

ML Developers for Real Estate Developers

Ekapol Chuangsuwanich

Joint work with Parichat Chonwiharnphan, Pipop Thienprapasith, Proadpran Punyabukkana, Atiwong Suchato, Naruemon Pratanwanich, Ekkalak Leelasornchai, Nattapat Boonprakong, Panthon Imemkamon

About me

Lecturer at Chulalongkorn University

CHULA Σ ENGINEERING
Foundation toward Innovation

COMPUTER

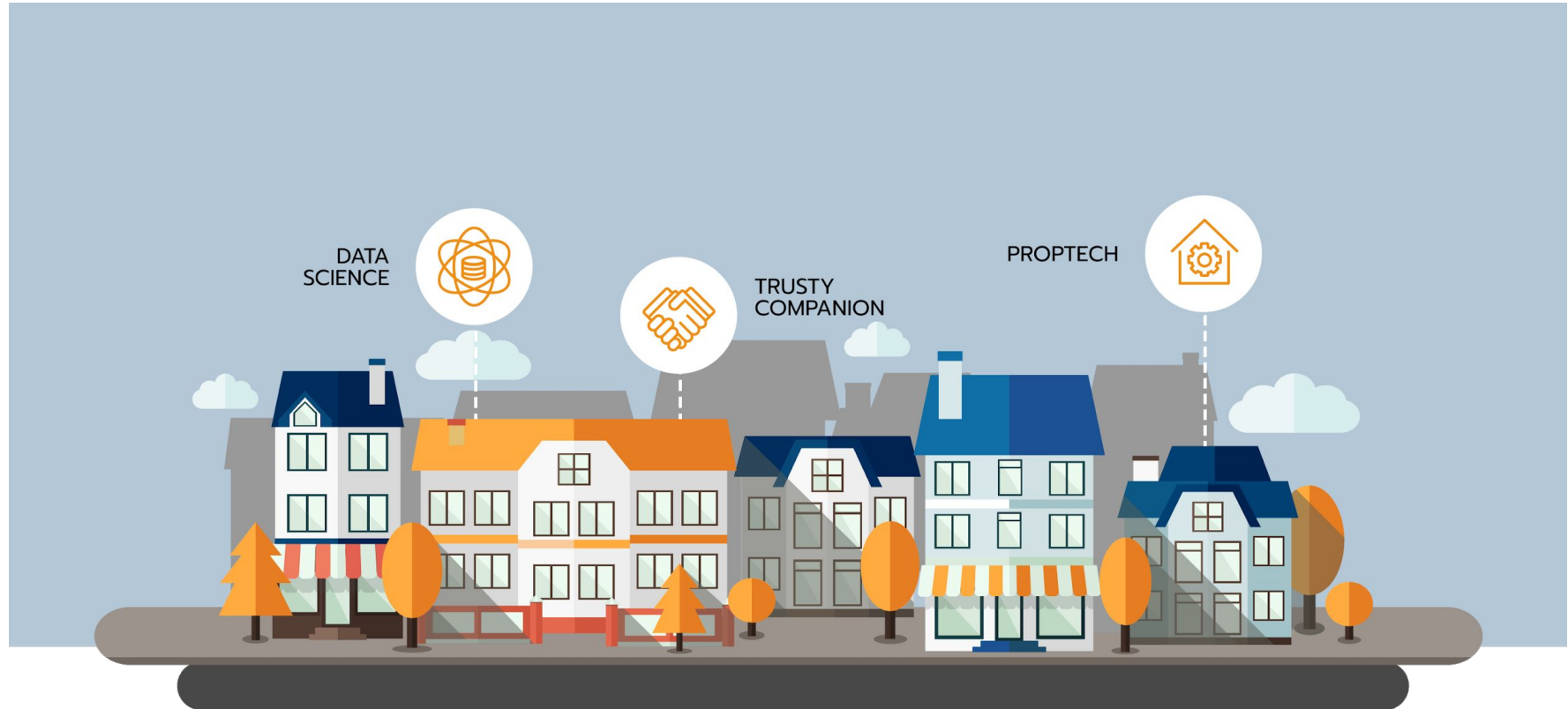
Research focus: ASR, NLP, Bioinformatics, or anything interesting

Various industry collaborations

Ex-intern Google Speech team, a tensorflow fanboy



About HomeDotTech



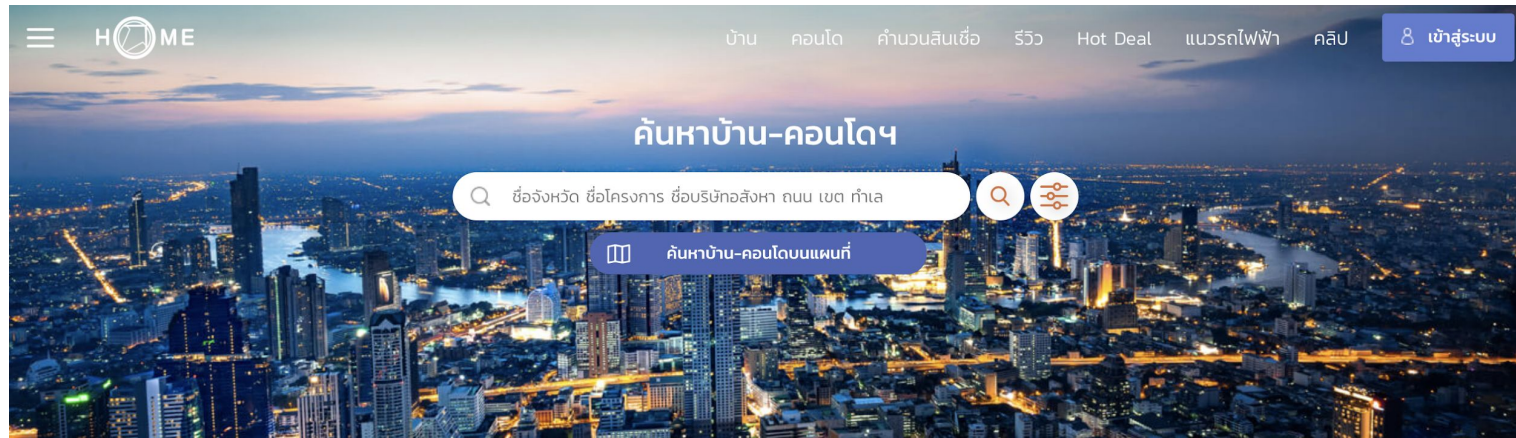
About HomeDotTech

Part of Home Buyer's Group

<http://home.co.th>

One of the most visited Real Estate Listings website in Thailand

~2,000,000 page views per month



Real Estate

The most expensive purchase for most people

Little prior experience

Top complaints to the Office of the Consumer Protection Board (สคบ.)

Homedottech's mission is to help with the home buying process.



Data science for Real Estate

Consumer

Matching

Social listening

(Real Estate) Developers

Lead generation and smart marketing

Social listening

Project development

Customer segmentation

Trend prediction

Pricing



Data science for Real Estate

Consumer

Matching

Social listening

(Real Estate) Developers

Lead generation and smart marketing

Social listening

Project development

Customer segmentation

Trend prediction

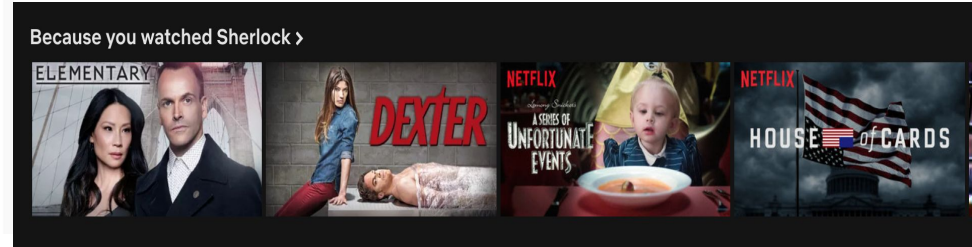
Pricing



Recommendation systems

Goal: predict user's preference toward an item

Related to items you've viewed [See more](#)













Top Picks for Panthon



Information for recommendation systems

There's many information available

Product and info

| |  |  |  |  |  |
|---|---|---|---|---|--|
|  | 1 | 1 | 1 | 1 | 0 |
|  | 0 | 1 | 0 | 1 | 0 |
|  | 1 | 1 | 0 | 1 | 0 |
|  | 0 | 0 | 1 | 0 | 1 |
|  | 1 | 1 | 1 | 0 | 1 |








Project features

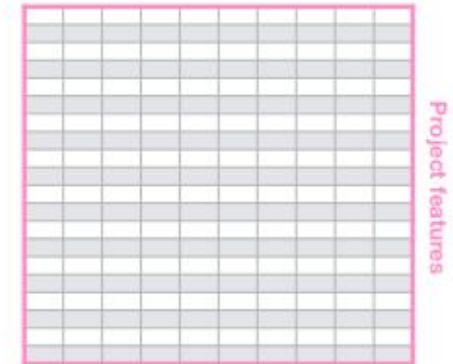
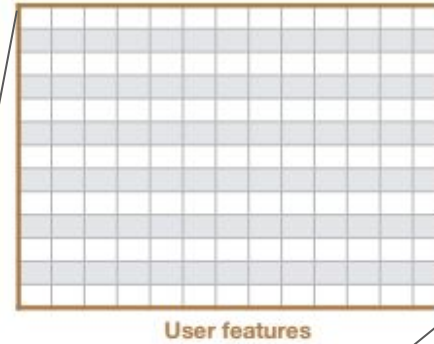
Information for recommendation systems

There's many information available

Product and info

User and info

| | | | | |
|---|---|---|---|---|
|  |  |  |  |  |
| | 1 | 1 | 1 | 0 |
|  | 0 | 1 | 0 | 1 |
|  | 0.1 | 0.8 | 0.9 | 0.5 |



Information for recommendation systems

There's many information available

Product and info

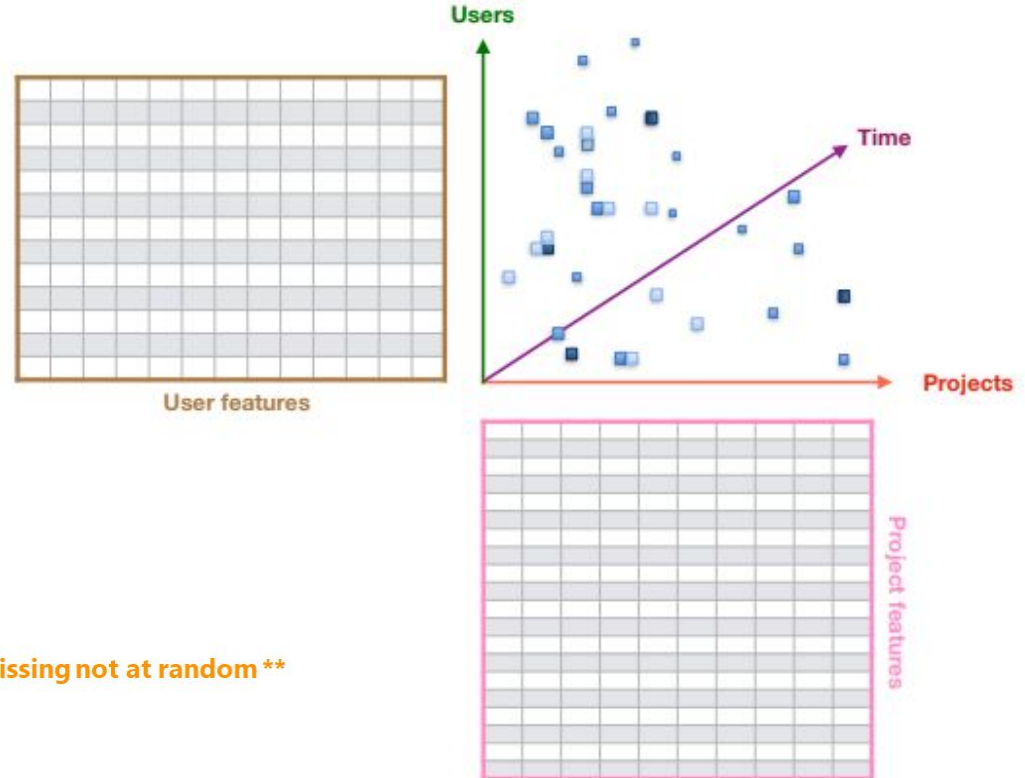
User and info

Interactions between product and user

Rating

Time

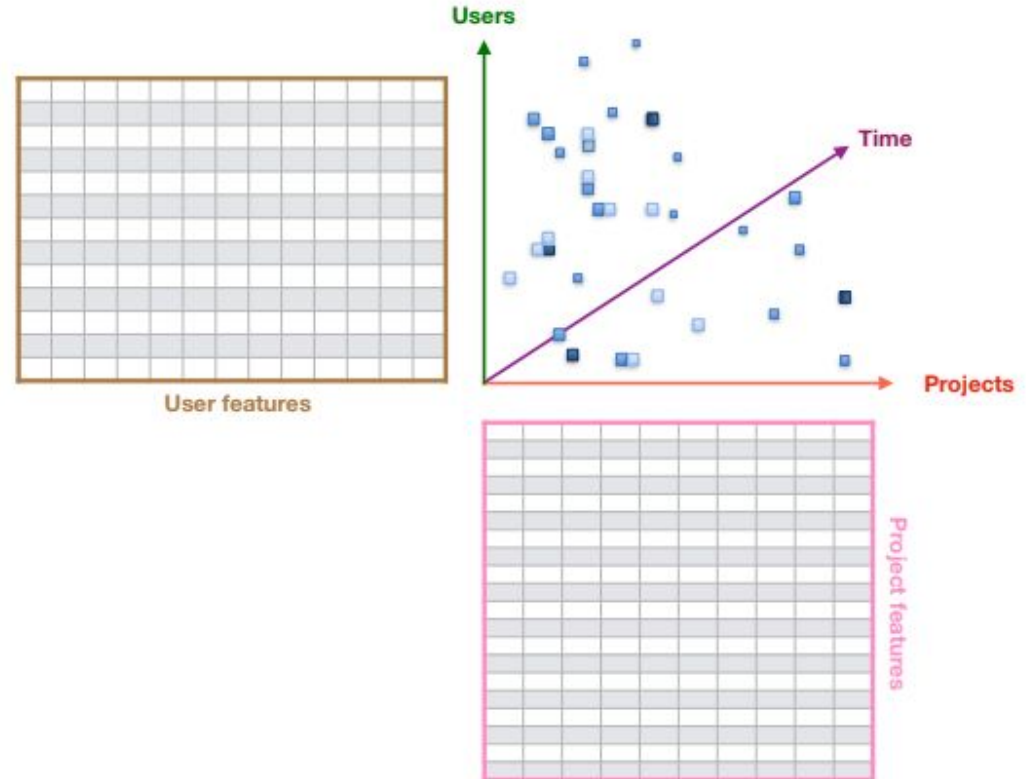
Missing interactions



**** Missing not at random ****

Information in the clicks

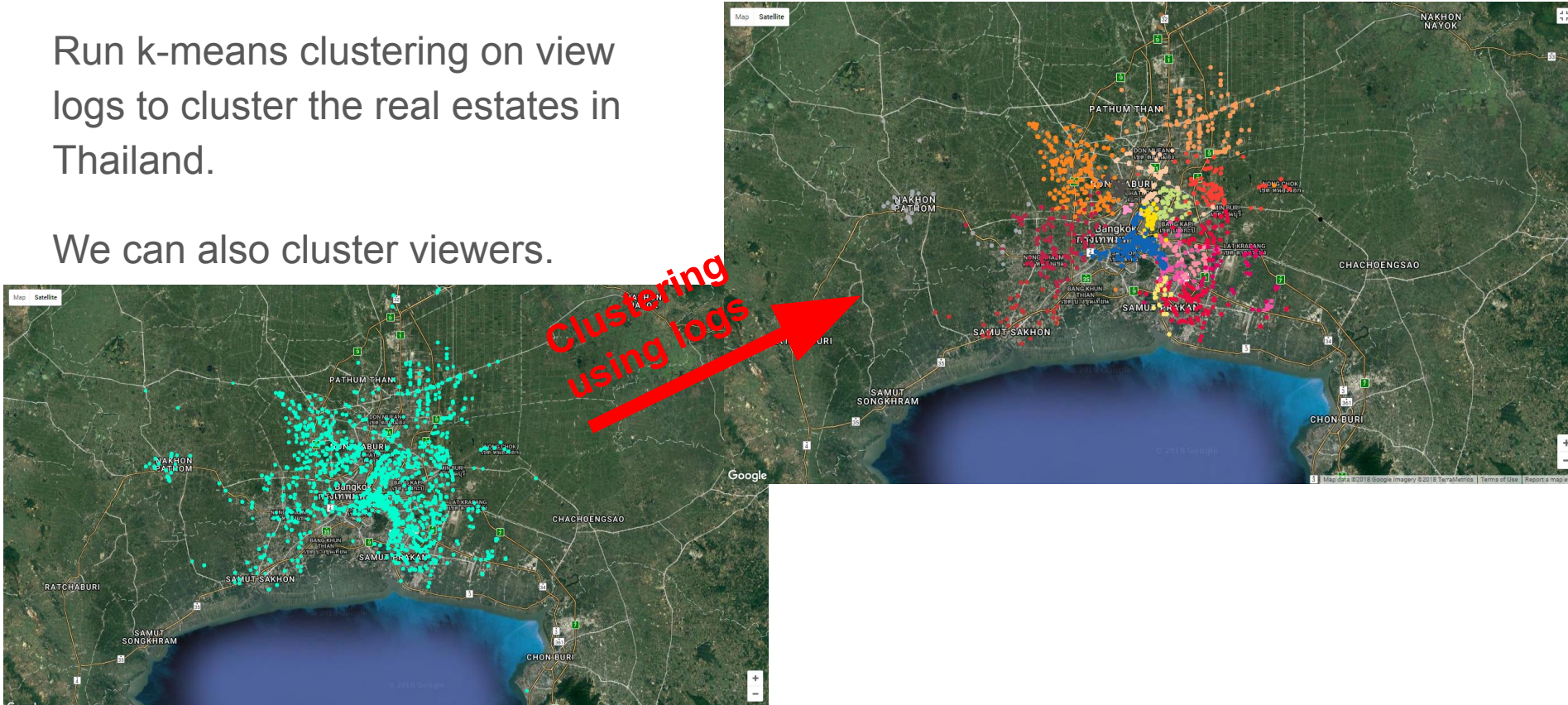
User interactions (views of the projects) can provide interesting insights



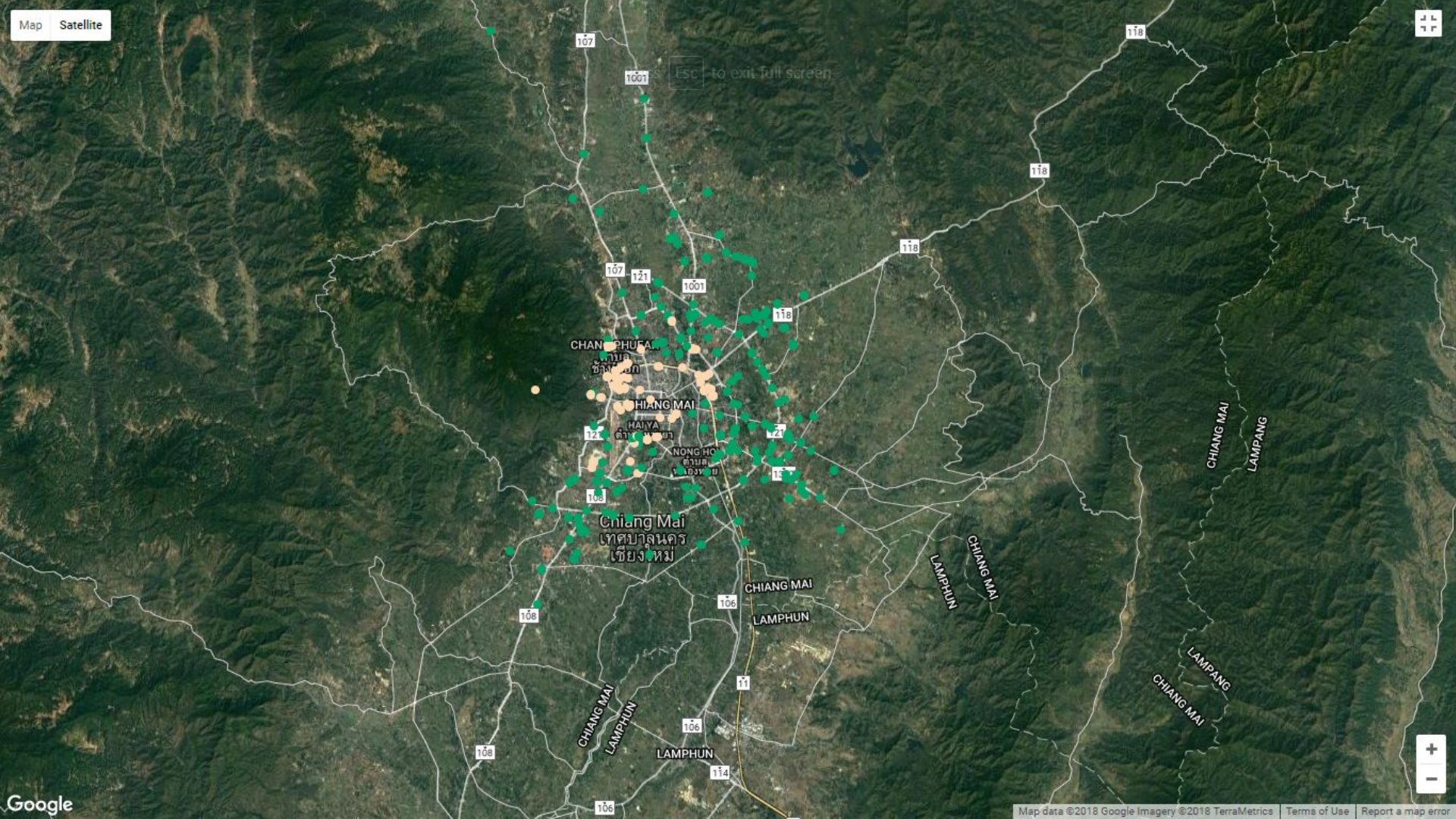
Product segmentation from user interactions

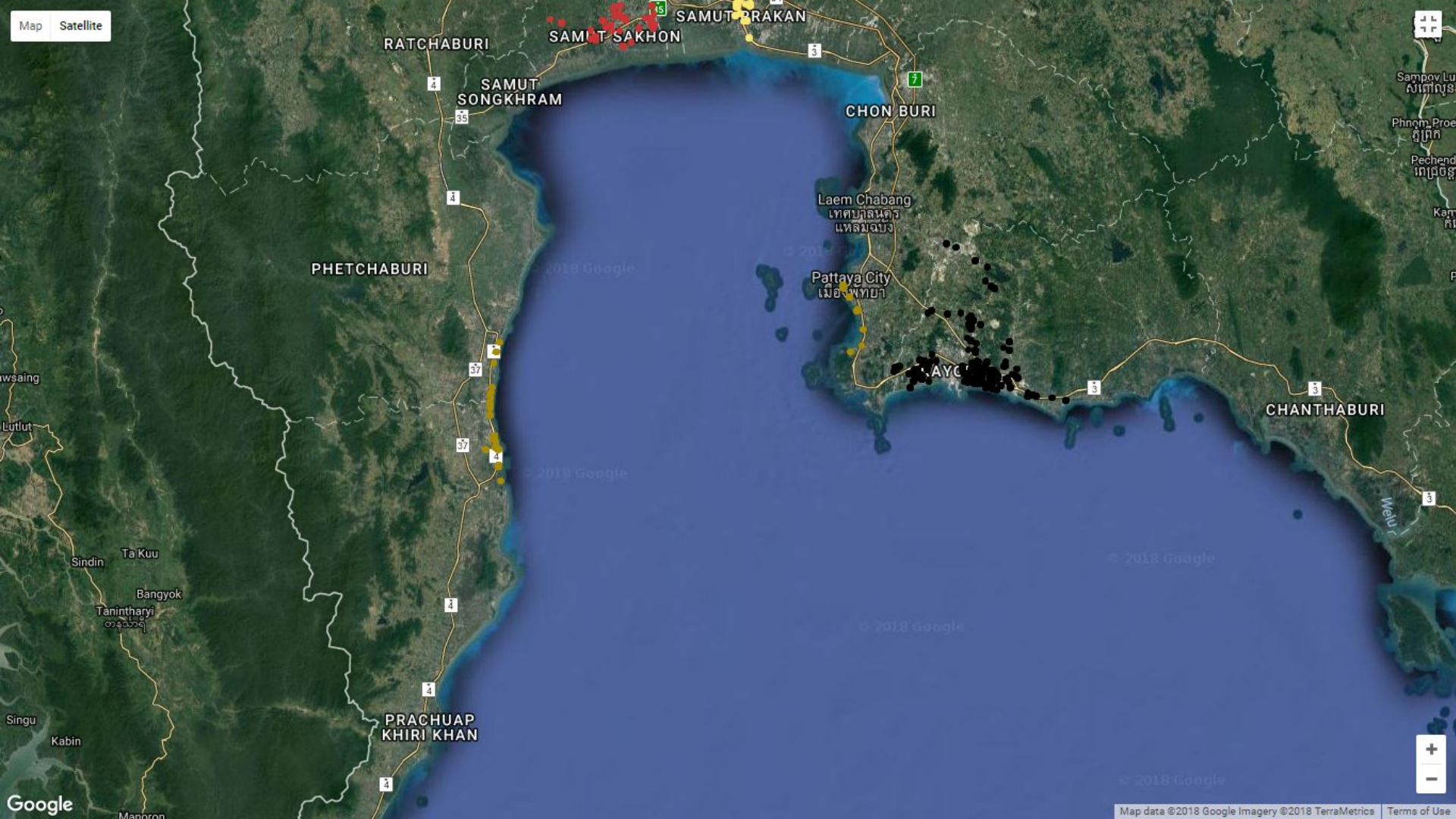
Run k-means clustering on view logs to cluster the real estates in Thailand.

We can also cluster viewers.



Esc to exit full screen





Map Satellite

RATCHABURI

SAMUT SAKHON

SAMUT PRAKAN

SAMUT SONGKHRAM

CHON BURI

PHETCHABURI

Laem Chabang
เทศบาลนคร
แหลมฉบัง

Pattaya City
เมืองพัทยา

PAYC

CHANTHABURI

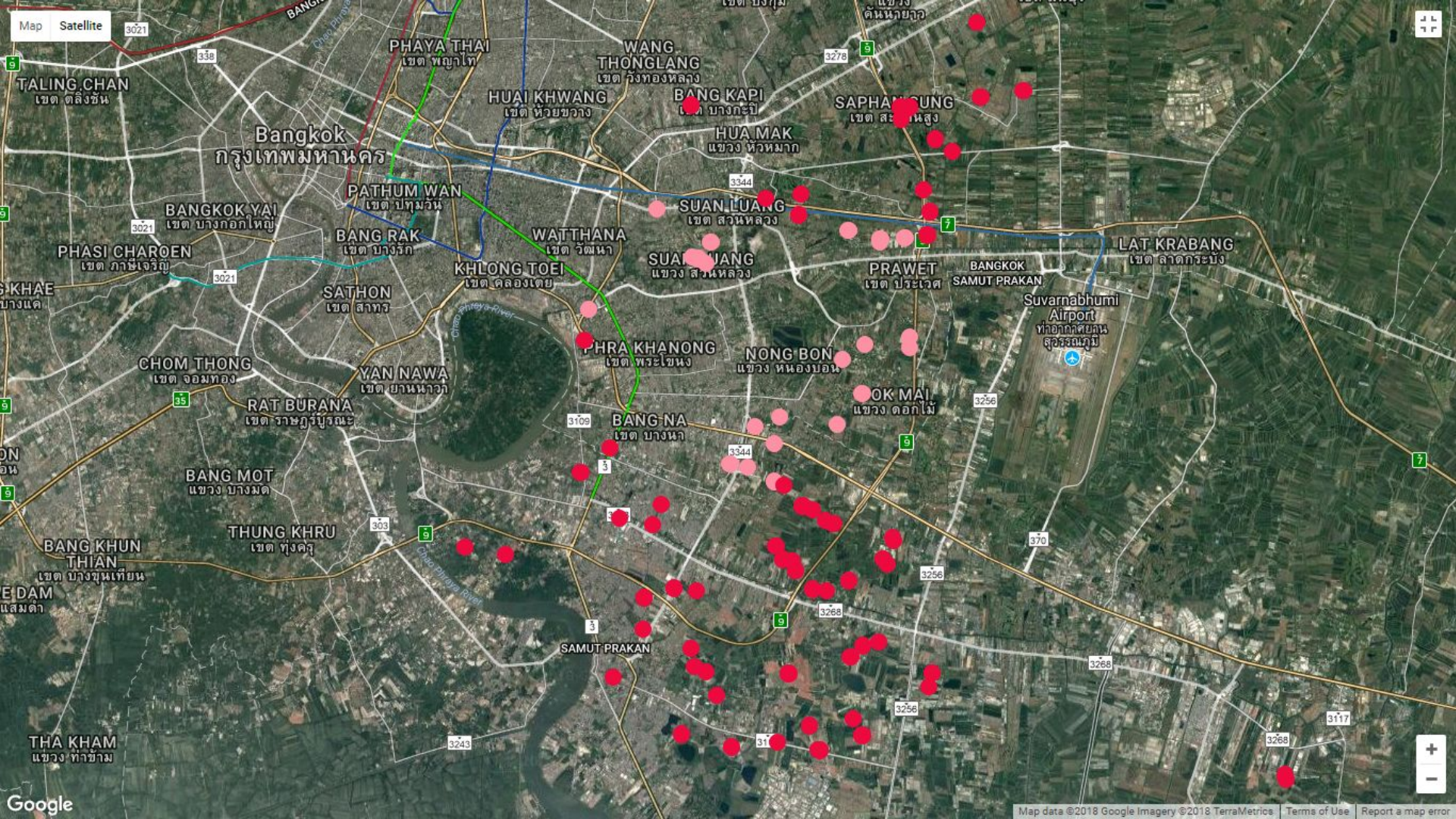
PRACHUAP
KHIRI KHAN

Google

© 2018 Google

Map data ©2018 Google Imagery ©2018 TerraMetrics

Terms of Use



Map Satellite

TALING CHAN
เขต ดลิ่งชัน

Bangkok
กรุงเทพมหานคร

PHAYA THAI
เขต พญาไท

WANG THONGLANG
เขต วังทองหลาง

HUAI KWANG
เขต ห้วยขวาง

BANG KAPI
เขต บางกะปิ

SAPHAN SONG
เขต สะพานสูง

BANGKOK YAI
เขต บางกอกใหญ่

PATHUM WAN
เขต ปทุมวัน

BANG RAK
เขต บางรัก

WATTHANA
เขต วัฒนา

SUAN LUANG
เขต สวนหลวง

SUAN PHONG
เขต สวนพหลวง

PRAWET
เขต ประเวศ

LAT KRABANG
เขต ลาดกระบัง

BANGKHAE
เขต บางแค

PHASI CHAROEN
เขต ภาษีเจริญ

SATHON
เขต สาทร

KHLONG TOEI
เขต คลองเตย

PHRA KHANONG
เขต พระโขนง

NONG BON
เขต นongบอน

Suvarnabhumi
Airport
ท่าอากาศยาน
สุวรรณภูมิ

CHOM THONG
เขต จอมทอง

YAN NAWA
เขต ยานนาวา

RAT BURANA
เขต ราชบุรีบูรณะ

BANG NA
เขต บางนา

OK MAI
เขต ดอกไม้

BANG MOT
เขต บางมด

THUNG KHRU
เขต ทุ่งครุ

BANG KHUN THIAN
เขต บางขุนเทียน

WANG BANG
เขต บางบาง

THA KHAM
เขต ท่าข้าม

SAMUT PRAKAN
เขต สมุทรปราการ

Context information

There's many information available

Product and info

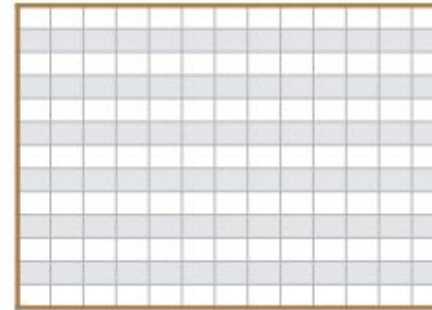
User and info

Interactions between product and user

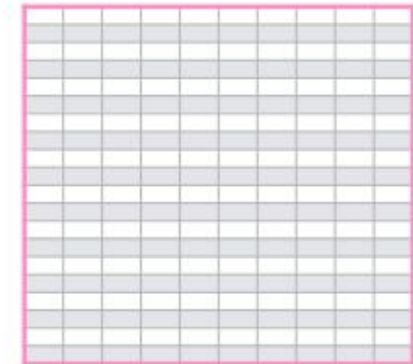
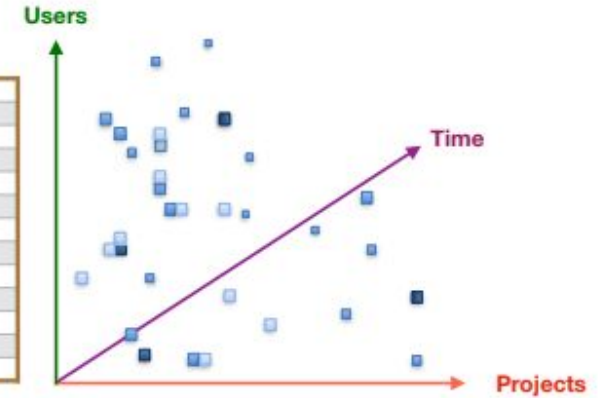
Rating

Time

Missing interactions



User features



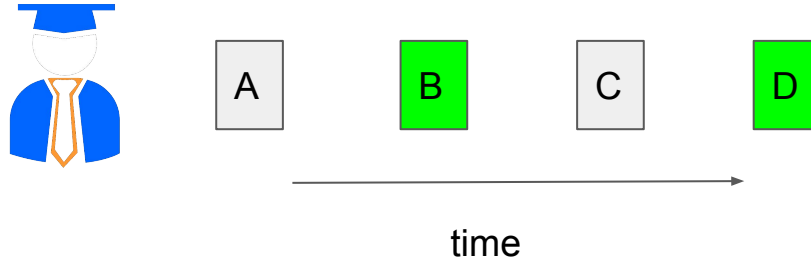
Project features

**** Missing not at random ****

Autoregressive recommendation model

Modeling time information (sequence)

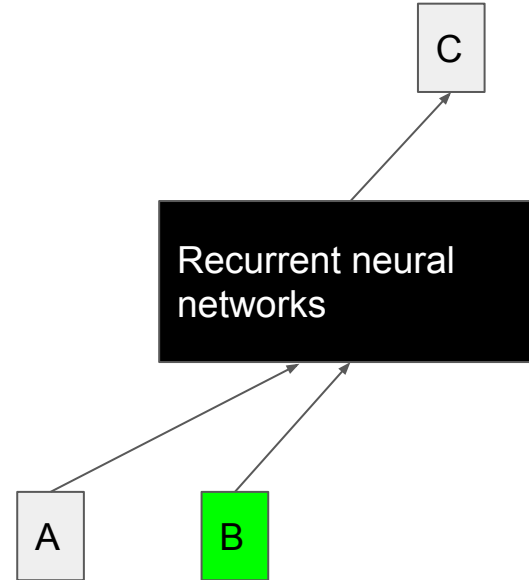
Recurrent Neural Networks



Autoregressive model

Modeling time information (sequence)

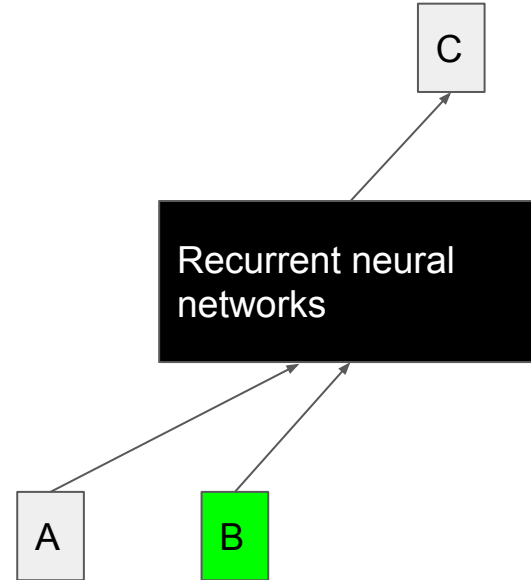
Recurrent Neural Networks



Autoregressive model

Modeling time information (sequence)

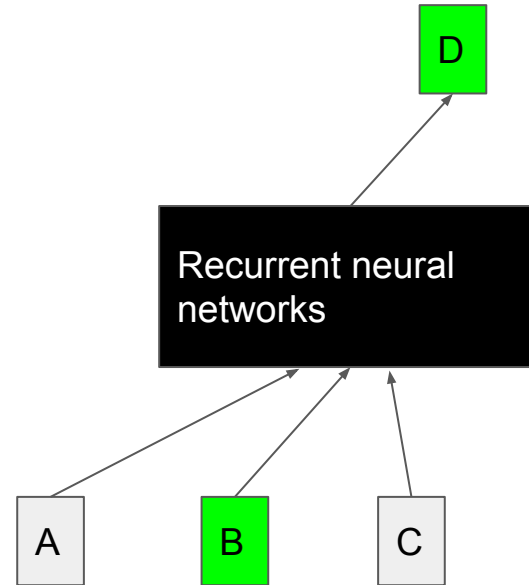
Recurrent Neural Networks



Autoregressive model

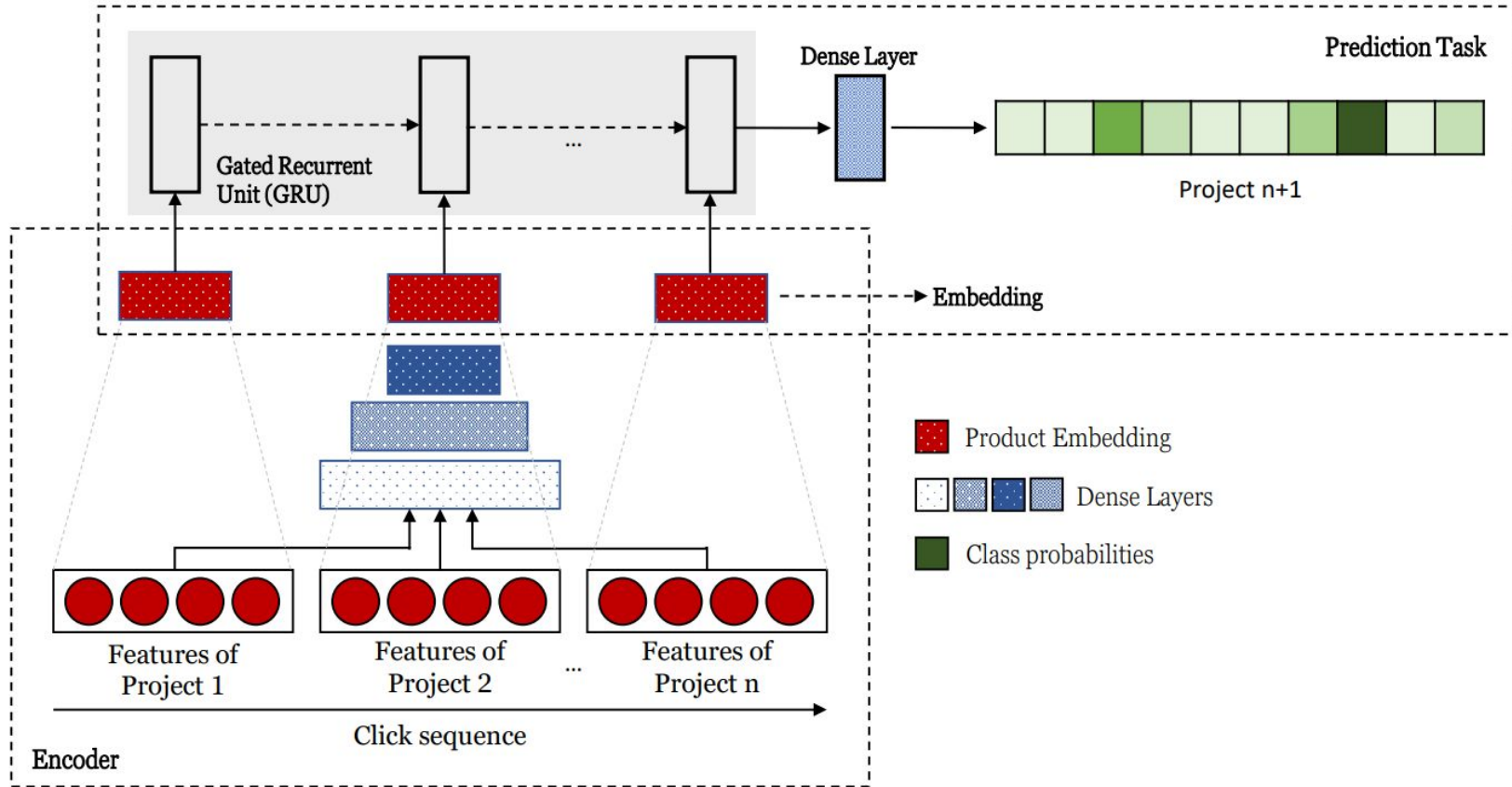
Modeling time information (sequence)

Recurrent Neural Networks



P Covington, Deep Neural Networks for YouTube Recommendations. 2016

A Beutel, Latent Cross: Making Use of Context in Recurrent Recommender Systems, 2018



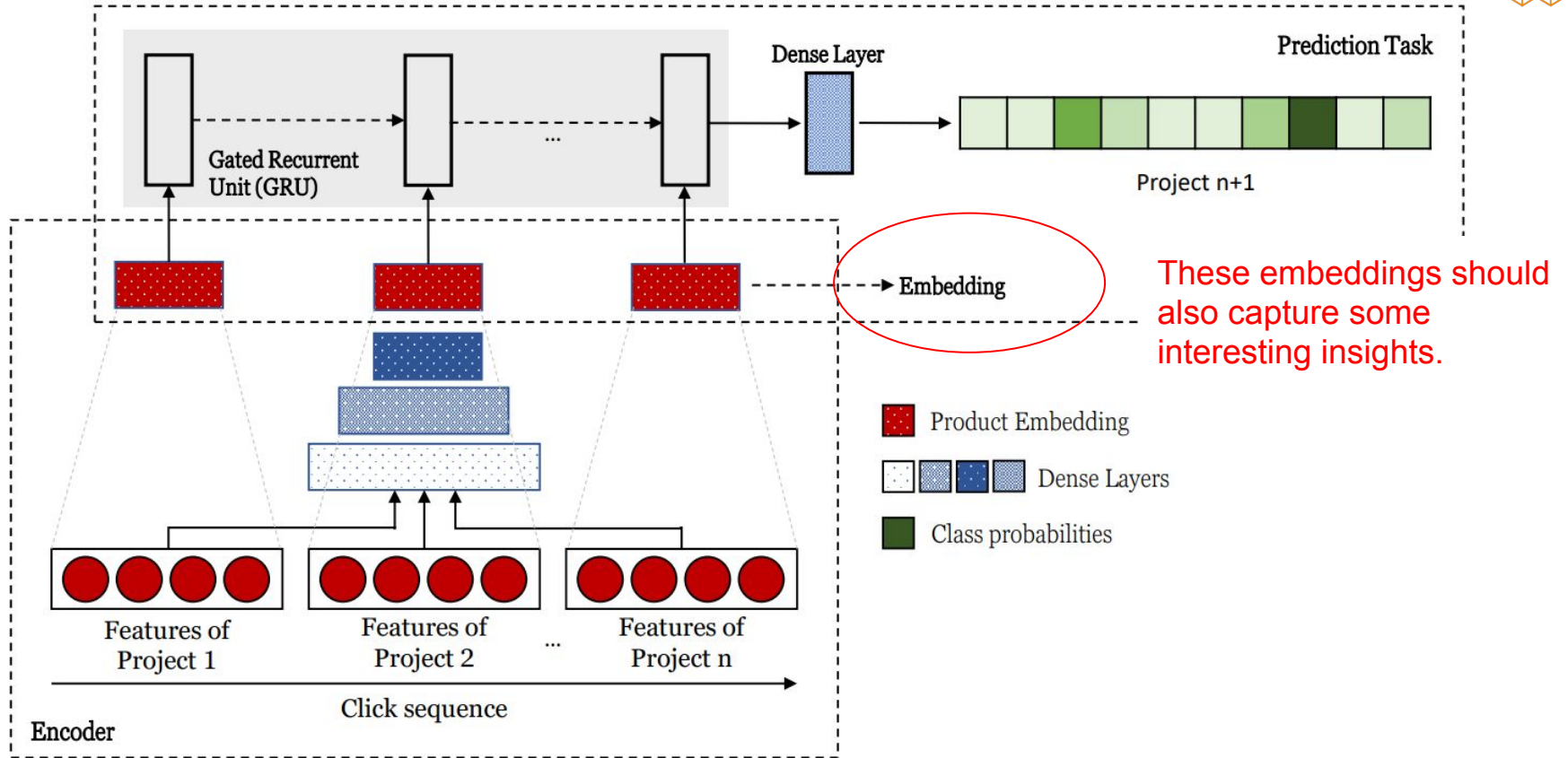
Price: 2,900,000
 Type: Home
 District: Nonthaburi
 Facility: [Security, Park]



Price: 3,900,000
 Type: Home
 District: Bangkok
 Facility: [Fitness, Security, Park]



Price: 5,900,000
 Type: Home
 District: Bangkok
 Facility: [Fitness, Security, Park, Pool]



Price: 2,900,000
 Type: Home
 District: Nonthaburi
 Facility: [Security, Park]



Price: 3,900,000
 Type: Home
 District: Bangkok
 Facility: [Fitness, Security, Park]



Price: 5,900,000
 Type: Home
 District: Bangkok
 Facility: [Fitness, Security, Park, Pool]

HOMEHOP

Home recommender app based on user's lifestyle and commute.

นวัตกรรม
เหนือประสบการณ์
หาบ้านด้วย AI

เวลาเดินทางน้อยลง
เวลาครอบครัวมากขึ้น

HOMEHOP

HOMEHOP Lifestyle

QR code

App Store

Google Play

HOME TECH

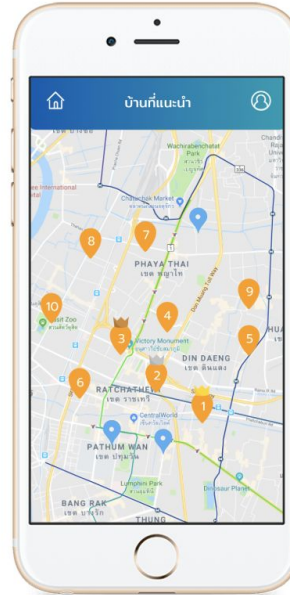
The image shows a family of four (a man, a woman, and two children) playing a game in a bright, modern bedroom. A young boy is blindfolded with a pink cloth and has his arms outstretched. A woman is standing on a bed, and a man is lying on the floor, both with their hands up. A smartphone is overlaid on the left side of the image, displaying the HOMEHOP app interface. The app shows a search bar, a map, and various home listings. A QR code is located in the bottom left corner of the smartphone overlay. The text 'เวลาเดินทางน้อยลง เวลาครอบครัวมากขึ้น' (Less travel time, more family time) is written in large, bold letters across the bottom of the image. The HOMEHOP logo is in the top left corner, and the HOME TECH logo is in the bottom right corner.

Persona

ไลฟ์สไตล์ (persona) ในการเลือกซื้อบ้าน ซึ่งแบ่งประเภท
โดย AI จากข้อมูลผู้ใช้กว่า 10 ล้านคน

Daily life travel

วิธีและเวลาปกติที่คุณเดินทางจากบ้านไปยังที่ทำงาน หรือ
สถานที่ต่างๆ ในชีวิตประจำวันของคุณ



Affordable Price

ช่วงราคาบ้านที่คุณต้องการ หรือสามารถจ่ายได้

Traffic data from iTIC

ข้อมูลการจราจร จากมูลนิธิศูนย์ข้อมูลจราจรอัจฉริยะไทย
เพื่อแนะนำโครงการที่จะใช้เวลาเดินทางน้อยที่สุด

1. เลือกแผนการเดินทาง

ระบุสถานที่ต่างๆ ที่คุณมักจะเดินทางไปในแต่ละวัน เช่น บ้าน โรงเรียนของลูก สถานที่ทำงาน ห้างสรรพสินค้าที่มักเดินทางไปบ่อยๆ เป็นต้น พร้อมทั้งระบุเวลาตั้งแต่ออกจากบ้าน จนถึงเวลากลับถึงบ้าน

2. เลือกวิธีการเดินทาง

เลือกวิธีการเดินทาง เช่น เดินทางโดยรถยนต์ส่วนตัว รถประจำทาง เรือ รถไฟฟ้า โดยระบบจะคำนวณเวลาการเดินทางจากข้อมูลจราจร ของมูลนิธิศูนย์ข้อมูลจราจรอัจฉริยะไทย

3. เลือกช่วงราคาบ้านที่คุณต้องการ

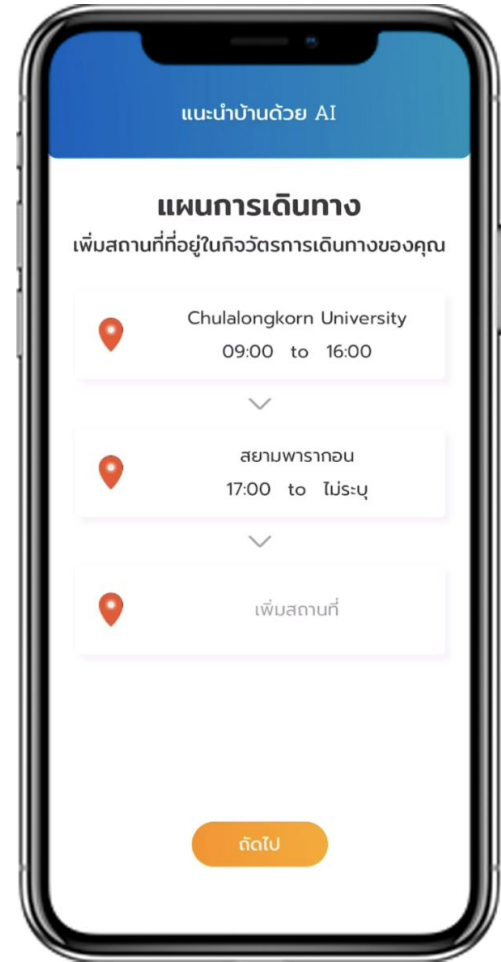
เลือกช่วงราคาบ้านที่คุณต้องการจะซื้อ

4. เลือกเพอร์โซนา

ระบุเพอร์โซนา (persona) หรือไลฟ์สไตล์ของคุณ เช่น เน้นประโยชน์ใช้สอย หรือเน้นความหรูหรา

5. ประมวลผล

ยืนยันข้อมูล แล้วสนุกไปกับการเลือกบ้านที่โปรแกรมแนะนำ ได้ทันที!



Data science for Real Estate

Consumer

~~Matching~~

Social listening

(Real Estate) Developers

Lead generation and smart marketing

Social listening

Project development

Customer segmentation

Trend prediction

Pricing

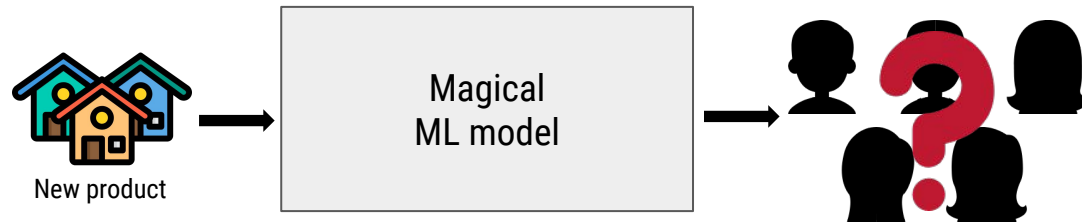
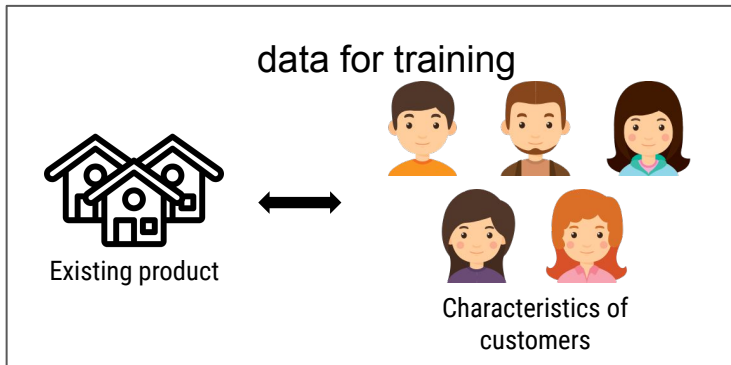


ML for product development

For real estates, no two products are the same. Development based on gut feeling.

Make some informative guess about a new product

- popularity
- the type of potential buyers
- whether to add or remove some features
- the best marketing channel



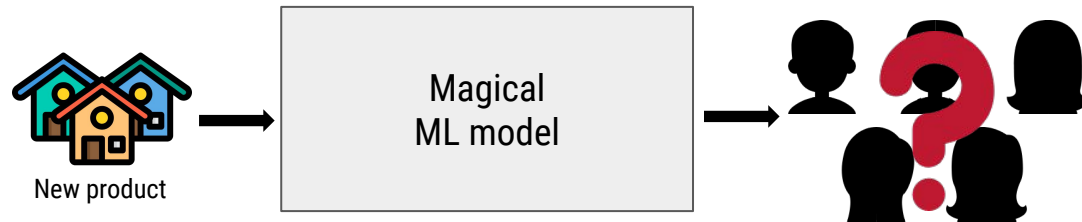
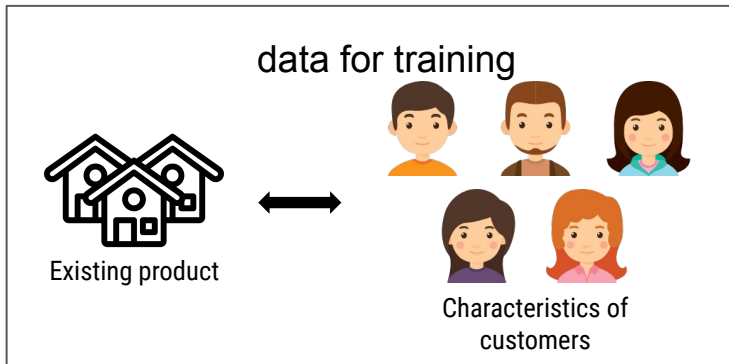
ML for product development

For real estates, no two products are the same. Development based on gut feeling.

Make some informative guess about a new product

- popularity
- the type of potential buyers
- whether to add or remove some features
- the best marketing channel

We want to learn the distribution of the user given some input.
How?



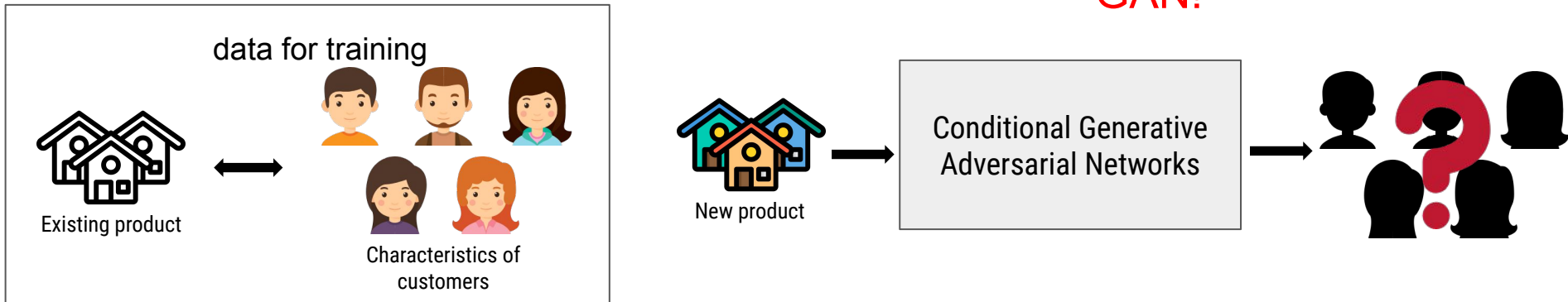
ML for product development

For real estates, no two products are the same. Development based on gut feeling.

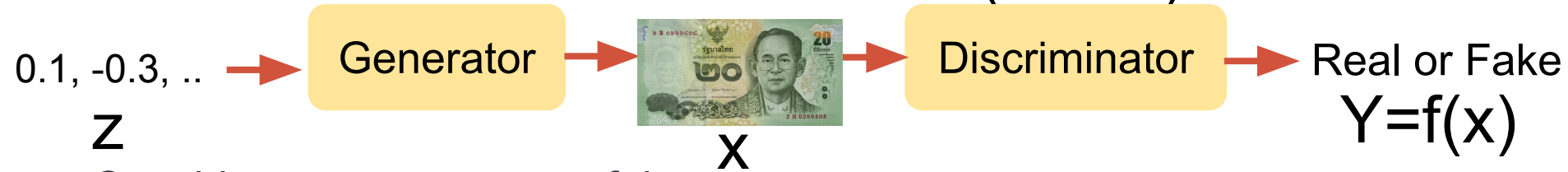
Make some informative guess about a new product

- popularity
- the type of potential buyers
- whether to add or remove some features
- the best marketing channel

GAN!



Generative Adversarial Networks (GANs)



Consider a money counterfeiter

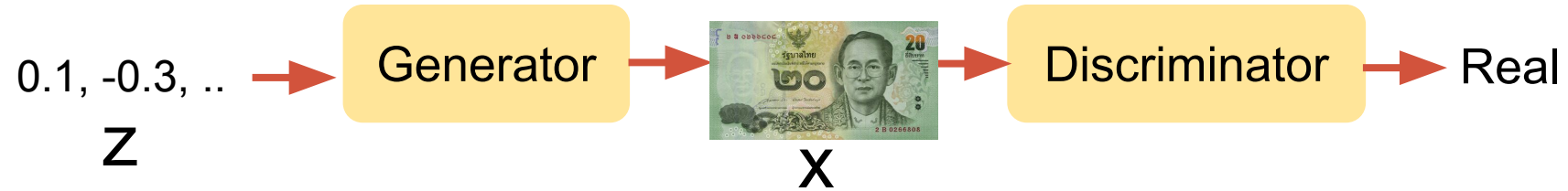
He wants to make fake money that looks real

There's a police that tries to differentiate fake and real money.

The counterfeiter is the **adversary** and is **generating** fake inputs. – Generator network

The police is try to discriminate between fake and real inputs. – Discriminator network

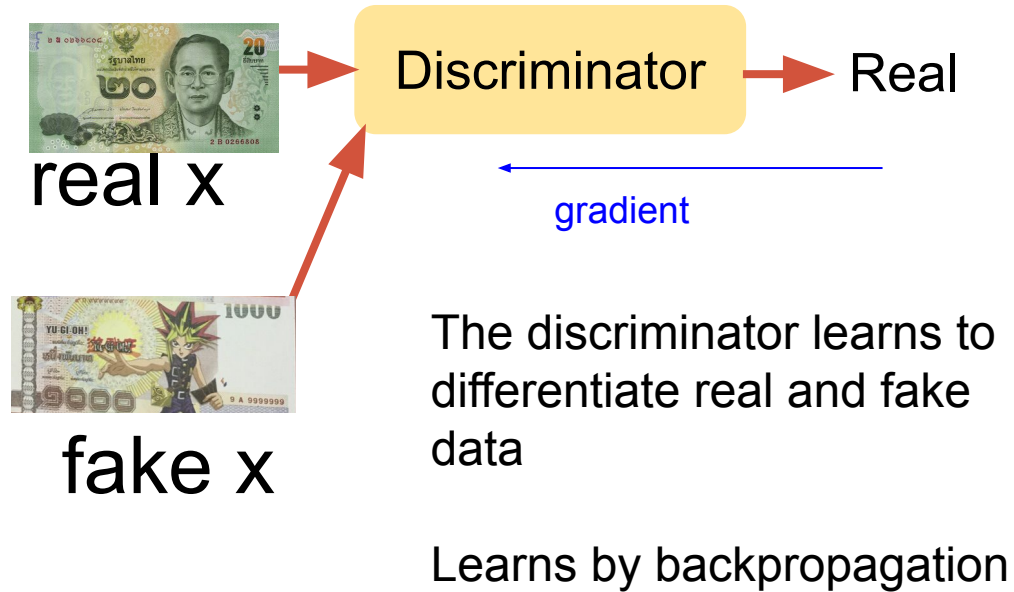
Generative Adversarial Networks (GANs)



Generative Adversarial Networks (GANs)



Generative Adversarial Networks (GANs)

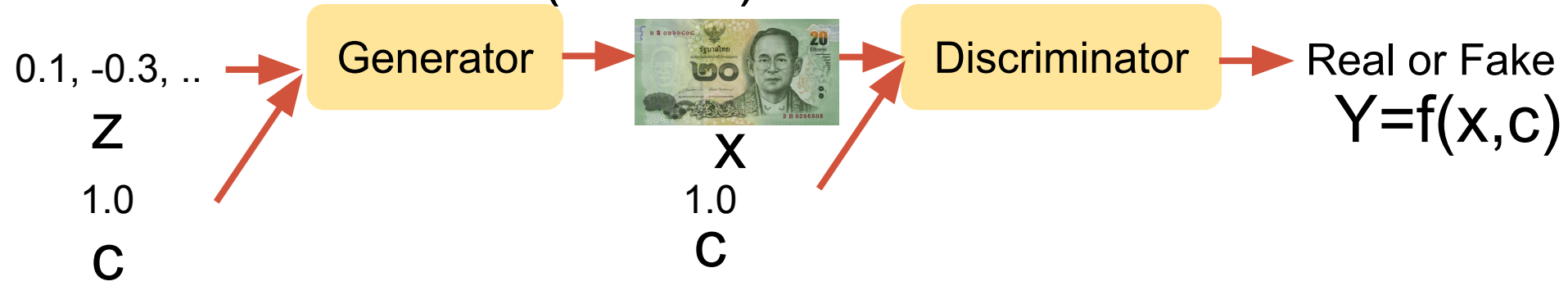


Generative Adversarial Networks (GANs)



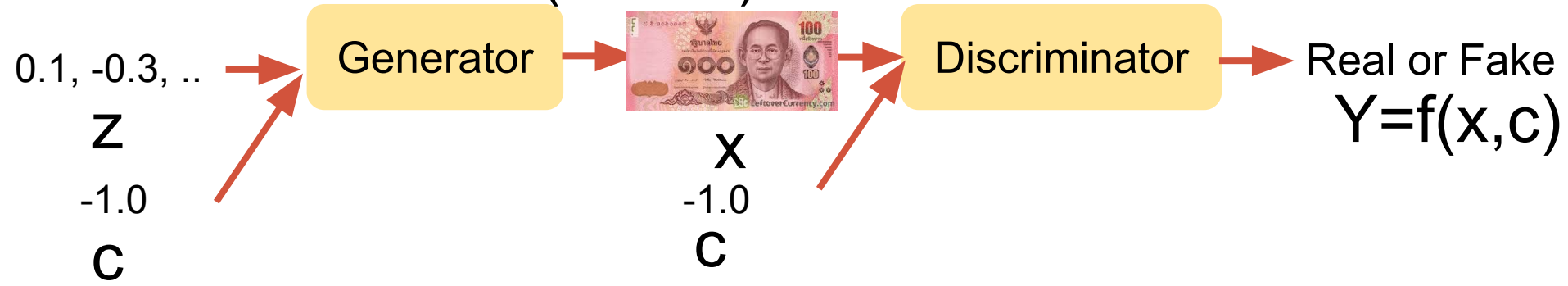
The generator learns to be better by the gradient given by the discriminator

Conditional GAN (CGAN)



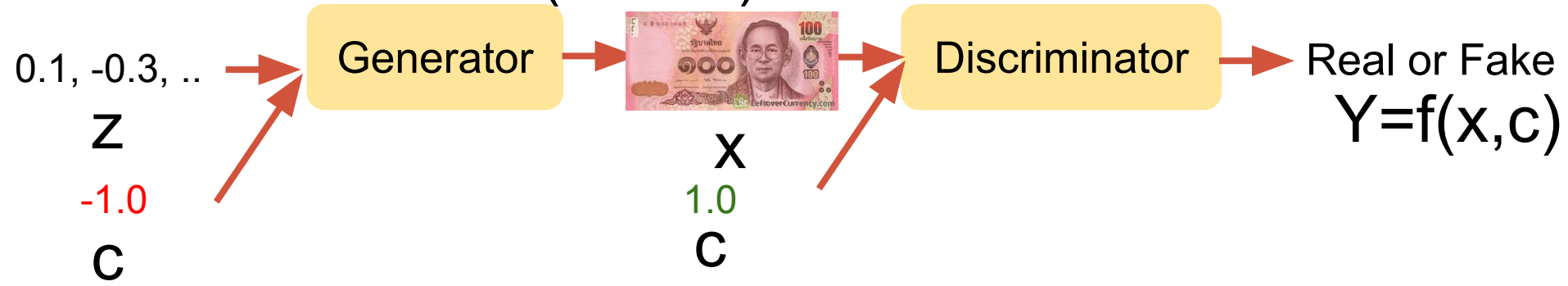
GAN can be conditioned (controlled) to generate things you want by concatenating additional information

Conditional GAN (CGAN)



GAN can be conditioned (controlled) to generate things you want by concatenating additional information

Conditional GAN (CGAN)



GAN can be conditioned (controlled) to generate things you want by concatenating additional information

Example of CGAN applications



This bird is white with some black on its head and wings, and has a long orange beak



This bird has a yellow belly and tarsus, grey back, wings, and brown throat, nape with a black face



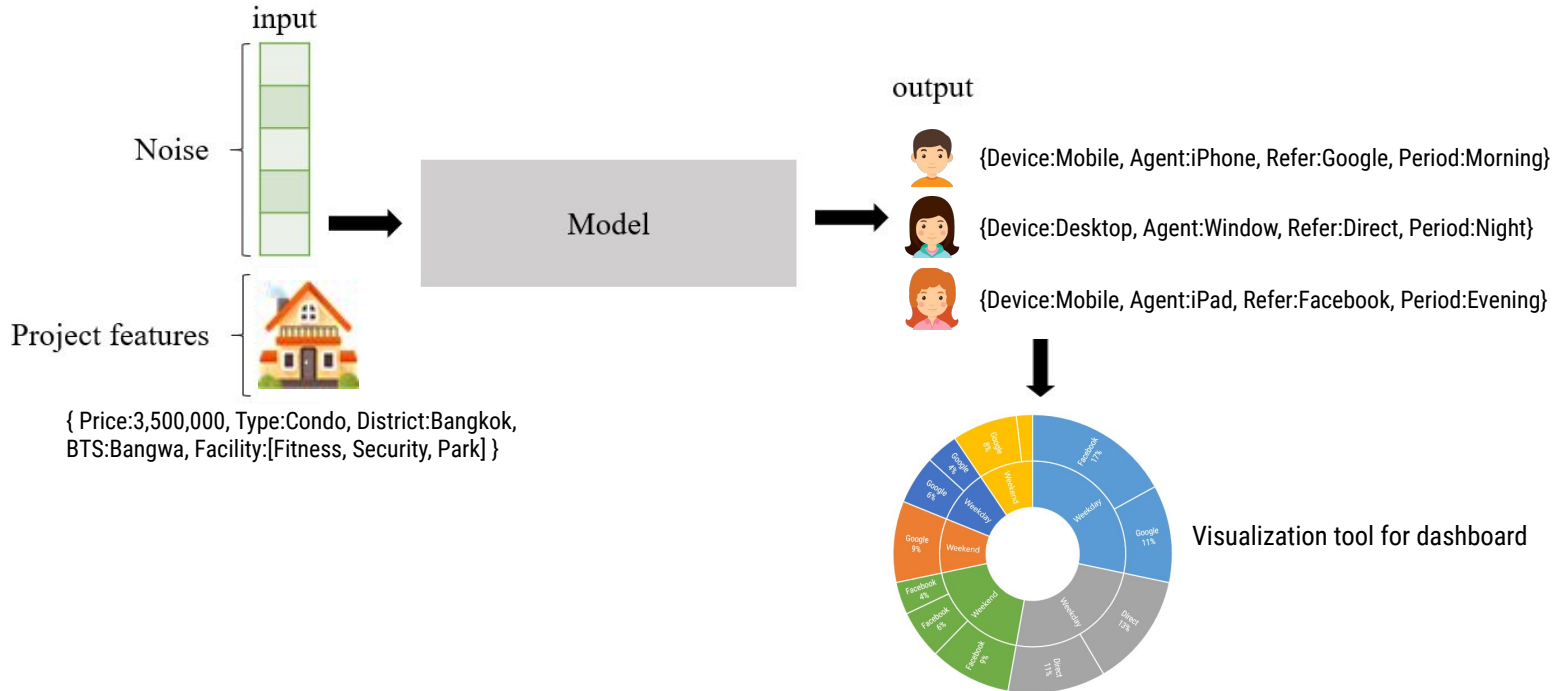
This flower has overlapping pink pointed petals surrounding a ring of short yellow filaments



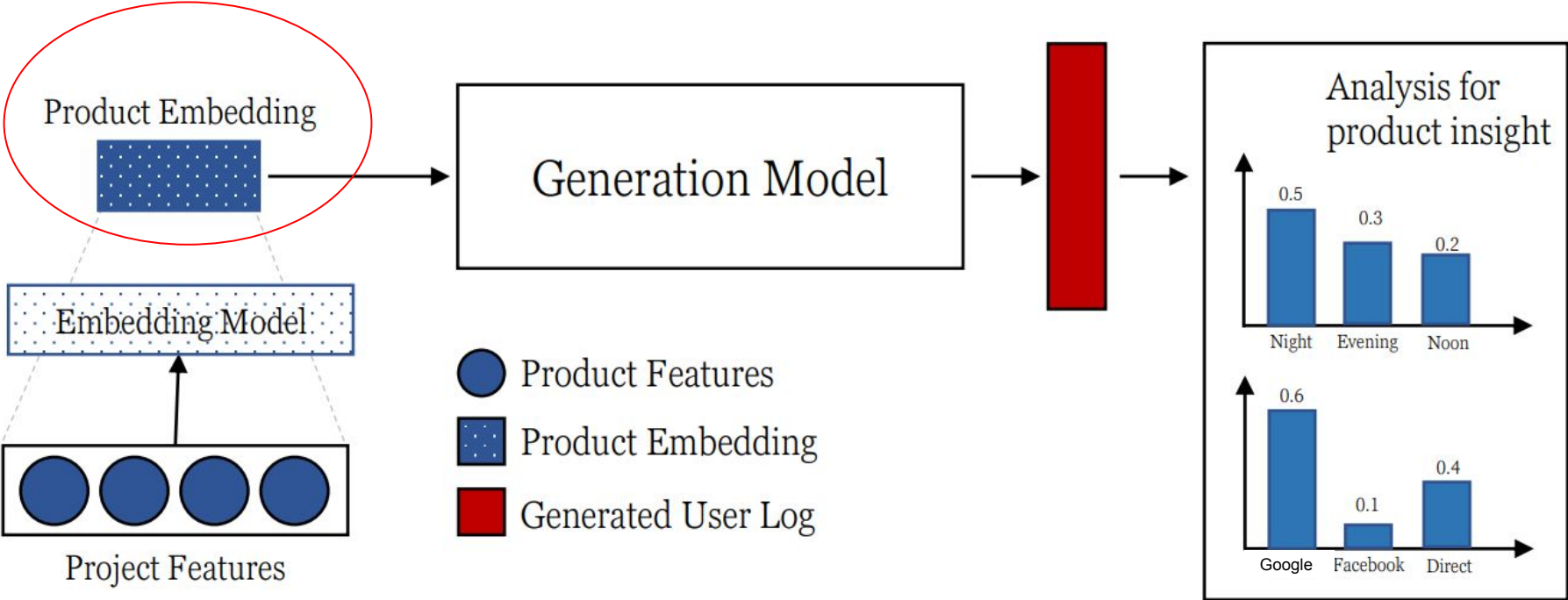
Globally and Locally Consistent Image Completion [Iizuka et al., 2017]

StackGAN: Text to Photo-realistic Image Synthesis with Stacked GANs [Zhang et al. 2017]

Overview of our system



Embedding learned from our recommender system



Why GAN?

vs supervised learning

- supervised learning yields one correct answer (not learning the distribution)
- cannot be used to generate examples

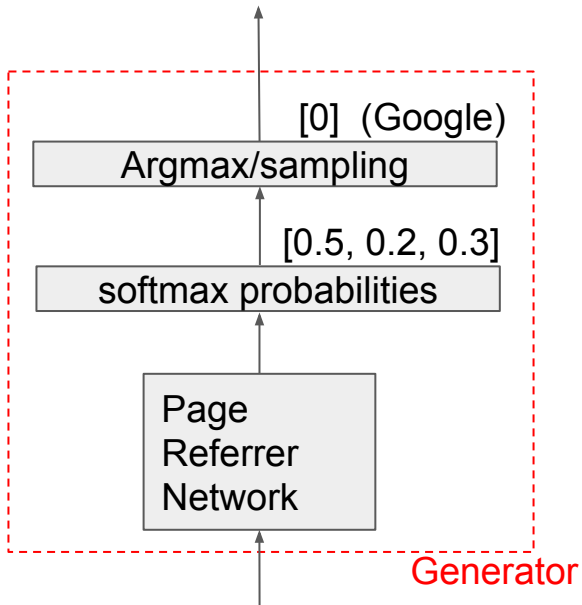
vs other distribution learning methods

- non-parametric
- better than other methods for multi-modal distributions
- generate things that differ from the training data but still “realistic”

GAN for discrete output

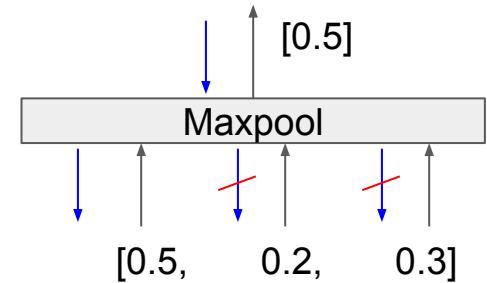
Unlike images, generating discrete output includes a sampling process

fake log for the discriminator



Gradient from discriminator

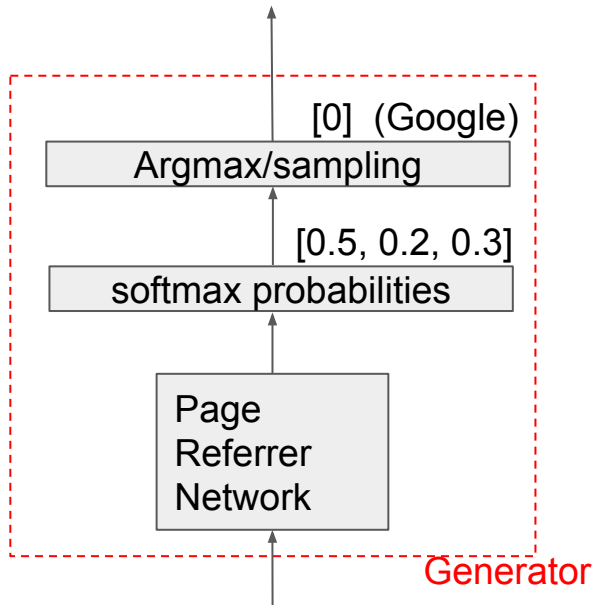
Cannot backprop through the argmax



GAN for discrete output

Unlike images, generating discrete output includes a sampling process

fake log for the discriminator



Gradient from discriminator

Cannot backprop through the argmax

Two popular methods: REINFORCE,
 Gumbel-Softmax approximation
 (<https://arxiv.org/abs/1611.01144>)

BRACE YOURSELVES



MATH IS COMING

A meme featuring a man with dark hair, wearing a light blue button-down shirt, looking slightly to the right with a neutral expression. The background is a blurred office hallway.

**IF YOU DON'T UNDERSTAND, DON'T
WORRY ABOUT IT**

Gumbel Softmax

Sampling from a softmax can be done via the Gumbel-max trick

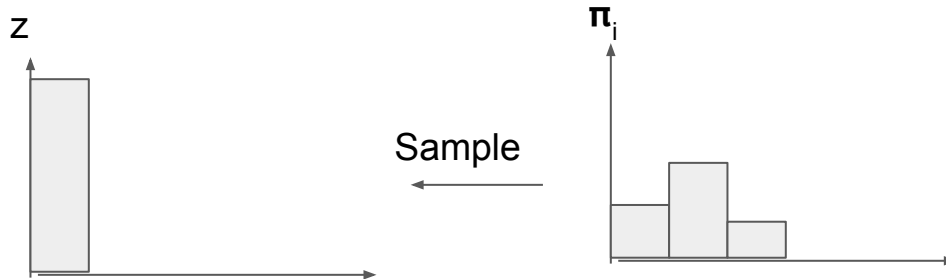
$$z = \text{one_hot} \left(\underset{i}{\text{arg max}} [g_i + \log \pi_i] \right)$$

random value generated from Gumbel dist. (points to g_i)

prob values from softmax (points to π_i)
 Ex. [0.5, 0.2, 0.3]

index for discrete output (points to i)

Final output
 Ex. [1, 0, 0]



Gumbel Softmax

Approximate the argmax term with \mathbf{y} (continuous)

$$z = \text{one_hot} \left(\underset{i}{\text{arg max}} [g_i + \log \pi_i] \right)$$

random value generated from Gumbel dist. (points to g_i)

prob values from softmax (points to $\log \pi_i$)

index for discrete output (points to i)

$$y_i = \frac{\exp((\log(\pi_i) + g_i)/\tau)}{\sum_{j=1}^k \exp((\log(\pi_j) + g_j)/\tau)} \quad \text{for } i = 1, \dots, k.$$

Gumbel Softmax

Approximate the argmax term with \mathbf{y} (continuous)

$$z = \text{one_hot} \left(\underset{i}{\text{arg max}} [g_i + \log \pi_i] \right)$$

random value generated from Gumbel dist.

prob values from softmax

index for discrete output

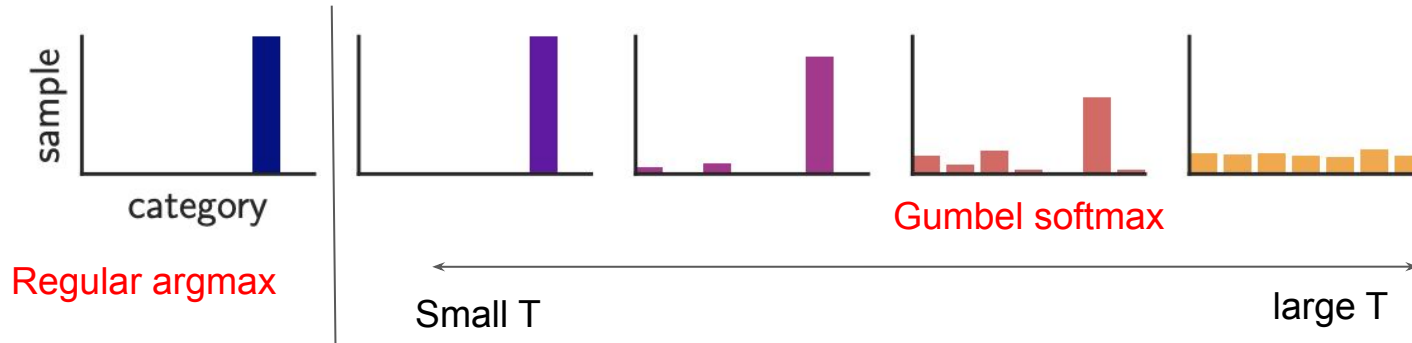
$$y_i = \frac{\exp((\log(\pi_i) + g_i)/\tau)}{\sum_{j=1}^k \exp((\log(\pi_j) + g_j)/\tau)} \quad \text{for } i = 1, \dots, k.$$

Temperature parameter

This rescales the distribution

Gumbel Softmax

y at small T is similar to an argmax but can be backpropagated through



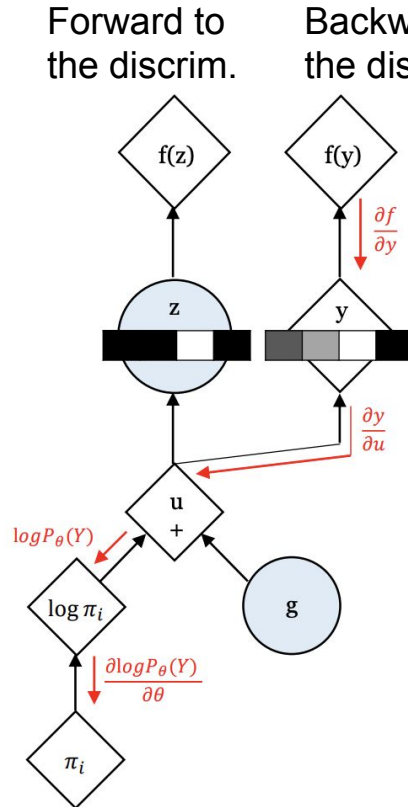
$$y_i = \frac{\exp((\log(\pi_i) + g_i)/\tau)}{\sum_{j=1}^k \exp((\log(\pi_j) + g_j)/\tau)}$$

Temperature parameter

for $i = 1, \dots, k$.

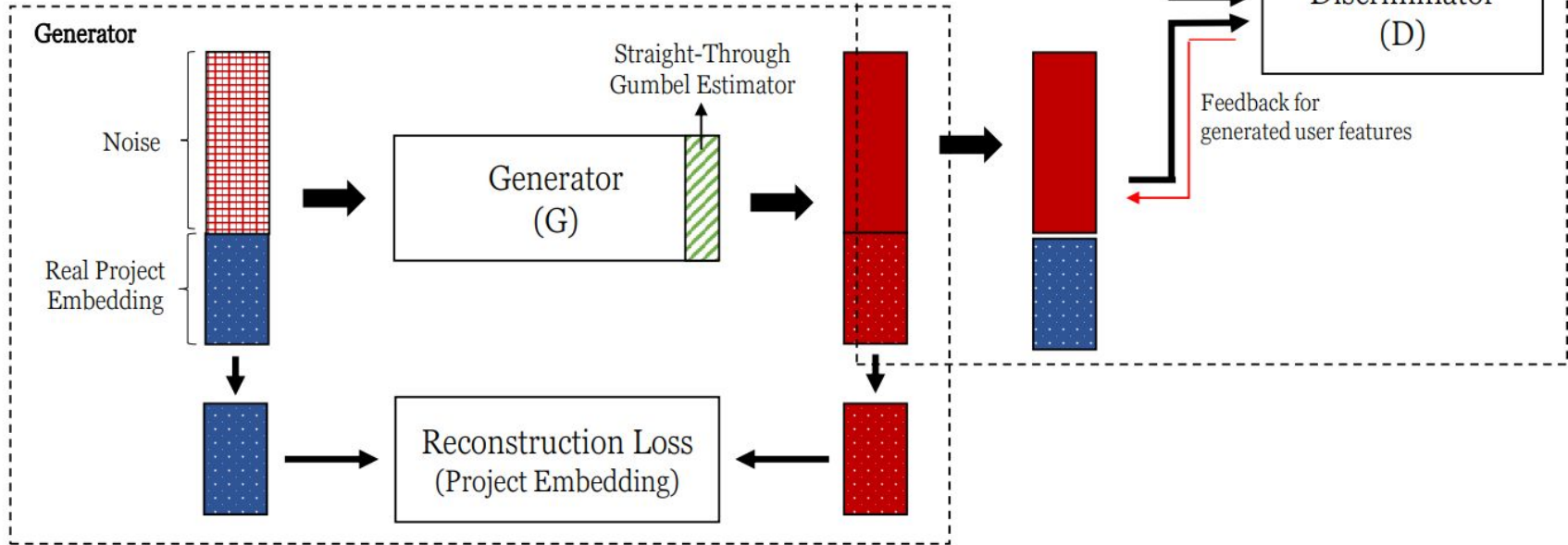
This rescales the distribution

Straight-through Gumbel estimator



The generator generates both the argmax and the Gumbel version. The discriminator uses the argmax version as input. However, the gradient is passed through the Gumbel version.

-  Real User Features
-  Real Product Embedding
-  Noise
-  Generated User Features
-  Generated Product Embedding
-  Straight-Through Gumbel Estimator

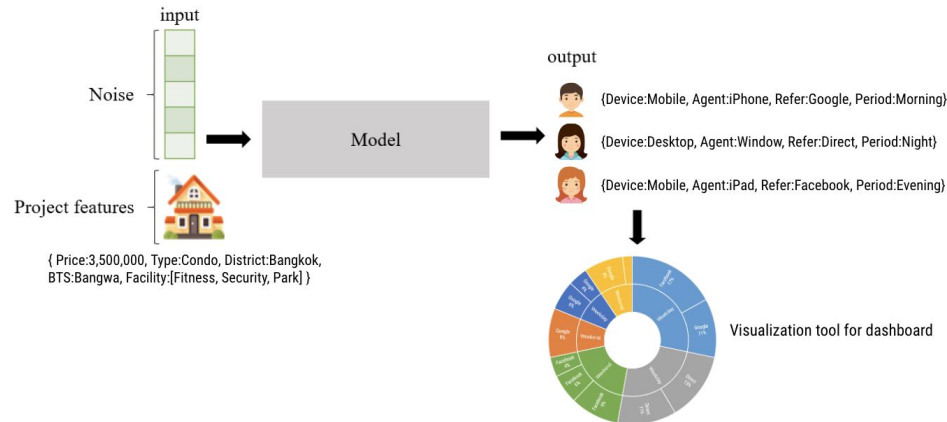


Experimental setup

~5000 projects, ~2 million log entries

- Held out 50 random projects as novel projects to generate
- Measure the distribution of generated logs vs real data

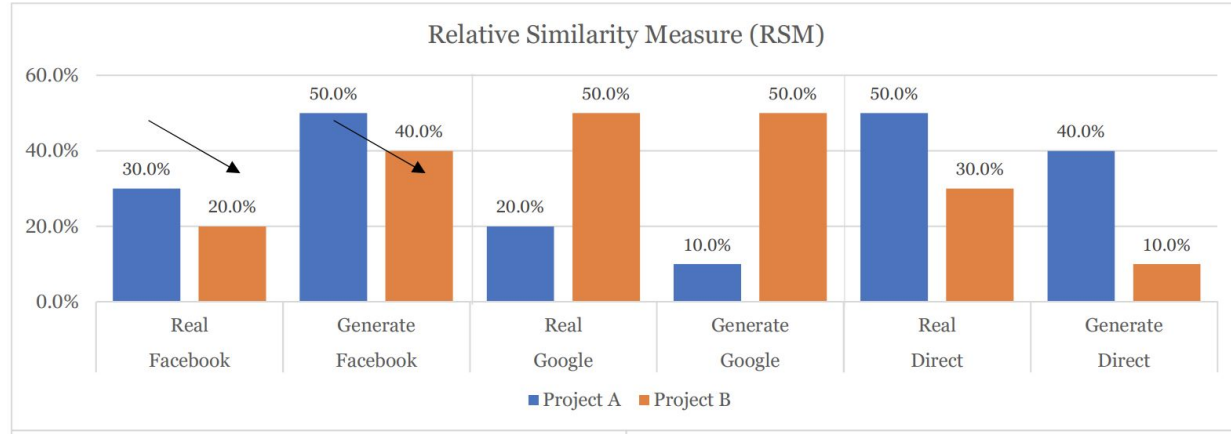
Average the performance over 10 runs



Metrics

RSM

Relative measure
Across project pairs



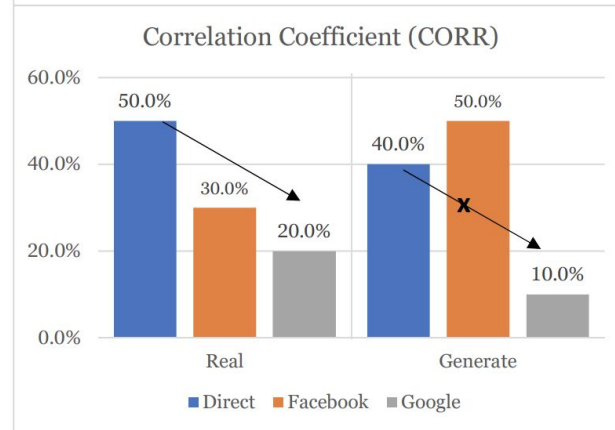
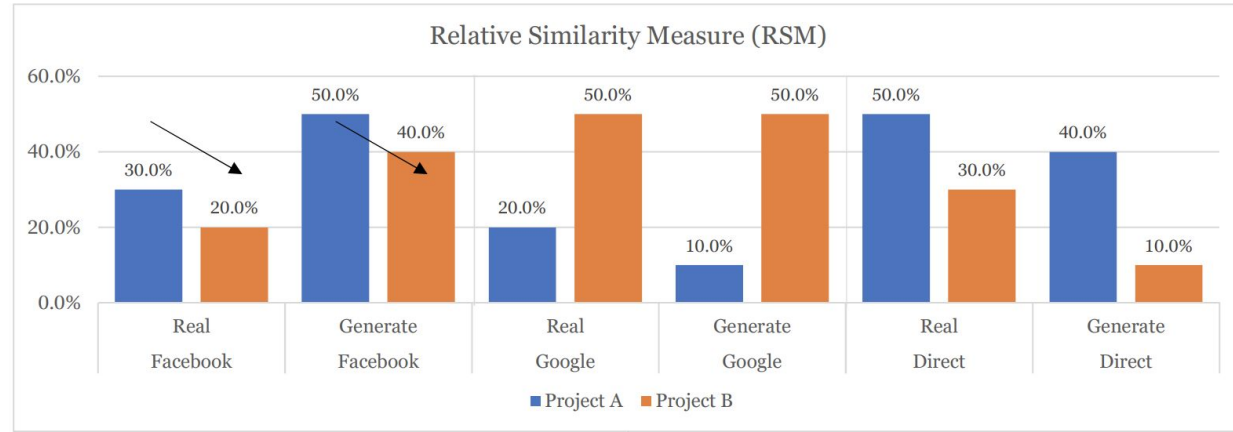
Metrics

RSM

Relative measure
 Across project pairs

Correlation

Relative measure
 Within a project



Metrics

RSM

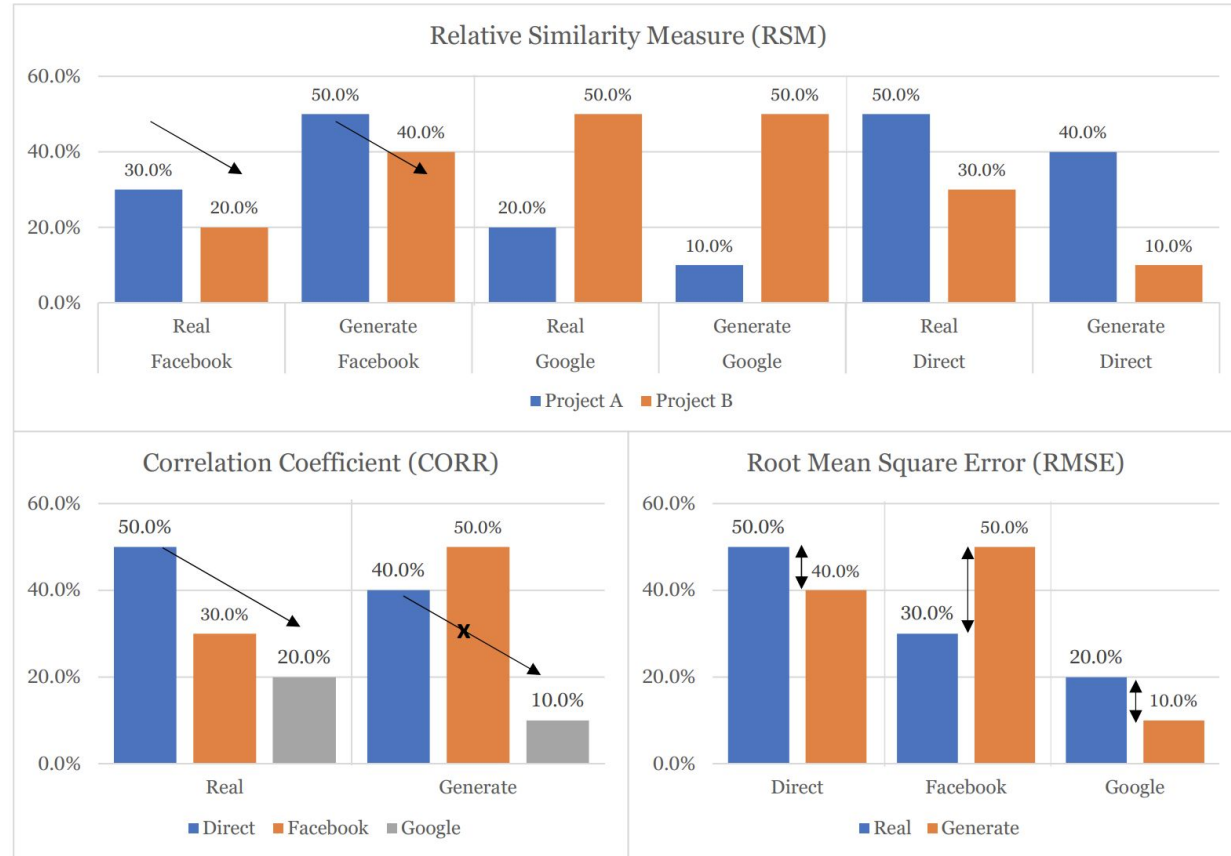
Relative measure
Across project pairs

Correlation

Relative measure
Within a project

RMSE

Absolute measure



Results

| Model | RSM | CORR | RMSE |
|-------------------|-------|-------|-------|
| GAN with Rec. Emb | 72.5% | 88.9% | 16.2% |
| | | | |
| NN with Rec. Emb | 54.7% | 71.6% | 28.0% |

Our model with recommender embedding

Use the most similar project in the training data based on recommendation embeddings

Results

| Model | RSM | CORR | RMSE |
|---------------------------|-------|-------|-------|
| GAN with Rec. Emb | 72.5% | 88.9% | 16.2% |
| GAN with AutoEncoder Emb | 69.7% | 87.8% | 18.1% |
| GAN with product features | 67.9% | 86.6% | 18.2% |
| | | | |
| NN with Rec. Emb | 54.7% | 71.6% | 28.0% |

Our model with recommender embedding

Our model with embeddings learned from Autoencoder
 Our model with product features instead of embedding

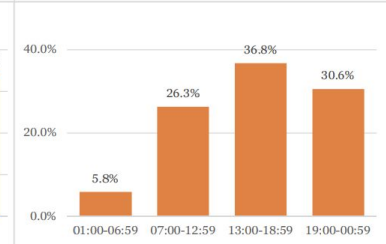
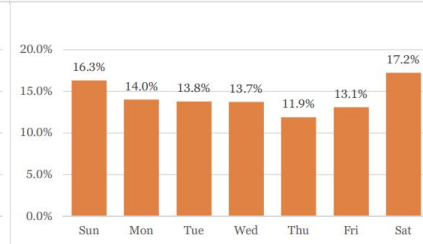
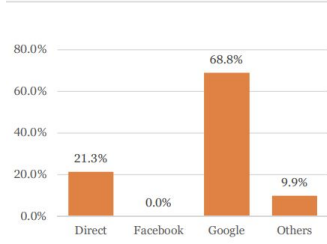
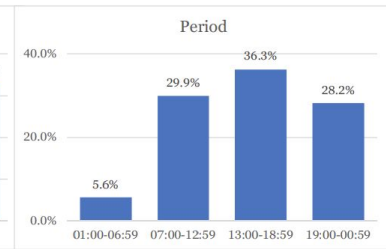
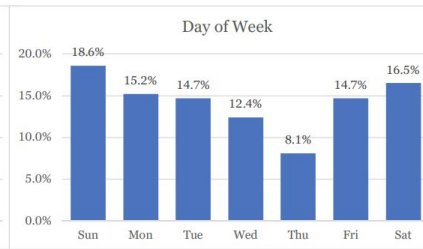
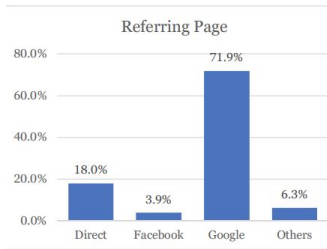
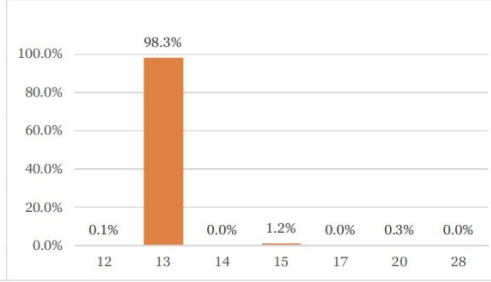
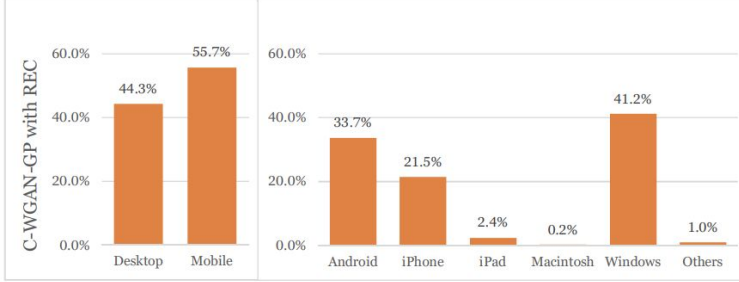
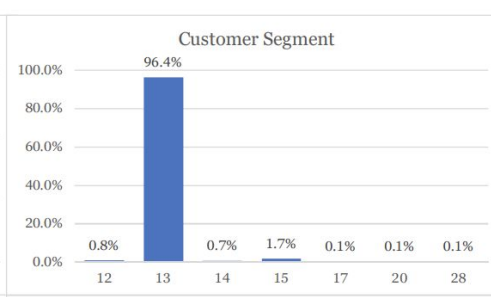
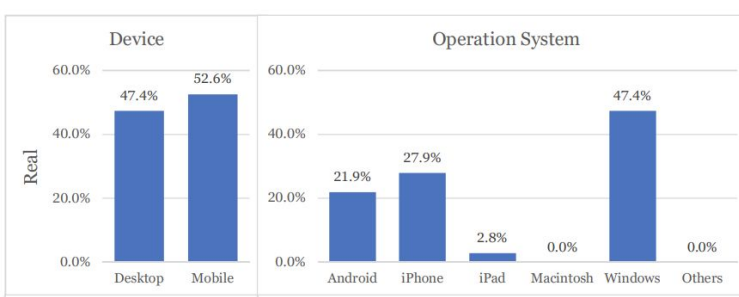
No knowledge about relationships between different products

Results

| Model | RSM | CORR | RMSE |
|---------------------------|-------|-------|-------|
| GAN with Rec. Emb | 72.5% | 88.9% | 16.2% |
| GAN with AutoEncoder Emb | 69.7% | 87.8% | 18.1% |
| GAN with product features | 67.9% | 86.6% | 18.2% |
| VAE with Rec. Emb | 65.3% | 85.6% | 20.3% |
| NN with Rec. Emb | 54.7% | 71.6% | 28.0% |

Our model with recommender embedding

Instead of GAN use VAE



Data science for Real Estate

Consumer

Matching

Autoregressive Recommender system

(Real Estate) Developers

Project development

GAN-based distribution learning



Team



Questions?

