# pCTR Challenge

Kenny He and Chris Yu

# Background

- Search advertising has been one of the major revenue sources of the Internet industry for years. A key technology behind search advertising is to predict the click-through rate (pCTR) of ads, which drives the pricing model.
- This year we will be looking at the dataset of image ads provided by Tencent Multimedia

GOAL:

- Given training dataset of online advertising system results, contestants must accurately predict the test data

# Dataset

The training data is published as a text file, where each line is a training instance derived from session log messages. (23,906,738)

| UserID | AdvID | AdID | CreativeID | Impression | Click |
|--------|-------|------|------------|------------|-------|

Evaluation DS (3,253,943)    Validation DS (3,236,631)


Additional Files:

titles.txt (34,163 lines)

| AdID | Title |
|------|-------|

images.zip: folder of all the image ads provided (24,125 jpg files)

users.txt

| UserID | Gender | Age |
|--------|--------|-----|

(23,439,495 lines, 491MB in size)

# Approach

- Classification of Titles
- Sort User Data by UID
- Split Training Data by Gender and Age
- Analyze Image Characteristics
- Build the model

# First Step: Classification of Ad Titles

- Put the different keywords from titles.txt into categories (for us it was 7 categories)

  - Men's wear, Women's wear, Gaming, Beauty, Household, Education, etc.
- Scan through each title, and keep weight of each Category that the keywords are from
- The category with the highest keywords weight will be the label for that title

# 2nd Step: Sort Users Data by UID

- Read records from Users.txt, and build a binary search tree of the user information
- The nodes in the tree are stored in an array
- A hash function to generate a number to represent a combination of age and gender:
  GA(age, gender) = 31 | (age*2 + gender) | 130

  age <=15      15<age<65      age>=65

# 3rd Step: Split Training Data

- We are sure that age and gender are the two most important factors related with the interest to certain ads
- Split the huge training dataset by the age and gender of the users into 100 groups
- Take advantage of the BST and hash function created in step 2 for higher performance
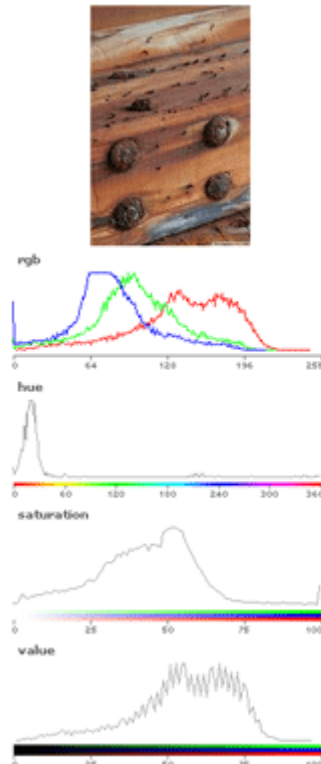
# 4th Step: Analyze Image Features

Using a Web-based API: [http://mkweb.bcgsc.ca/color_summarizer/](http://mkweb.bcgsc.ca/color_summarizer/)

- The api is capable of analyzing image's color palette and color statistics
- We will mine the api for image features to be added to the prediction model
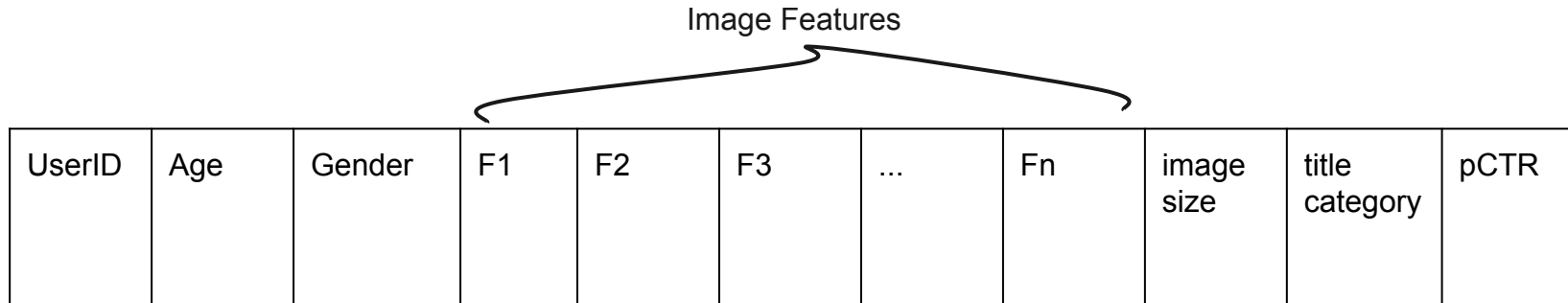
# Examples

# 5th Step: Build the model

- Building the model to predict pCTR for each ad for gender/age groups:
- Apply logistic regression for CTR predictions

Image Features

| UserID | Age | Gender | F1 | F2 | F3 | ... | Fn | image size | title category | pCTR |
|--------|-----|--------|----|----|----|-----|----|------------|----------------|------|
|        |     |        |    |    |    |     |    |            |                |      |

# Potential Trends

Gender-based ads vs Color trends (pink/red for women, black/blue for men)

Gaming-based ads vs saturation and contrast (high contrast, full colors)

CTR rate vs Text Area (what % of image area is text or non-text)

# Difficulties

- Huge amount of user data and training data
- No tool or conventional method for data splitting
- Need to do title classification, which requires lots of work to manually identify keywords and categories
- Time consuming on image data analysis
- Our goal is to build a model to predict the pCTR for different age/gender groups, rather than a model for predicting the individual clicks -- avoid overfitting problem

# Summary & Thanks!

- Classification / Data Split / Logistic Regression

- Q & A

- Contact us if you have good ideas about the image analysis