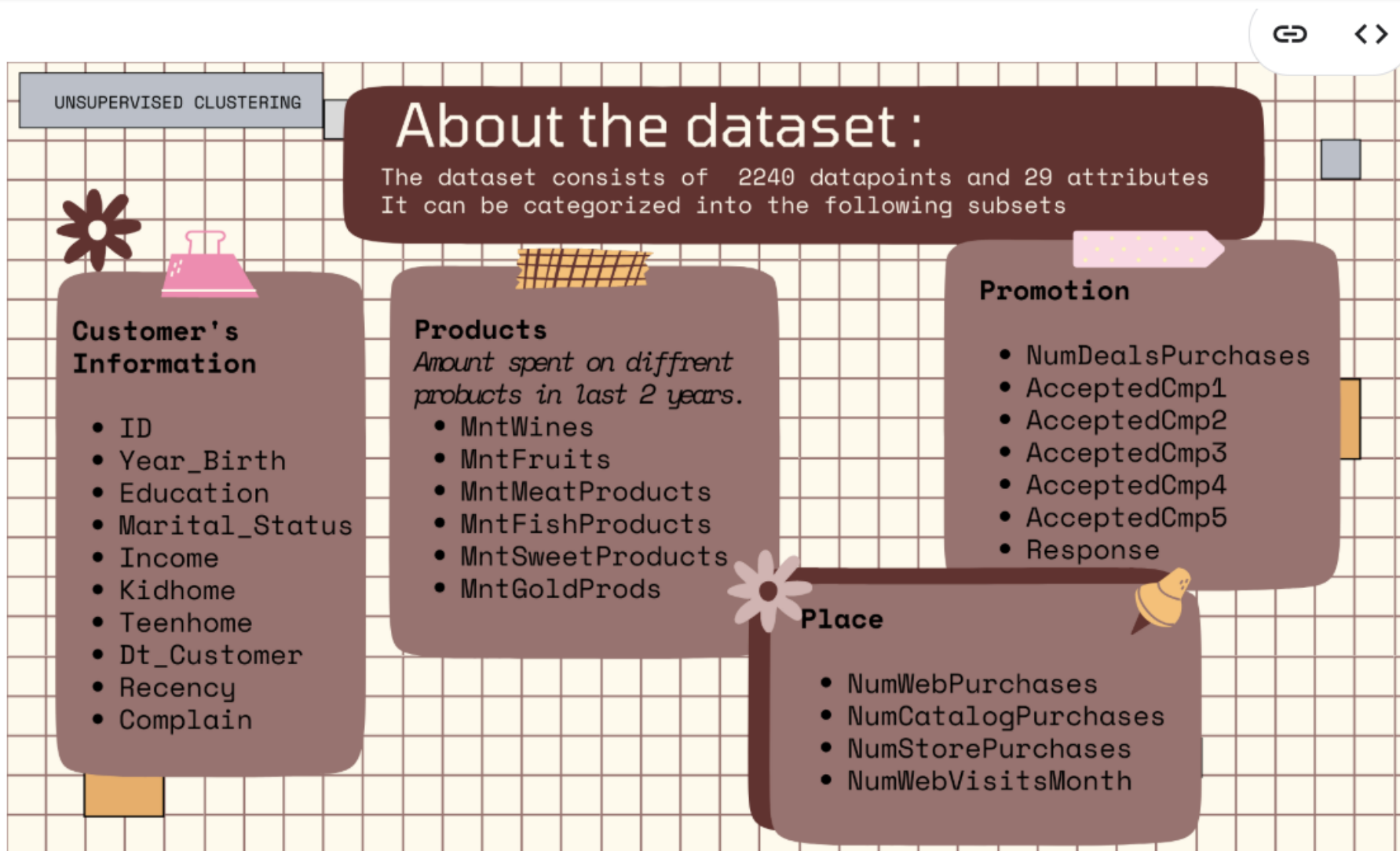
# Default Topic - Customer Segmentation



Topic website: <https://www.kaggle.com/code/karnikakapoor/customer-segmentation-clustering/notebook>

If you choose this default topic, you don't need to specify which topic it relates to.

- This project offers a basic approach to customer segmentation, aiming to understand customer groups better and support targeted marketing efforts.
- It serves as a **60-point baseline**, encouraging students to apply **innovative ideas for improved customer analysis**, such as exploring new clustering methods or models.
- All creative extensions should still be related to the core content of this course.



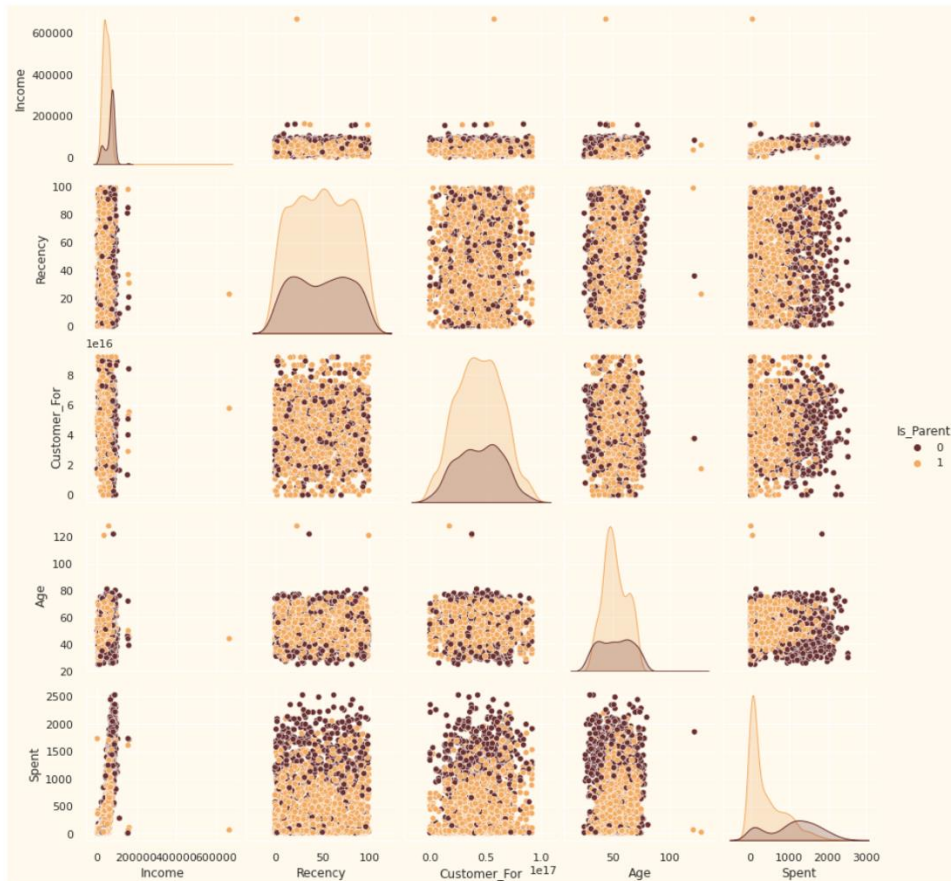


# Default Topic

## 1. Data Preparation & Feature Selection

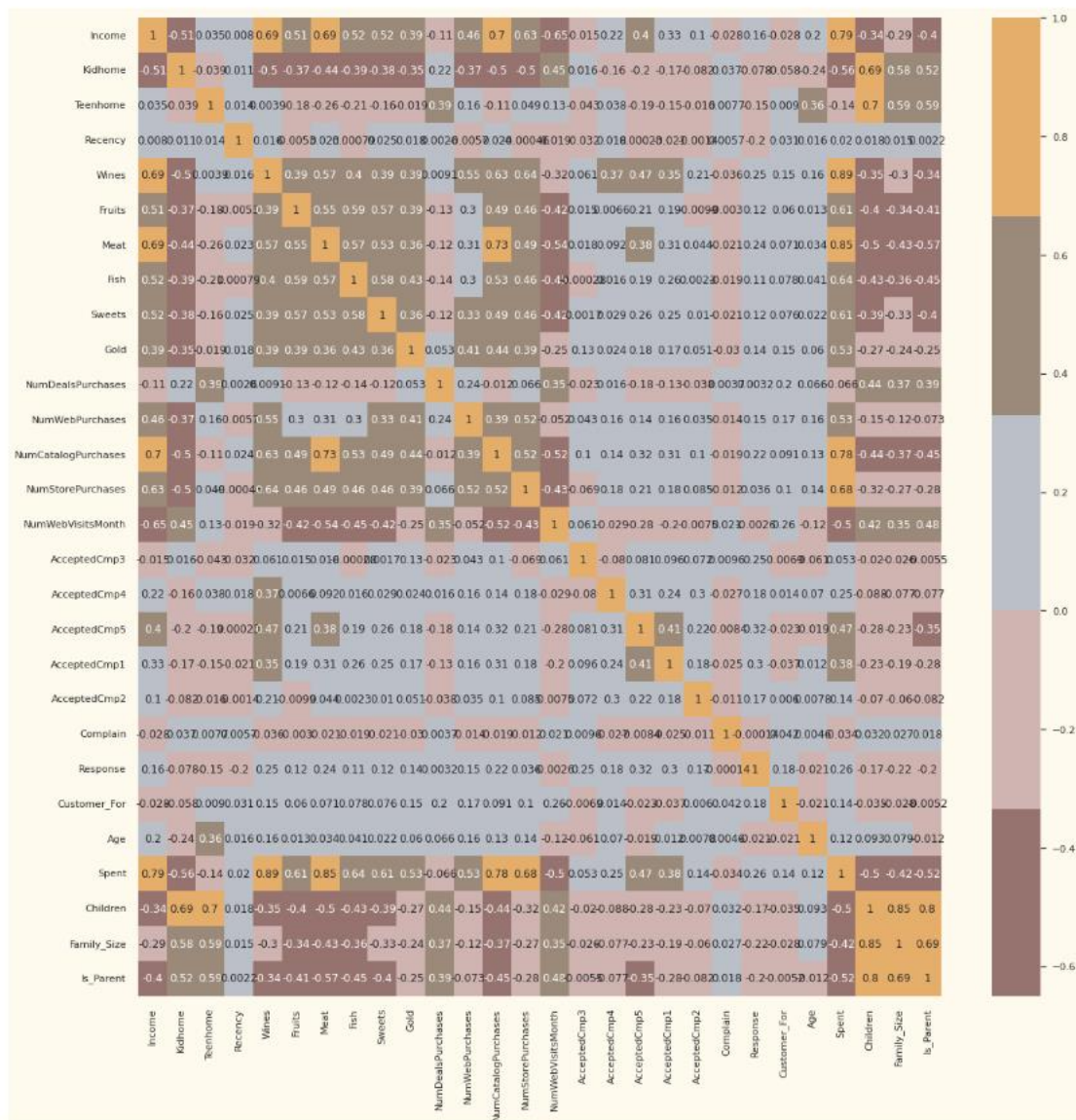
**Data Cleaning:** Handle missing and abnormal values.

**Feature Selection:** Use key indicators such as Recency (last purchase time), Frequency, and Monetary value as the basis for clustering.



Clearly, there are a few outliers in the Income and Age features, and deleting the outliers in the data.

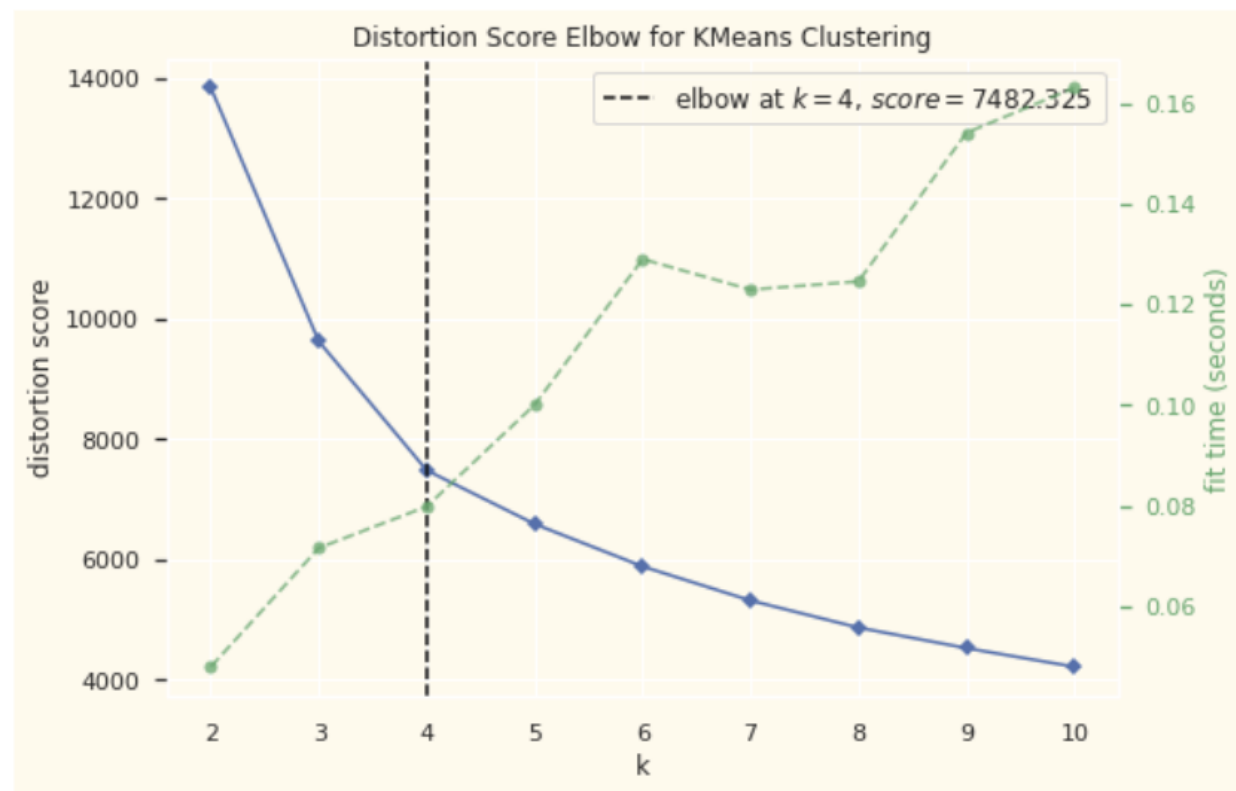
# Default Topic



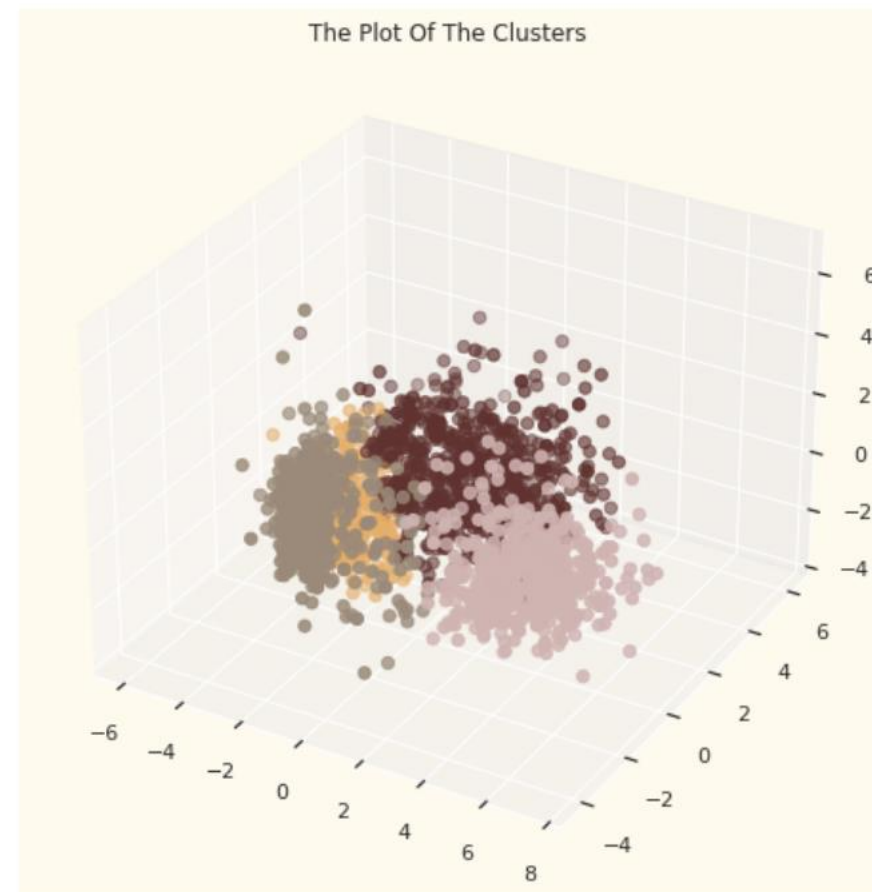
The data is quite clean and the new features have been included. We will proceed to the next step.

## 2. Clustering Method

**K-Means Clustering:** The most common partition-based clustering method used in this



Elbow Method to determine the number of clusters to be formed



To examine the clusters formed let's have a look at the 3-D distribution of the clusters.

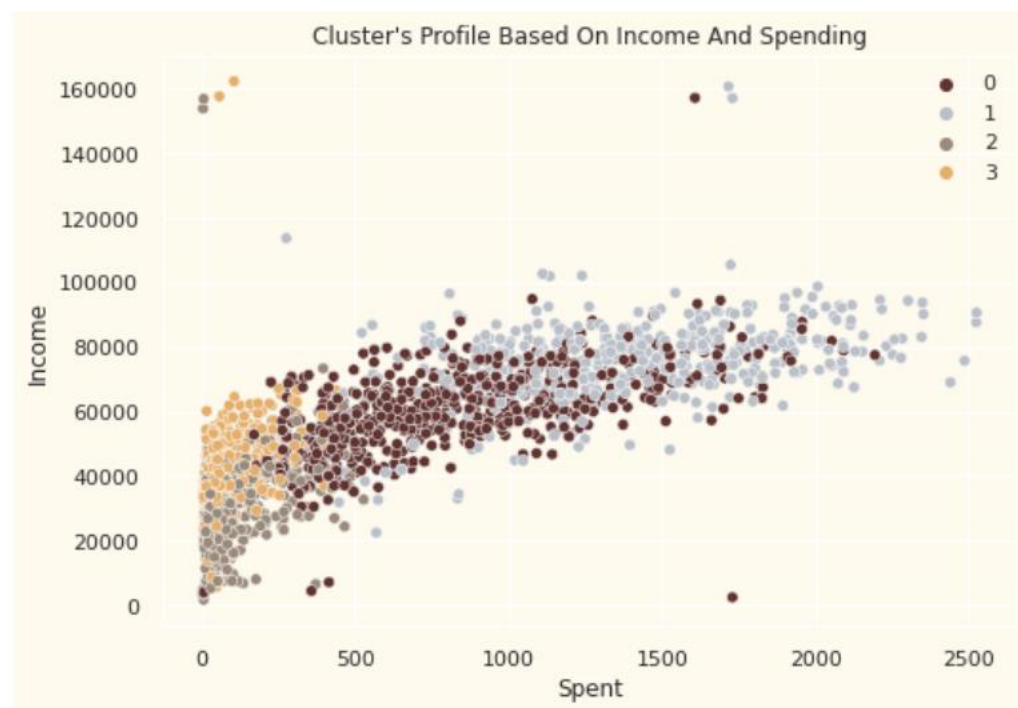
## 3. Evaluating Models

Since this is an unsupervised clustering. We do not have a tagged feature to evaluate or score our model. The purpose of this section is to study the patterns in the clusters formed and determine the nature of the clusters' patterns. For that, we will be having a look at the data in light of clusters via exploratory data analysis and drawing conclusions. **Firstly, let us have a look at the group distribution of clustering**



The clusters seem to be fairly distributed.

## 3. Evaluating Models

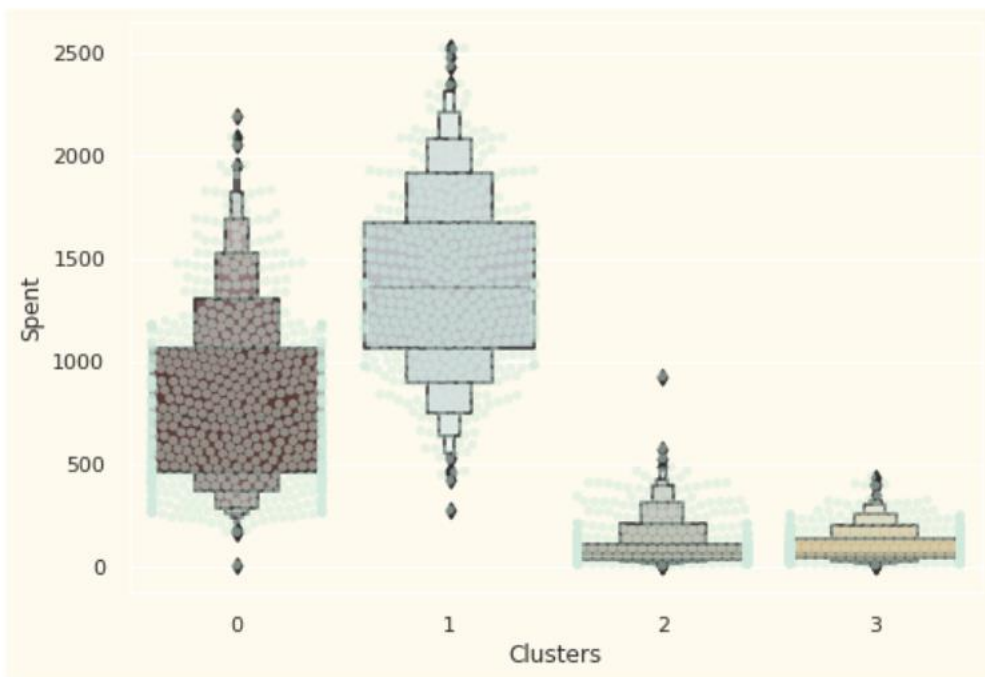


**Income vs spending plot shows the clusters pattern**

- group 0: high spending & average income
- group 1: high spending & high income
- group 2: low spending & low income
- group 3: high spending & low income

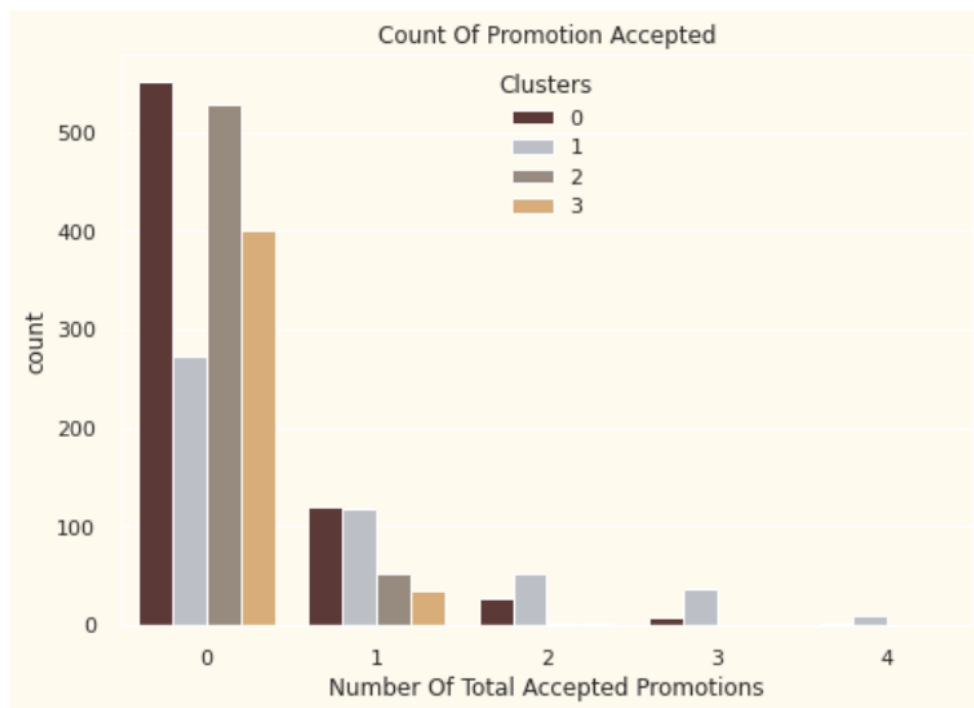


## 3. Evaluating Models



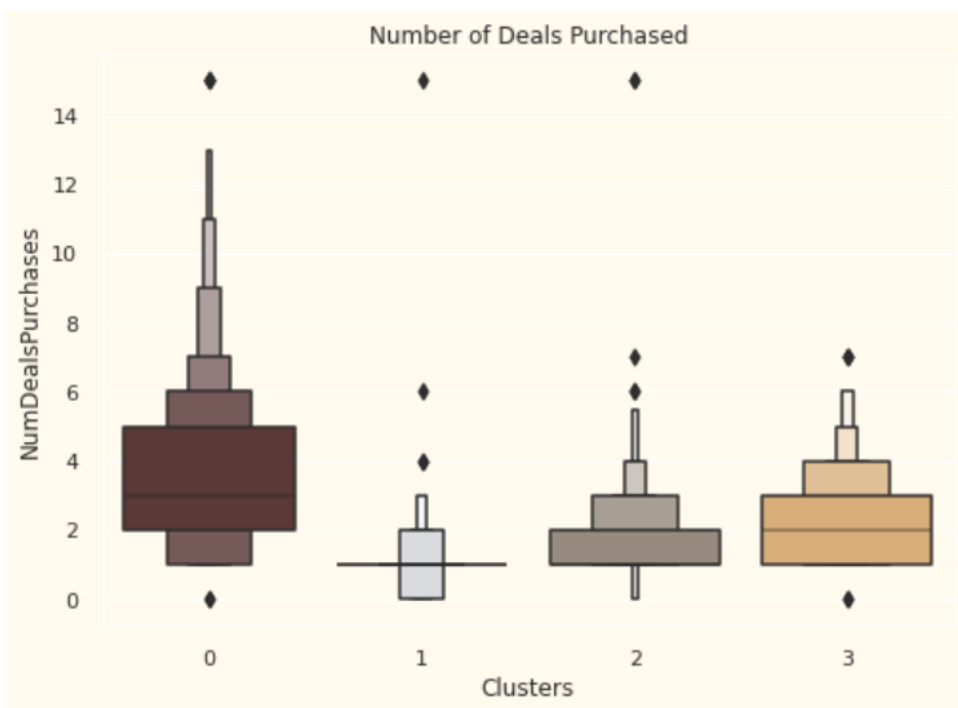
From the above plot, it can be clearly seen that cluster 1 is our biggest set of customers closely followed by cluster 0. We can explore what each cluster is spending on for the targeted marketing strategies.

## 3. Evaluating Models



There has not been an overwhelming response to the campaigns so far. Very few participants overall. Moreover, no one part take in all 5 of them. Perhaps better-targeted and well-planned campaigns are required to boost sales.

## 3. Evaluating Models



Unlike campaigns, the deals offered did well. It has best outcome with cluster 0 and cluster 3. However, our star customers cluster 1 are not much into the deals. Nothing seems to attract cluster 2 overwhelmingly



## 4. Profiling

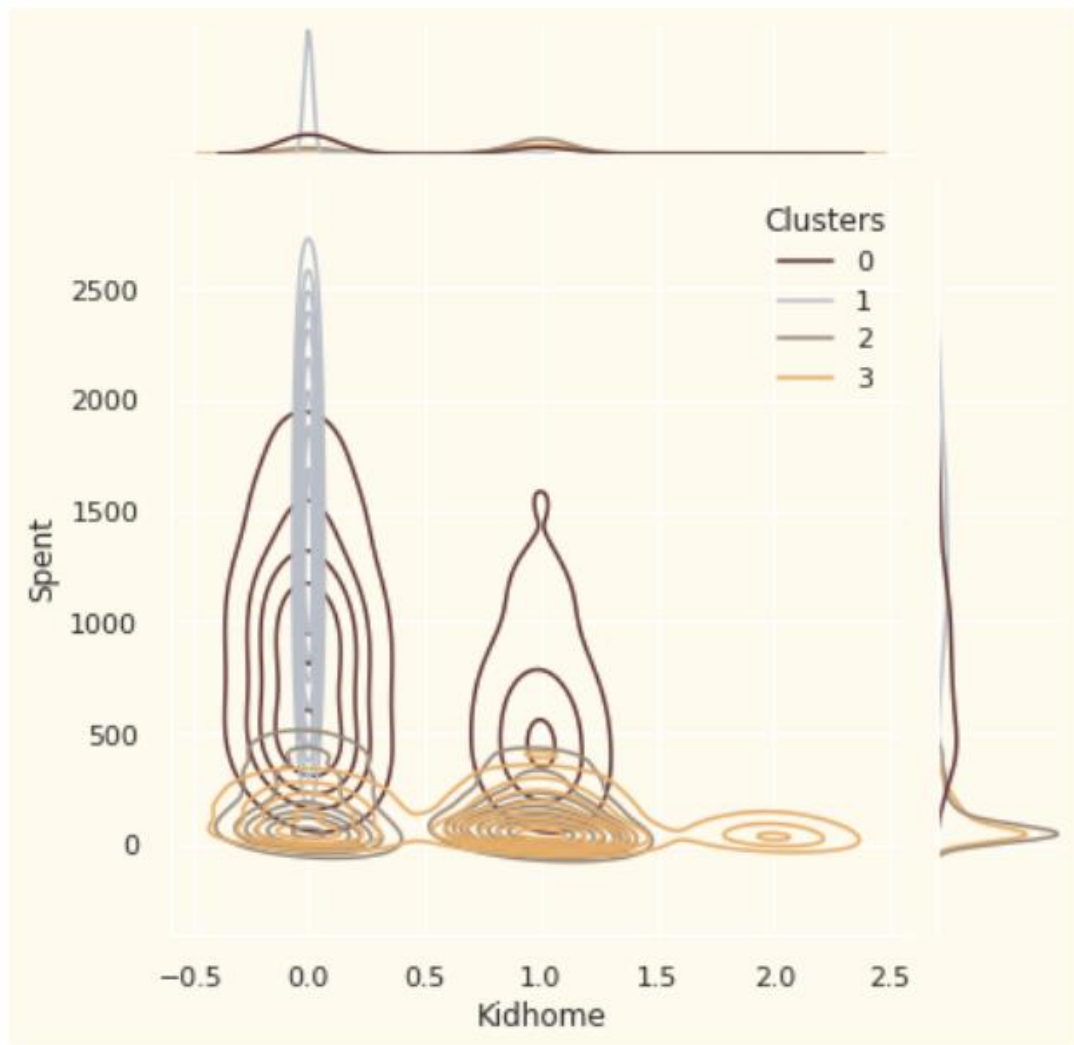
Now that we have formed the clusters and looked at their purchasing habits. Let us see who all are there in these clusters. For that, we will be profiling the clusters formed and come to a conclusion about who is our star customer and who needs more attention from the retail store's marketing team.

To decide that I will be plotting some of the features that are indicative of the customer's personal traits in light of the cluster they are in.

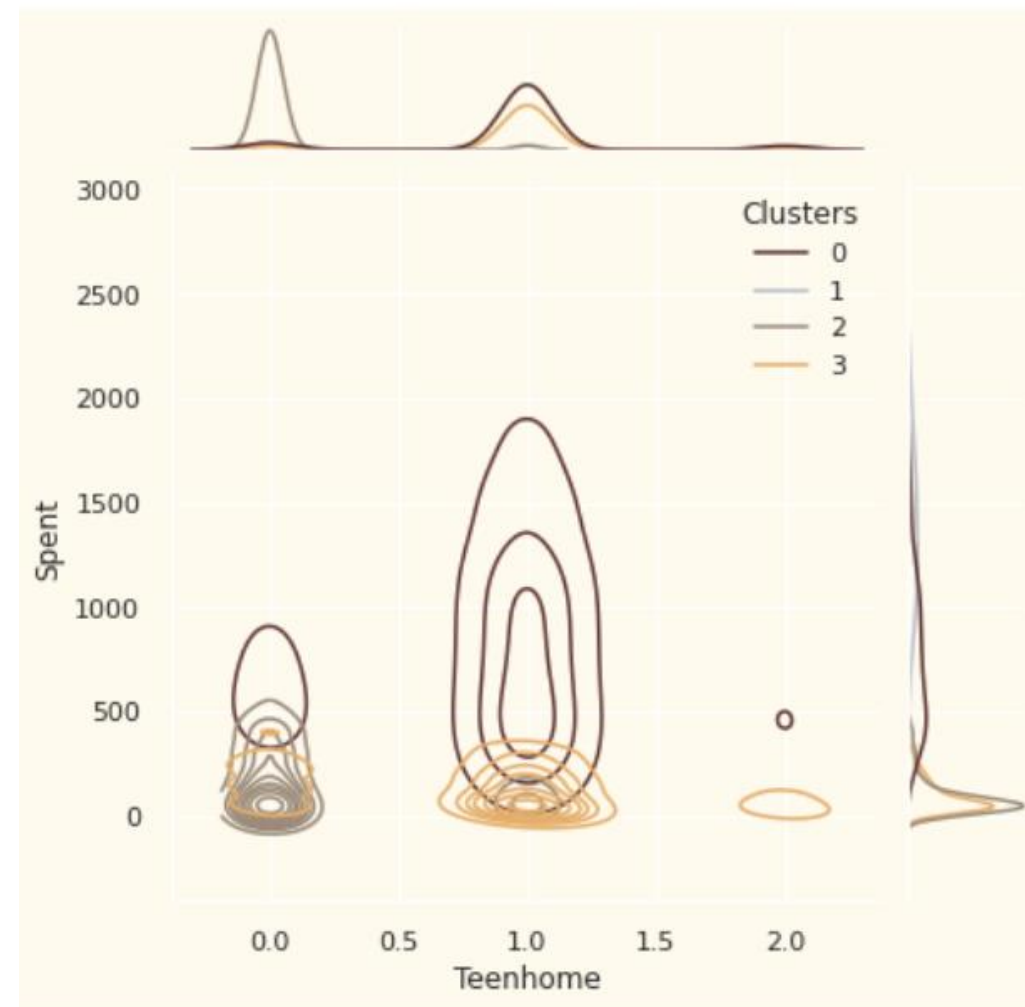
# Default Topic

## 4. Profiling

Kidhome vs Spent



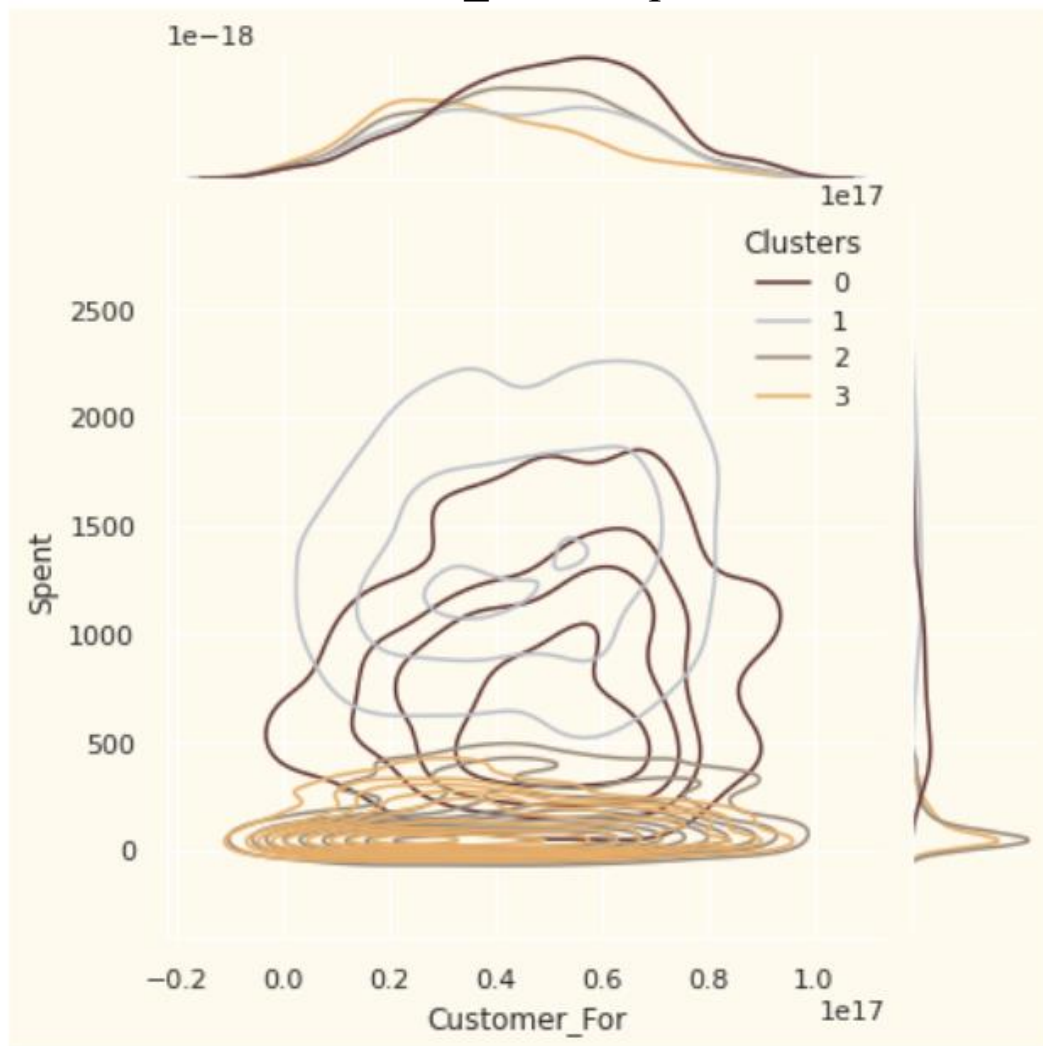
Teenhome vs Spent



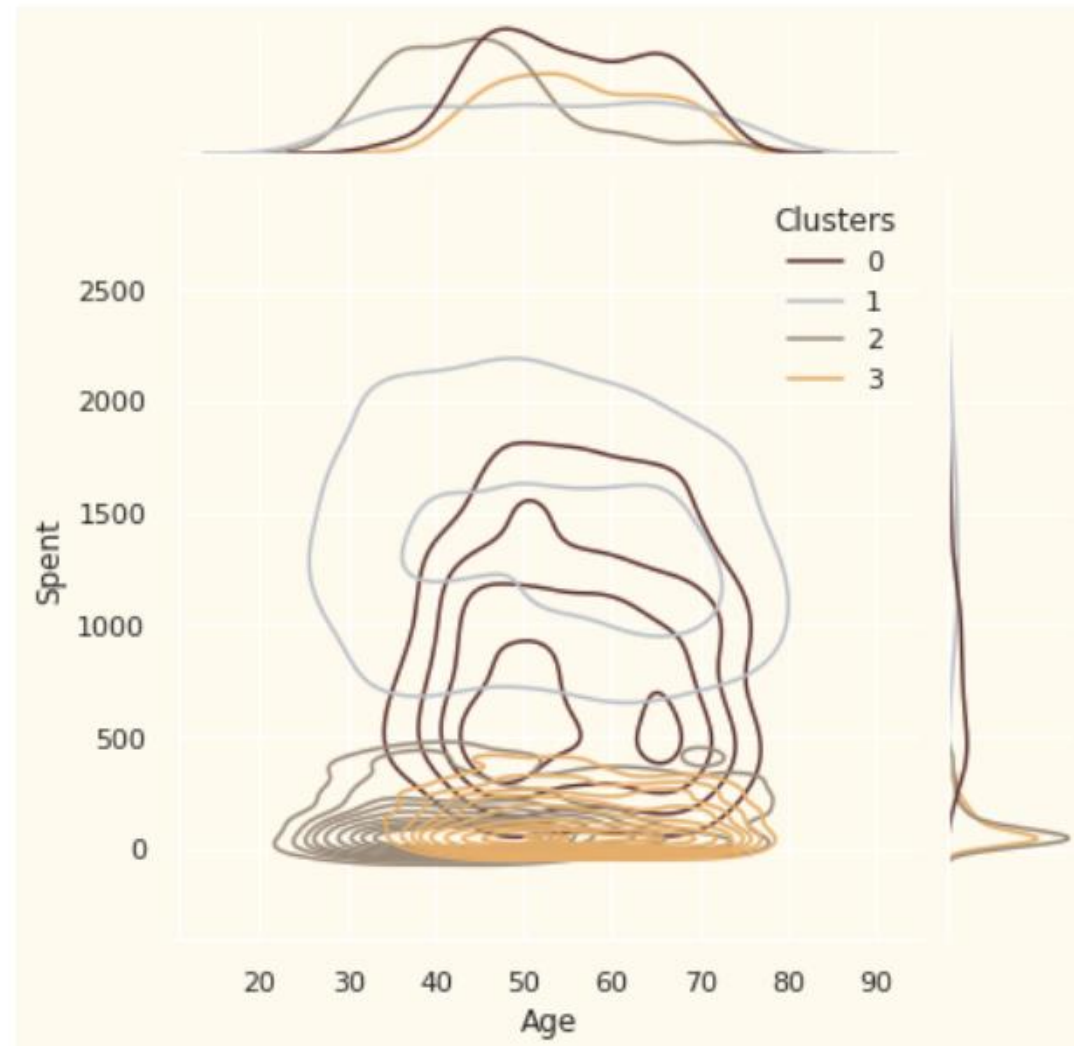
# Default Topic

## 4. Profiling

Customer\_For vs Spent



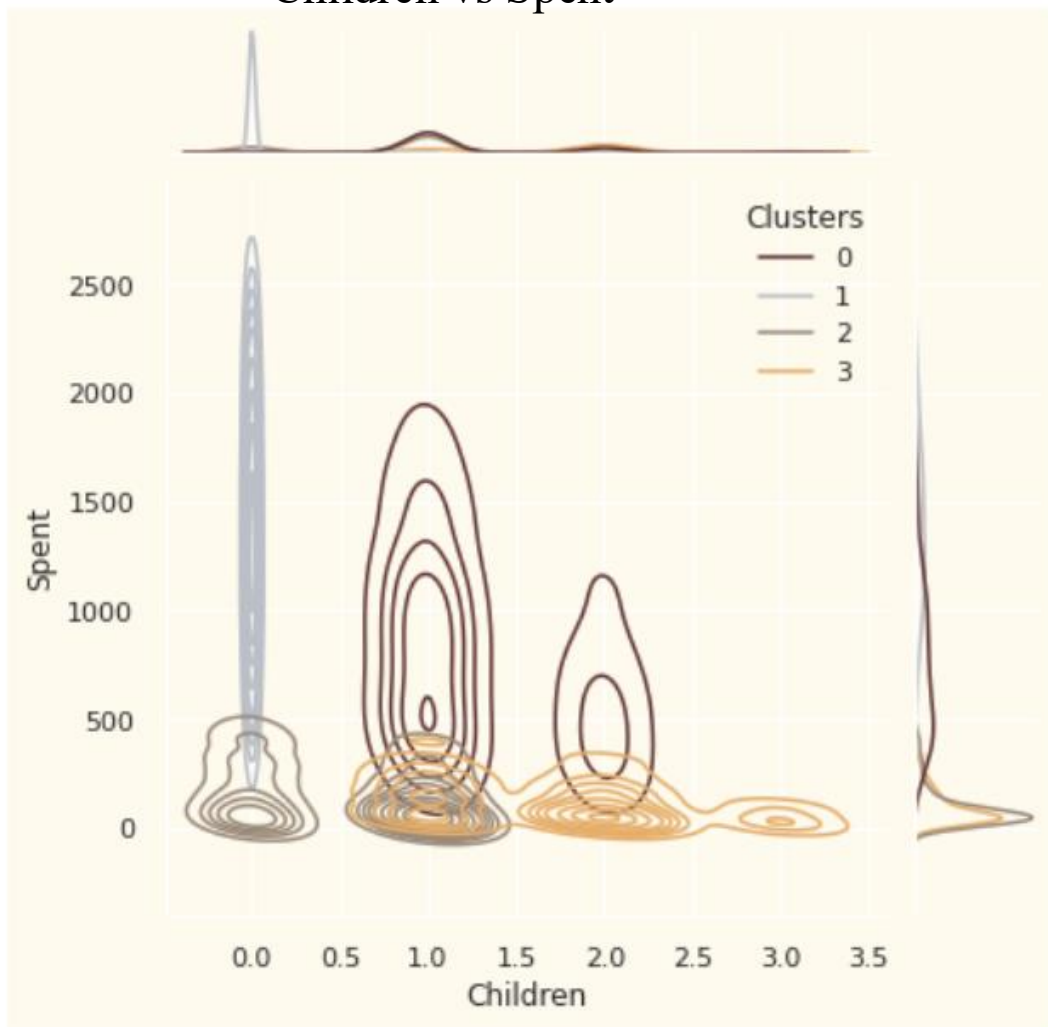
Age vs Spent



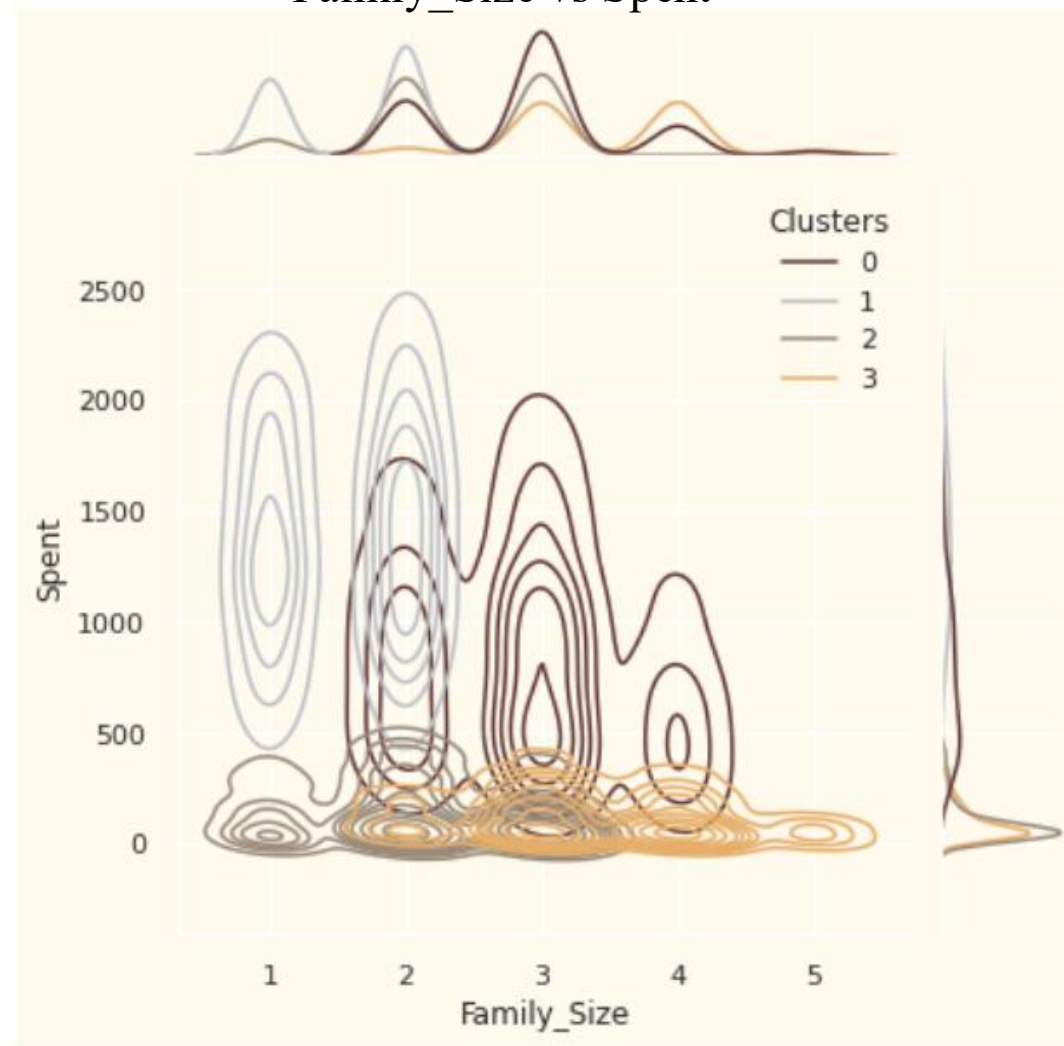
# Default Topic

## 4. Profiling

Children vs Spent



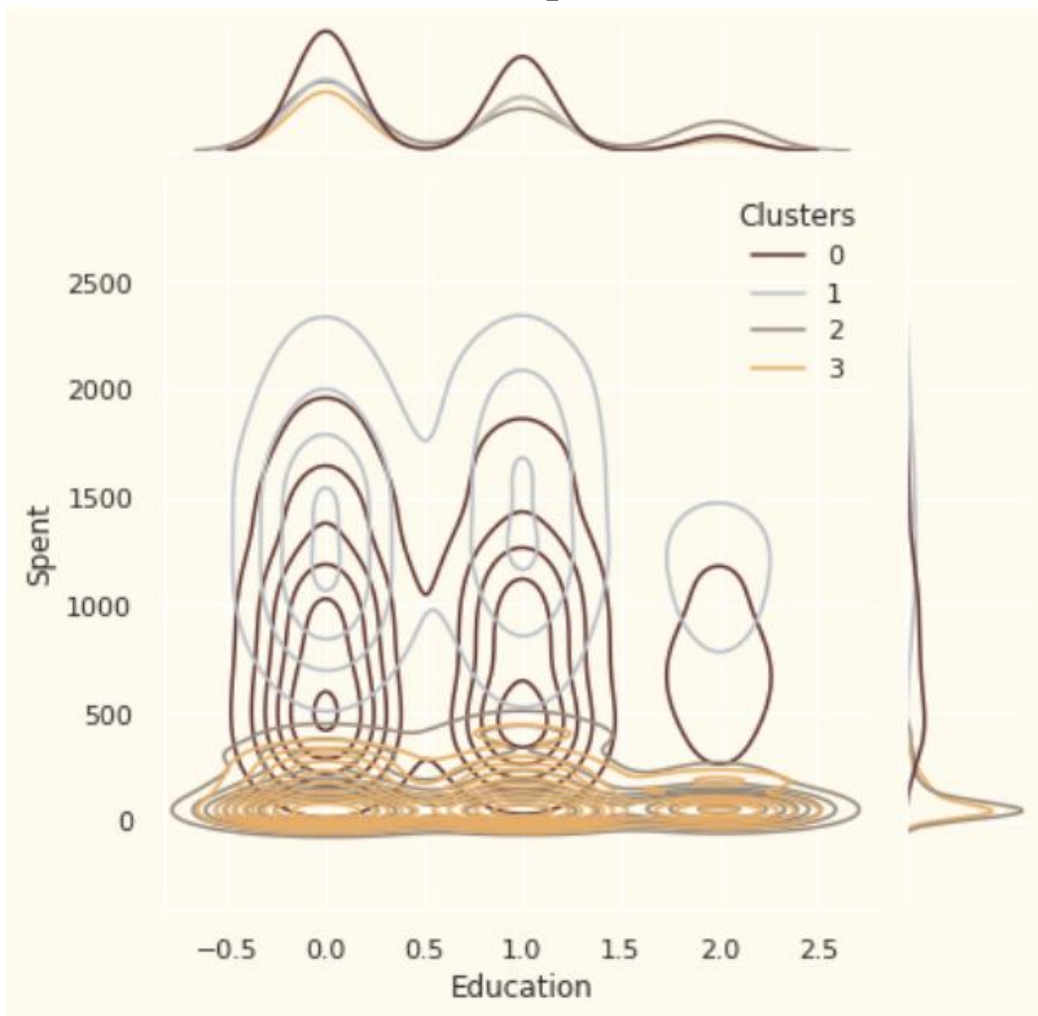
Family\_Size vs Spent



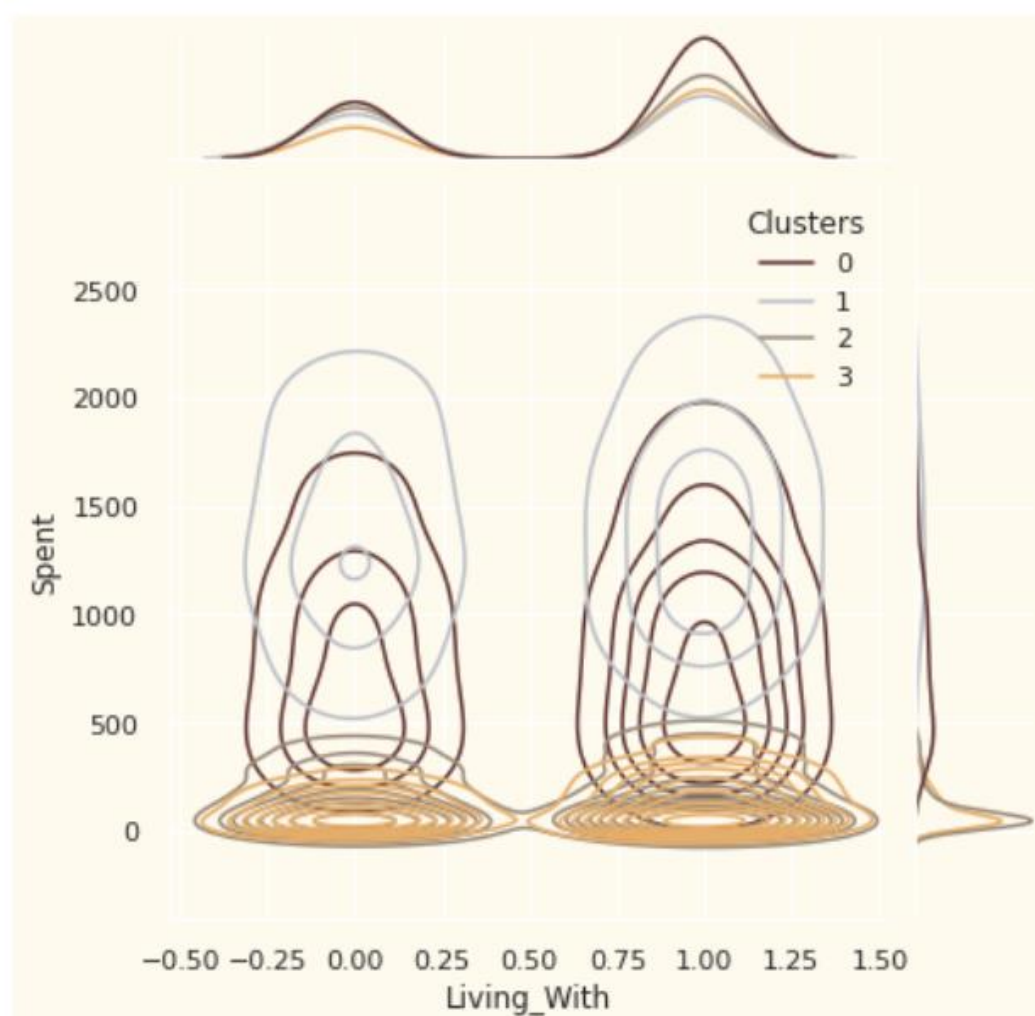
# Default Topic

## 4. Profiling

Education vs Spent

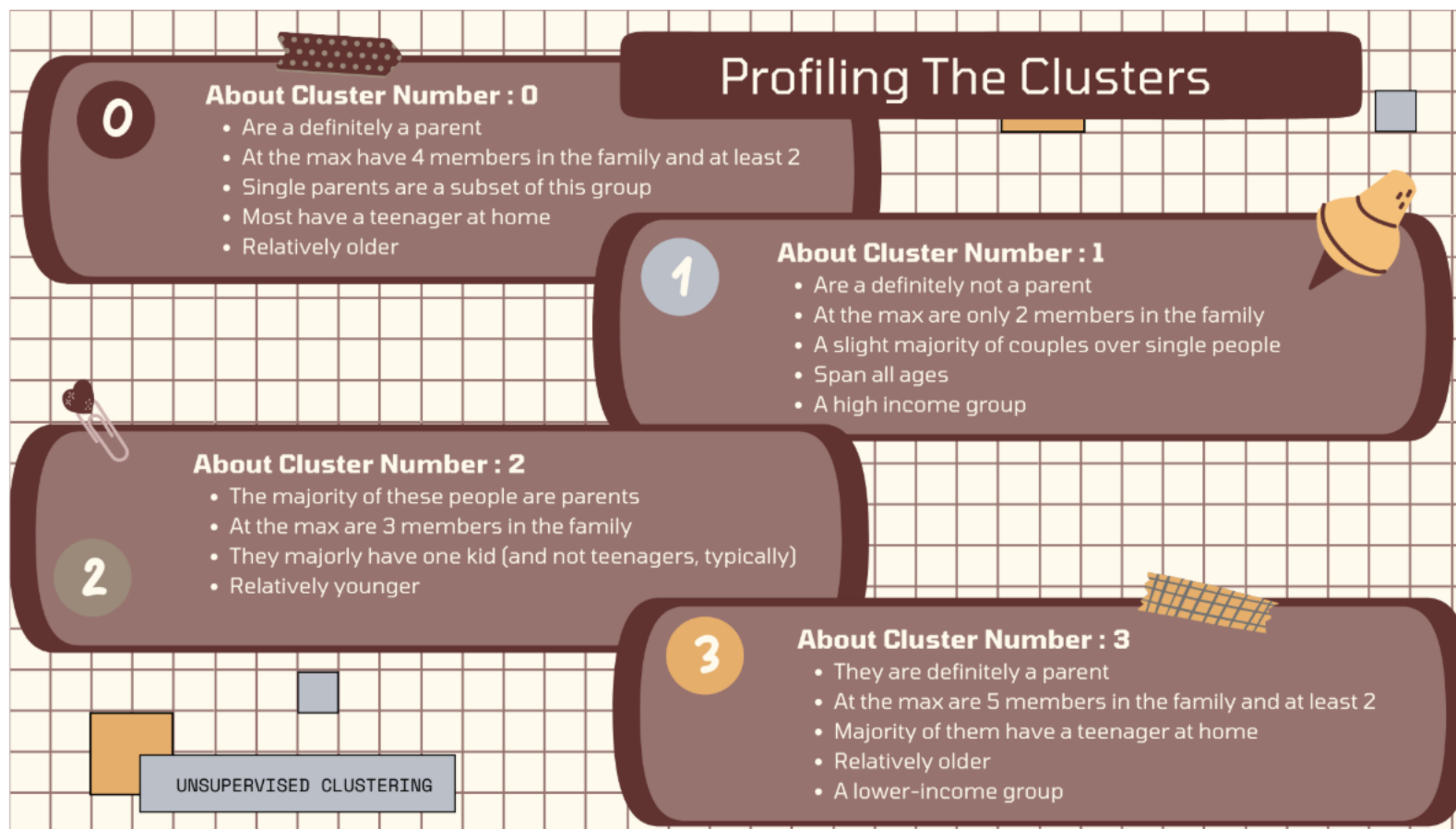


Living\_With vs Spent





## 4. Profiling



# Reminder

---

- This project is an opportunity to think critically and innovate. Please do **NOT** copy existing solution. Get your hands dirty and explore your own ideas.
- **Plagiarism** in any form is strictly prohibited. Violations will result in **0 points**.

## Zebra: When Temporal Graph Neural Networks Meet Temporal Personalized PageRank

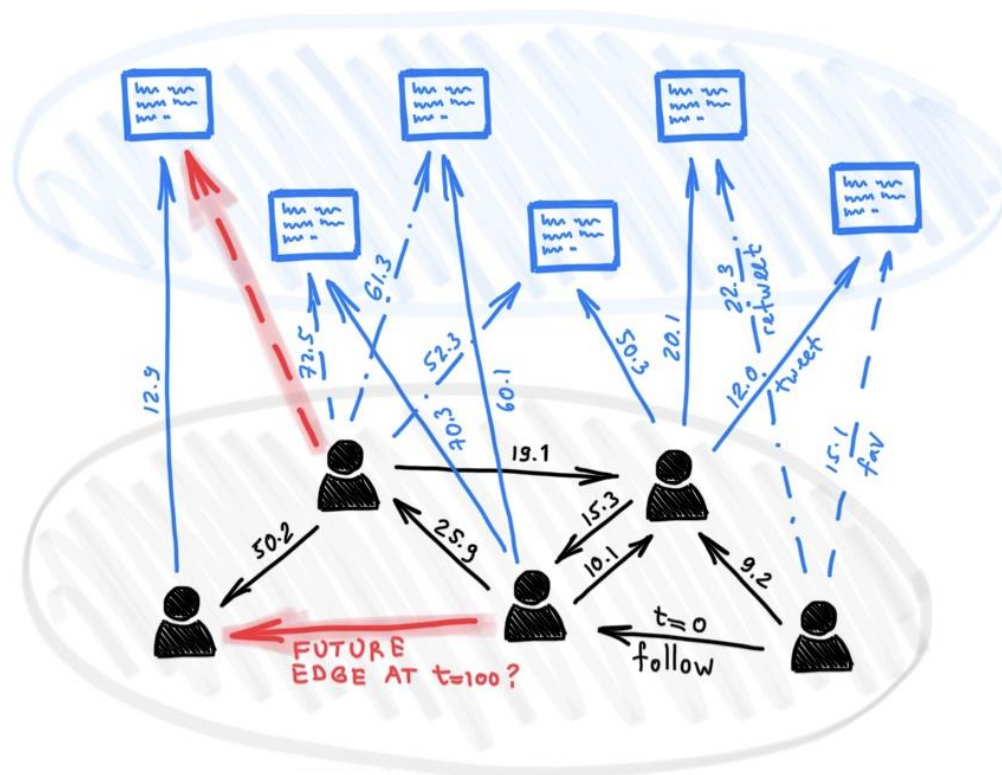
Yiming Li<sup>1</sup>, Yanyan Shen<sup>2</sup>, Lei Chen<sup>1</sup>, Mingxuan Yuan<sup>3</sup>

<sup>1</sup>Hong Kong University of Science and Technology, Hong Kong, China

<sup>2</sup>Shanghai Jiao Tong University, Shanghai, China

<sup>3</sup>Huawei Noah's Ark Lab, Hong Kong, China



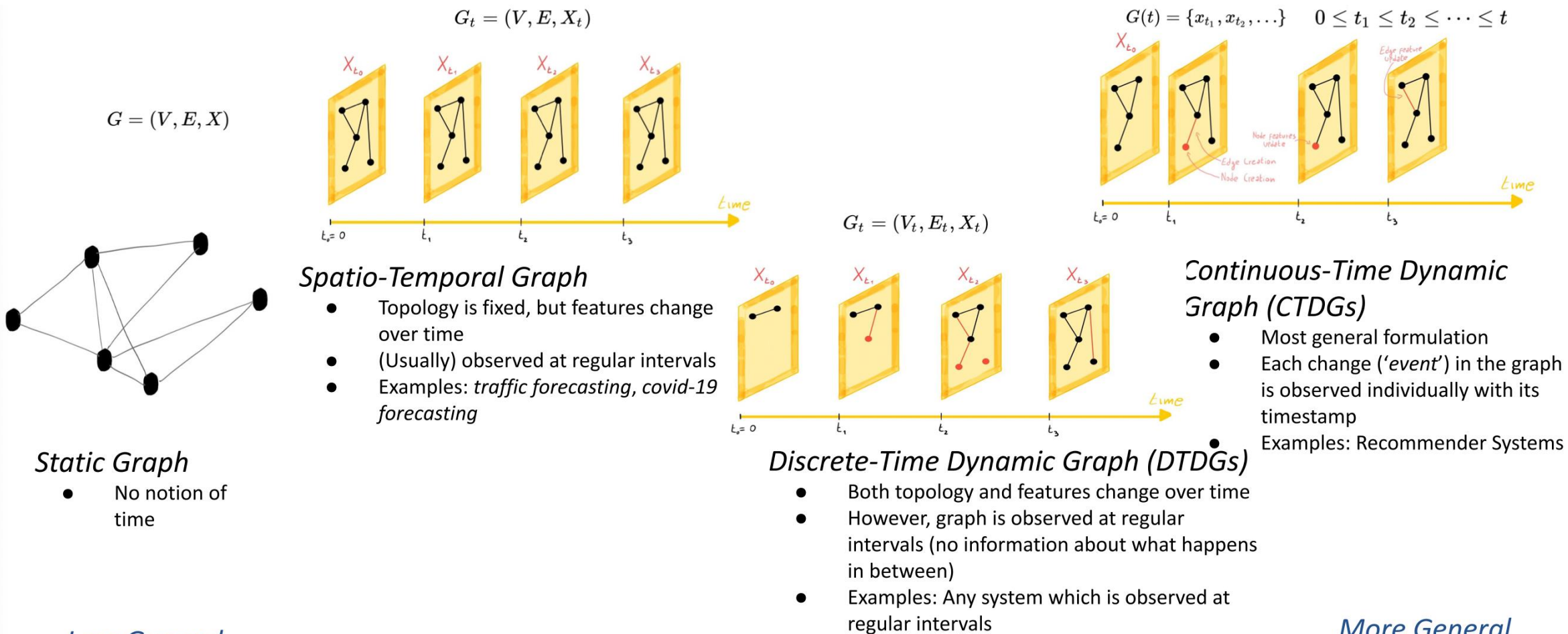


## Applications on dynamic graphs

- Temporal link prediction
- Dynamic node classification
- Social network
- Recommendation
- E-commerce
- Fraud detection

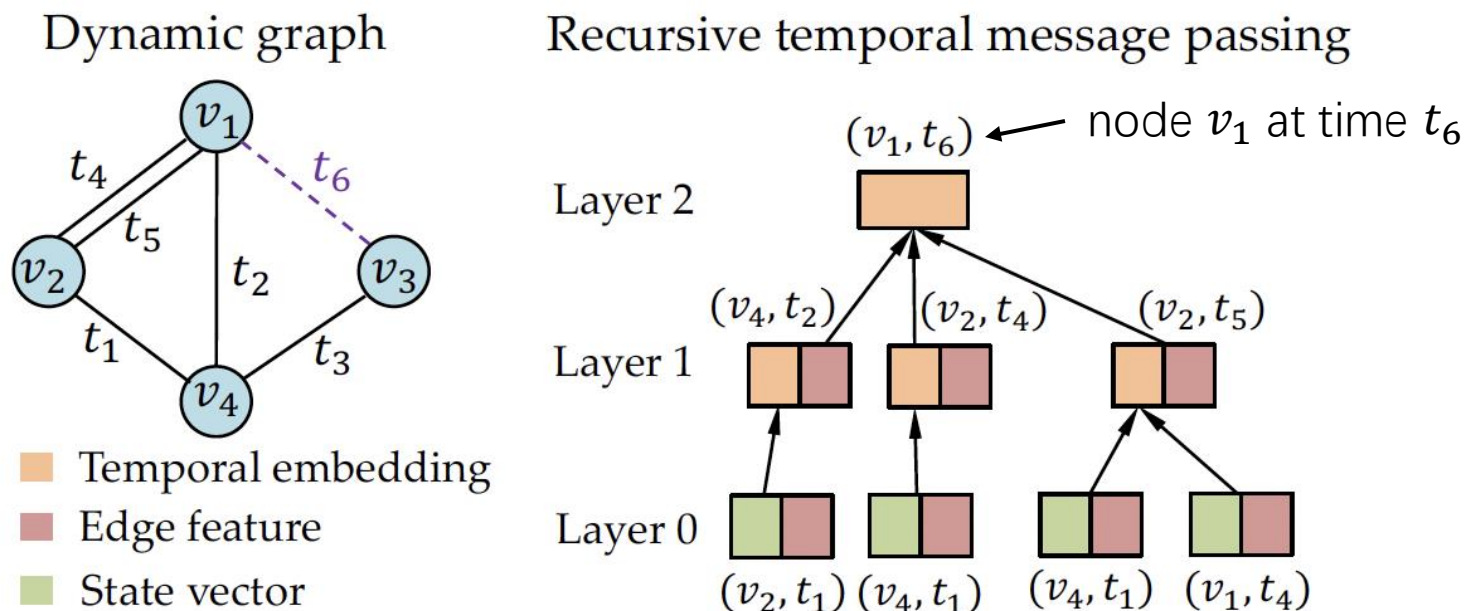
A dynamic network of Twitter users interacting with tweets and following each other.

# Background



- Using a static GNN for dynamic graphs would mean
  - Inefficiency: computation is repeated each time we want to make prediction
  - Loss of information: not able to take into account how the graph evolved
- T-GNNs learn embeddings on **continuous-time** dynamic graph
  - Support addition/deletion of nodes, edges, as well as feature changes
  - Make predictions (e.g., classify a node) at any time point

- T-GNN: given a node  $v$  at timestamp  $t$ , an  $L$ -layer T-GNN computes its temporal embedding via **recursive temporal message passing**

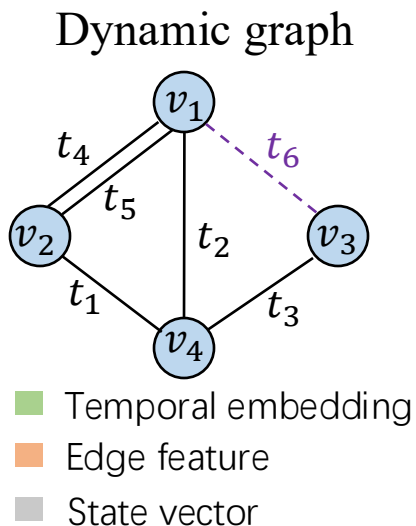


Example of a 2-layer vanilla T-GNN.

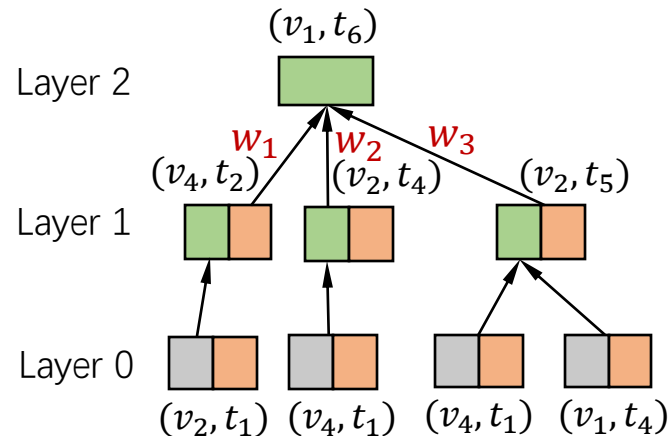
- T-GNN identify influential neighbor nodes via **tedious recursive temporal neighborhood aggregation**
- Do there exist influential nodes that we can utilize to accelerate the computation of T-GNNs without compromising model accuracy?

# T-GNN and Temporal Random Walk

- T-GNN mimics the temporal random walk process on dynamic graphs
- Key observation: **1-hop neighbor weight** implicitly defines the **1-step temporal random walk probability** on dynamic graphs



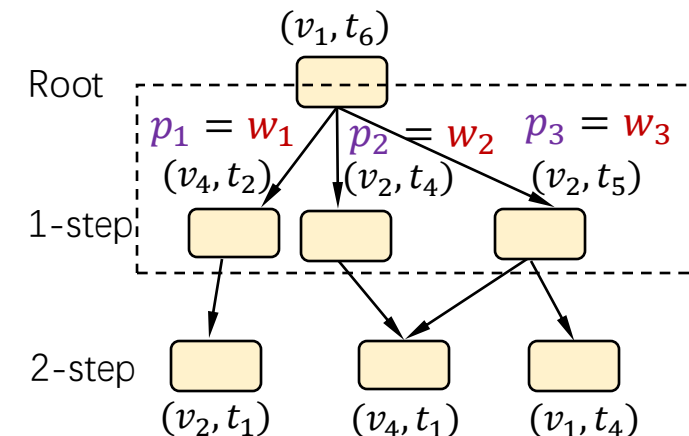
Recursive temporal message passing



1-hop neighbor weights learned by T-GNN

$$w_1 + w_2 + w_3 = 1$$

Temporal random walk

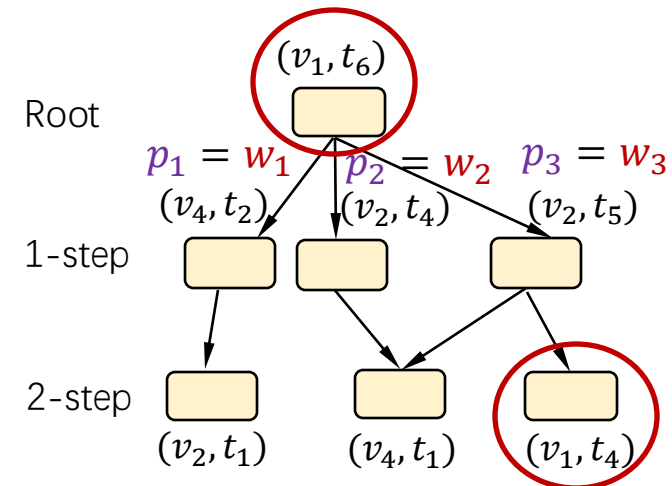
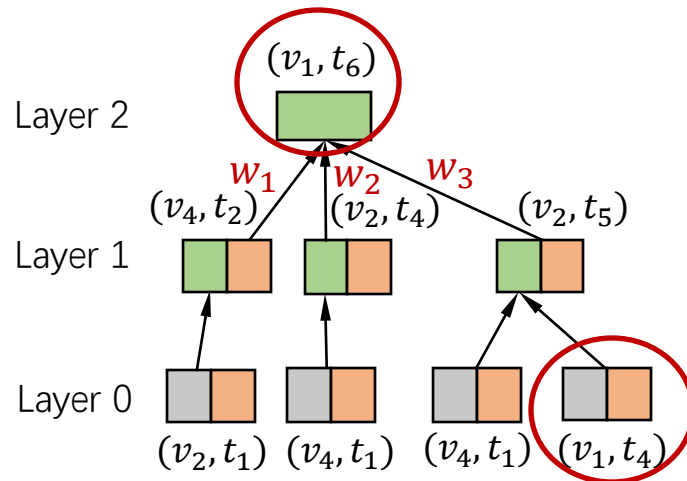


1-step temporal random walk probabilities

$$p_1 + p_2 + p_3 = 1$$

- Implication

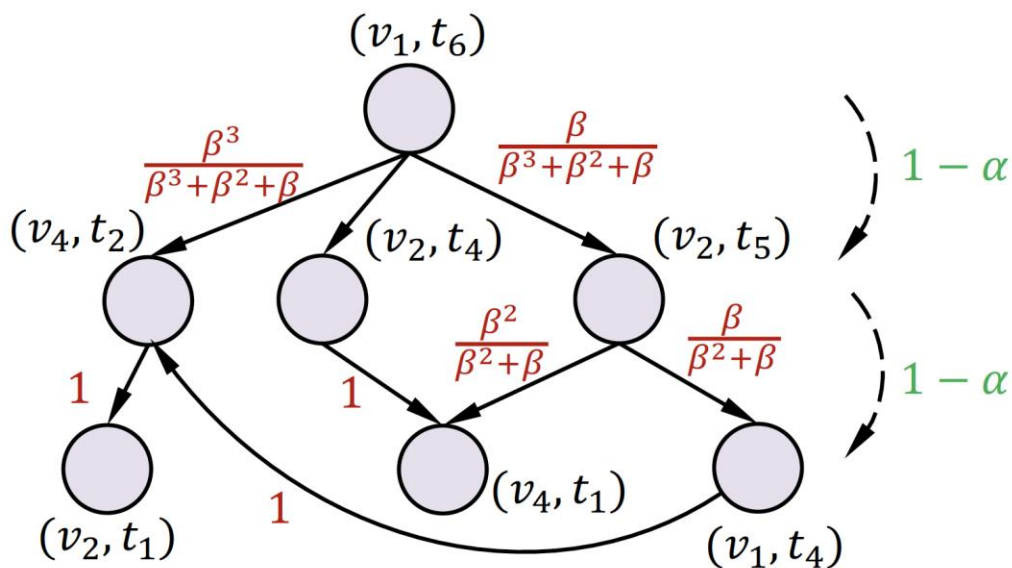
- We can identify influential neighbors if there exists a method that is able to sufficiently capture the **diffusion process** (temporal random walk probabilities) on dynamic graphs
- In this way, we can **prune insignificant neighbors** and avoid the laborious process of recursive temporal message passing



Challenge: As we cannot know the neighbor weights in advance without running T-GNN, how to estimate the corresponding temporal random walk probabilities?



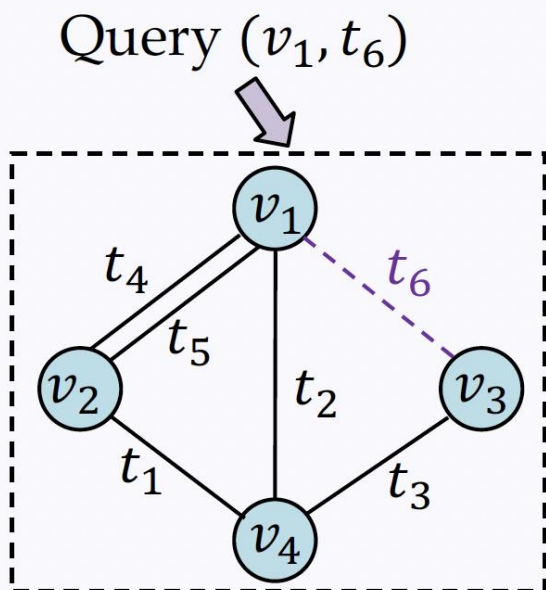
- Temporal Personalized PageRank (T-PPR) is a parameterized metric for estimating the influence score of temporal nodes on dynamic graphs
  - The T-PPR value of  $(j, \tau)$  with respect to  $(i, t)$  is defined as the probability that an  $\alpha$ -temporal random walk with  $\beta$ -exponential decay starting from  $(i, t)$  terminates at  $(j, \tau)$



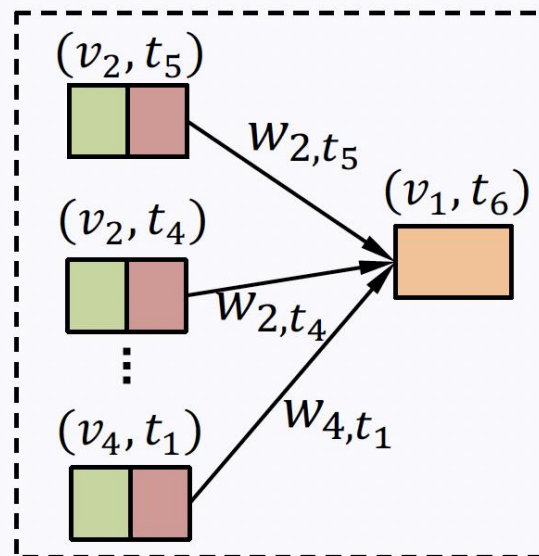


- Zebra framework avoids the tedious recursive temporal message passing

Step 1: Top- $k$  T-PPR query



Step 2: Aggregation



- Existing techniques for PPR cannot be adapted to efficiently solve the top- $k$  T-PPR query due to the highly imbalanced weight distribution
- A straightforward solution — —power iteration
  - Recursively traverses the graph and then selects the top- $k$  temporal nodes
  - However, this approach is computationally prohibitive

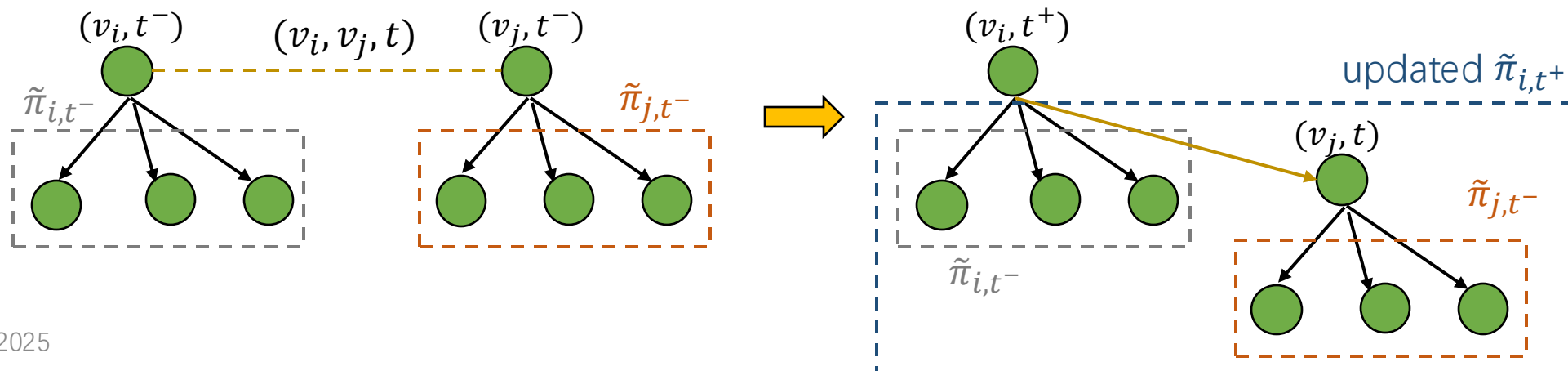
## Single-scan Top-k T-PPR Algorithm (SANTA)

- Data structure

- For each node  $v_i$  SANTA maintains  $\tilde{\pi}_{i,t^-}$ , the latest approximated top-k T-PPR dictionary for node  $v_i$  before the time current timestamp  $t$

- Single-edge update

- Given a new interaction  $(v_i, v_j, t)$ , SANTA updates the top-k T-PPR dictionary for node  $v_i$  by **merging**  $\tilde{\pi}_{i,t^-}$  and  $\tilde{\pi}_{j,t^-}$  in a streaming fashion



- Additive error bound: the additive errors are negligible when the temporal random walk probabilities are **highly skewed** or a **large  $k$**  is adopted

THEOREM 7 (ADDITIVE ERROR BOUND). *Given a new interaction  $\gamma(t) = (v_i, v_j, e_{ij}(t), t)$ , let  $\tilde{\pi}_{i,t^+}$  be the top- $k$  T-PPR dictionary returned by SANTA. For any temporal node  $(z, \tau)$ , we have*

$$\tilde{\pi}_{i,t^+}(z, \tau) \leq \pi_{i,t^+}(z, \tau) \leq \tilde{\pi}_{i,t^+}(z, \tau) + \frac{\theta_{\max}^{(k)}}{\alpha(1-\beta)}, \quad (16)$$

*where  $\theta_{\max}^{(k)}$  is the maximum possible top- $k$  threshold for the T-PPR metric parameterized by  $\alpha$  and  $\beta$ .*

Statistics of dynamic graphs.

| Dataset        | $ V $     | $ E $     | $d_v$ | $d_e$ | Timespan  |
|----------------|-----------|-----------|-------|-------|-----------|
| MOOC [22]      | 7,144     | 411,749   | 172   | 4     | 30 days   |
| Wikipedia [22] | 9,227     | 157,474   | 172   | 172   | 30 days   |
| Reddit [22]    | 10,984    | 672,447   | 172   | 172   | 30 days   |
| AskUbuntu [1]  | 159,316   | 964,437   | 172   | 0     | 2613 days |
| SuperUser [2]  | 194,085   | 1,443,339 | 172   | 0     | 2773 days |
| Wiki-Talk [4]  | 1,140,149 | 7,833,140 | 172   | 0     | 2320 days |

- **Zebra** identifies essential neighbors via SANTA (Algorithm 1) and computes node embeddings through 1-layer neighborhood aggregation (Eq.11). In addition, Zebra updates the state vectors of nodes upon new interactions as vanilla T-GNN [36].
- JODIE [22] uses RNNs to propagate the information of temporal interactions to update node representations.
- TGAT [54] mimics the message passing scheme of static GNNs and encodes time information through random Fourier features.
- TGN [36] dynamically maintains a state vector for each node and includes previous methods [22, 41, 54] as special cases.
- CAW [50] encodes temporal neighborhood information using anonymized random walk and attention modules.
- APAN [49] accelerates model inference by decoupling graph computation and message propagation.



- Zebra can be orders of magnitude faster than the state-of-the-art baselines without sacrificing model performance

Table 4: Comparison of T-GNNs on small dynamic graphs. Zebra is implemented as an ensemble of two top-20 T-PPR metrics. We report model performance in transductive average precision (%), inductive average precision (%), and total training time (s). In addition, we report the number of epochs required for model convergence in parentheses. The best and second-best results in each metric are marked in bold and underlined, respectively. “\*” indicates that the performance improvement of Zebra over the best-performing baseline is statistically significant with the significance level set to be 0.05.

| Model | Wikipedia                         |                                    |                    | Reddit                             |                                    |                     | MOOC                               |                                   |                    |
|-------|-----------------------------------|------------------------------------|--------------------|------------------------------------|------------------------------------|---------------------|------------------------------------|-----------------------------------|--------------------|
|       | Trans AP                          | Induct AP                          | Training (#)       | Trans AP                           | Induct AP                          | Training (#)        | Trans AP                           | Induct AP                         | Training (#)       |
| JODIE | 95.16 $\pm$ 0.4                   | 93.13 $\pm$ 0.5                    | 2356.4 (18)        | 95.83 $\pm$ 0.3                    | 93.20 $\pm$ 0.4                    | <u>9243.64 (14)</u> | 83.26 $\pm$ 0.5                    | 81.77 $\pm$ 0.4                   | 5160.30 (15)       |
| TGAT  | 94.26 $\pm$ 0.1                   | 92.88 $\pm$ 0.3                    | 2881.4 (29)        | 97.80 $\pm$ 0.2                    | 96.08 $\pm$ 0.3                    | 11933.8 (21)        | 70.22 $\pm$ 0.4                    | 70.83 $\pm$ 0.5                   | 7838.5 (25)        |
| TGN   | <u>98.58 <math>\pm</math> 0.1</u> | 98.05 $\pm$ 0.1                    | 2182.7 (26)        | <u>98.66 <math>\pm</math> 0.1</u>  | 97.55 $\pm$ 0.1                    | 11528.9 (26)        | <u>88.88 <math>\pm</math> 1.7</u>  | <u>88.17 <math>\pm</math> 2.1</u> | 5086.83 (21)       |
| APAN  | 96.41 $\pm$ 0.5                   | 96.06 $\pm$ 0.4                    | <u>1605.0 (21)</u> | 98.50 $\pm$ 0.2                    | 97.62 $\pm$ 0.7                    | 16431.8 (18)        | 87.02 $\pm$ 0.3                    | 86.74 $\pm$ 0.5                   | <u>3374.1 (14)</u> |
| CAW   | 98.18 $\pm$ 0.1                   | <u>98.24 <math>\pm</math> 0.1</u>  | 10175.5 (13)       | 98.54 $\pm$ 0.1                    | <u>97.97 <math>\pm</math> 0.1</u>  | 75585.1 (16)        | 80.60 $\pm$ 0.4                    | 80.18 $\pm$ 0.4                   | 34063.9 (14)       |
| Zebra | <b>98.67 <math>\pm</math> 0.1</b> | <b>98.59 <math>\pm</math> 0.1*</b> | <b>302.6 (34)</b>  | <b>98.76 <math>\pm</math> 0.1*</b> | <b>98.28 <math>\pm</math> 0.1*</b> | <b>1342.0 (31)</b>  | <b>92.45 <math>\pm</math> 0.2*</b> | <b>89.56 <math>\pm</math> 0.3</b> | <b>619.5 (25)</b>  |

- Zebra is also memory efficient
  - It reduces the computational cost of T-GNNs and thus save the GPU memory required to perform the corresponding computations

Table 5: Comparison of T-GNNs on large dynamic graphs. Zebra is implemented as an ensemble of two top-20 T-PPR metrics. The evaluation metrics are the same as those in Table 4. Particularly, “OOM” denotes out-of-memory error, and “TLE” denotes time limit exceed such that we cannot finish one epoch of model training in 12 hours. “\*” indicates that the performance improvement of Zebra over the best-performing baseline is statistically significant with the significance level set to be 0.05.

| Model | AskUbuntu                         |                                     |                    | SuperUser                         |                                     |                    | Wiki-Talk                           |                                     |                    |
|-------|-----------------------------------|-------------------------------------|--------------------|-----------------------------------|-------------------------------------|--------------------|-------------------------------------|-------------------------------------|--------------------|
|       | Trans AP                          | Induct AP                           | Training (#)       | Trans AP                          | Induct AP                           | Training (#)       | Trans AP                            | Induct AP                           | Training (#)       |
| JODIE | OOM                               | OOM                                 | OOM                | OOM                               | OOM                                 | OOM                | OOM                                 | OOM                                 | OOM                |
| TGAT  | $87.57 \pm 0.3$                   | $84.21 \pm 0.4$                     | 11728.5 (21)       | $86.40 \pm 0.5$                   | $83.12 \pm 0.5$                     | 17129.3 (19)       | $91.72 \pm 0.2$                     | $85.38 \pm 0.3$                     | 189420 (34)        |
| TGN   | <b><math>94.51 \pm 0.2</math></b> | <u><math>92.73 \pm 0.2</math></u>   | 36202.6 (20)       | <u><math>93.18 \pm 0.3</math></u> | <u><math>91.76 \pm 0.2</math></u>   | 81747.6 (24)       | OOM                                 | OOM                                 | OOM                |
| APAN  | $89.17 \pm 0.1$                   | $88.43 \pm 0.1$                     | 21606.1 (14)       | $87.07 \pm 0.3$                   | $85.50 \pm 0.5$                     | 48724.2 (19)       | TLE                                 | TLE                                 | TLE                |
| CAW   | $90.87 \pm 0.2$                   | $90.63 \pm 0.2$                     | 61903.7 (14)       | $88.92 \pm 0.2$                   | $88.32 \pm 0.1$                     | 111744.8 (15)      | TLE                                 | TLE                                 | TLE                |
| Zebra | <u><math>94.47 \pm 0.1</math></u> | <b><math>97.91 \pm 0.1^*</math></b> | <b>1362.2 (24)</b> | <b><math>93.21 \pm 0.3</math></b> | <b><math>97.93 \pm 0.1^*</math></b> | <b>2095.1 (23)</b> | <b><math>95.45 \pm 0.1^*</math></b> | <b><math>97.96 \pm 0.1^*</math></b> | <b>9909.6 (16)</b> |

# Report Contents

---

- Abstract
- Introduction
- Preliminary
- Methodology
- Related Work
- Experiments
- Conclusion



- Content you need to submit:
  - Source code
  - Report (in required format)
  - Presentation video & PPT
- **Pay attention:**
  - Clearly state your **project title** and the **specific course topic** your project relates to.
  - Arrive **on time** on your presentation day (even ID/odd ID)
  - Plagiarism in any form is strictly prohibited.

Thanks!

Q & A