

# Introduction to kaggle

*[where you get to solve “real-world” problems]*

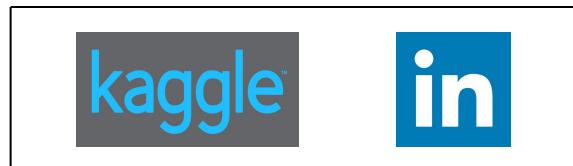
| Date : 15th June , 2019 | Venue : Grofers , Gurugram | Speaker : Akarsh Raj |



# WHO AM I?

---

- **Akarsh Raj**
  - Data Scientist @ [ZS Associates](#)
  - Former Data Scientist @ [Xceedance](#)
  - B.Tech | CSE | University of Petroleum & Energy Studies
  - Kaggle Competitions Expert - Top 2 percent
    - 5 Competition medals – 1 Silver, 4 Bronze (Still waiting for my first gold)
    - 12 Discussion medals – 2 Gold, 10 Bronze
- Hobbies – Kaggle, Cryptography, Hacking, Soccer, Poker
- Where can you find me? (Click on the thumbnails below)



# POLL

How much do you know about Kaggle?

- **Category 1** : I'm just curious about it
- **Category 2** : I registered to have a look at data, kernels and discussions
- **Category 3** : I have joined a real world data competition
- **Category 4** : It has already turned to an addiction



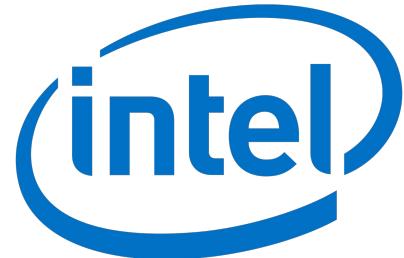
# Why Kaggle?

- Leading platform for ML related competitions, discussions since 2010
- Top companies post their real world problems which can be solved using ML, DL or some magic maybe.
- Google acquired Kaggle in 2017
- One can find implementations of almost all the algorithms, in most creative ways possible (Refer to the [awesome kernel](#) by Shivam Bansal who tried to collate implementations of various algorithms)
- Gets to meet the best data scientists from all over the world



## FEW TOP SPONSORS

---



Quora



# WHAT ALL YOU CAN DO AT KAGGLE ?



COMPETE

(Competition Tier)



SHARE

(Kernels Tier)



DISCUSS

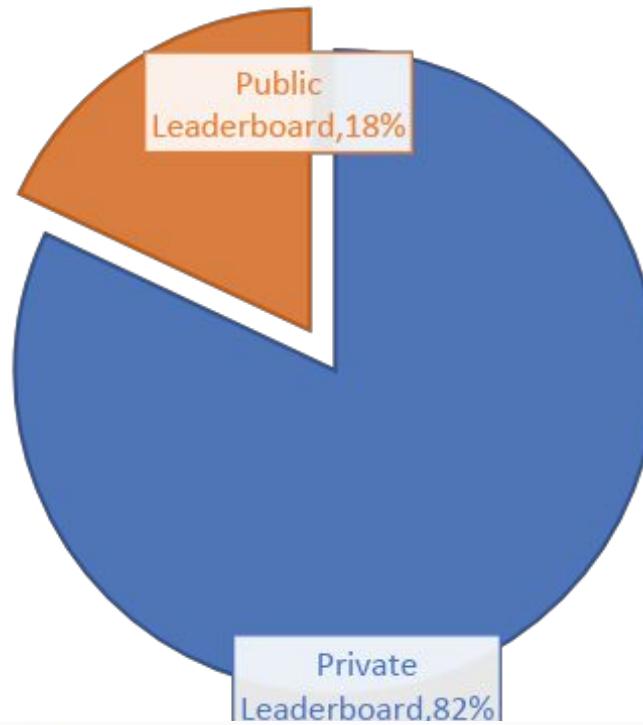
(Discussions Tier)

# COMPETITION

Public Leaderboard – **Test Data**

Private leaderboard – **Unseen Data**

## Submission file in a competition



No submission of code/model, until one is in the money zone.

[Public Leaderboard](#)

[Private Leaderboard](#)

This leaderboard is calculated with approximately 18% of the test data.

The final results will be based on the other 82%, so the final standings may be different.

[Raw Data](#)

[Refresh](#)



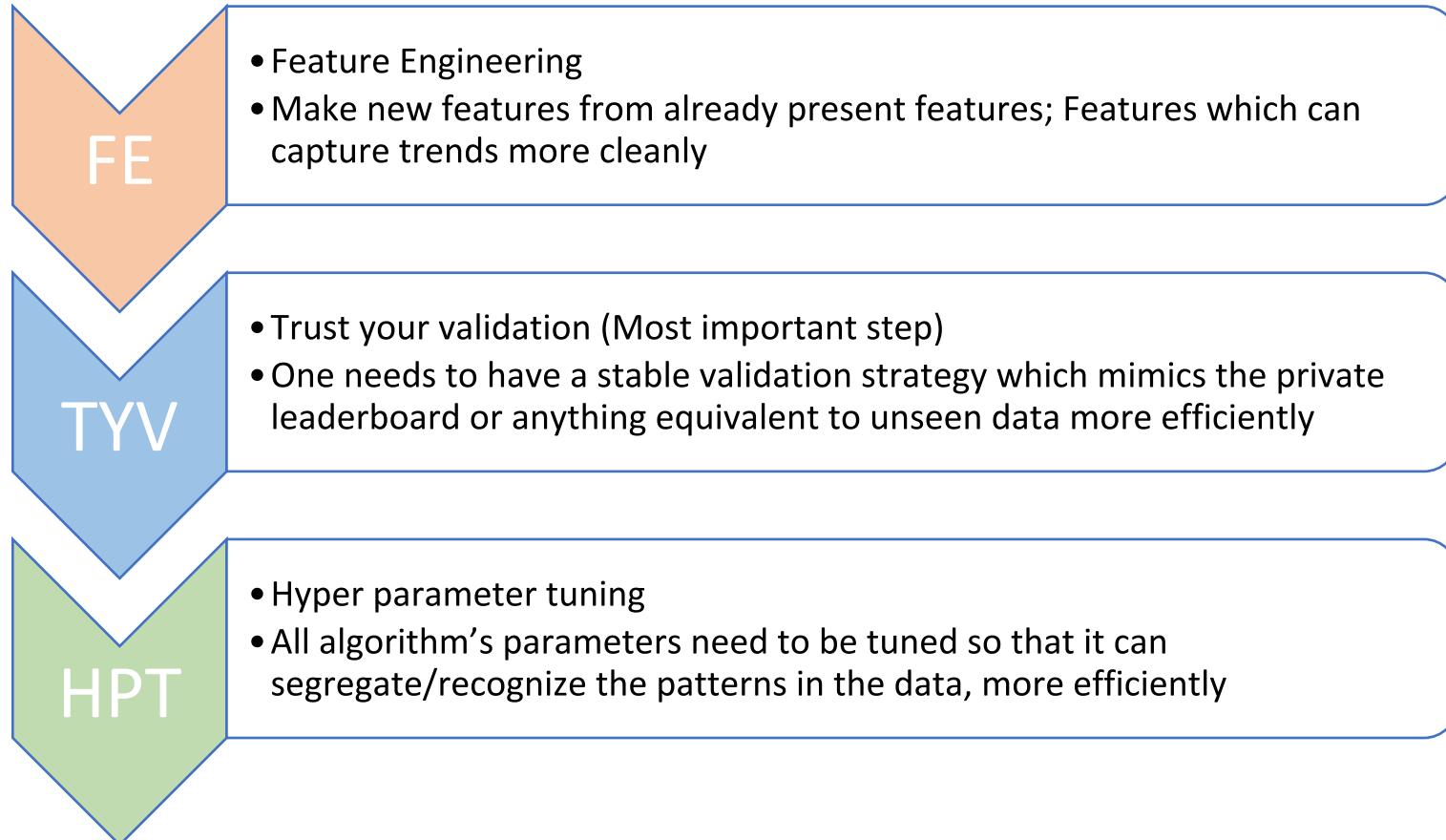
Google

LOGICO

kaggle

CODING  
NINJAS  
www.codingninjas.in

# MOST IMPORTANT THINGS TO BE DONE IN ANY COMPETITION



## HOW I APPROACH A KAGGLE COMPETITION ?

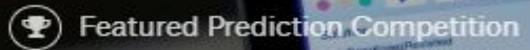


Google

LOGICO

kaggle

CODING  
NINJAS  
www.codingninjas.in



## TalkingData AdTracking Fraud Detection Challenge

Can you detect fraudulent click traffic for mobile app ads?

\$25,000

Prize Money

T TalkingData · 3,951 teams · a year ago

kaggle  
DAYS  
MEETUP

### Data Description

### Data fields

Each row of the training data contains a click record, with the following features.

- ip : ip address of click.
- app : app id for marketing.
- device : device type id of user mobile phone (e.g., iphone 6 plus, iphone 7, huawei mate 7, etc.)
- os : os version id of user mobile phone
- channel : channel id of mobile ad publisher
- click\_time : timestamp of click (UTC)
- attributed\_time : if user download the app for after clicking an ad, this is the time of the app download
- is\_attributed : the target that is to be predicted, indicating the app was downloaded



@电饭煲饭  
weibo.com/u/1936051607



Google

LOGICO

kaggle

CODING  
NINJAS  
www.codingninjas.in

## BIGGEST CHALLENGES

- **Data Size** – More than 175 millions rows which corresponds to 175 million clicks
- It was a **binary classification problem** where number of 1s i.e. positive samples (identified frauds) constitute of **less than 0.5 percent** of the total training data.

## WHAT WORKED FOR ME

- **Magic feature** – Any feature which give more than average boost to the model is mostly termed as magic feature on Kaggle. (In this case, delta of timestamps was the one)
- **Down-sampling** worked wonders for me (Maybe we got lucky, we need that luck sometimes, on Kaggle)



Featured Prediction Competition

kaggle  
DAYS

## Santander Customer Transaction Prediction

Can you identify who will make a transaction?

\$65,000  
Prize Money

MEETUP



Banco Santander · 8,802 teams · 2 months ago

- 200 columns/features
- All column names were redacted as var0, var1, var2, var3, .....
- No context or background of any columns
- More of a **HIT AND TRIAL** at the start
- Hunt for the Magic feature began again



Google

LOGICO

kaggle™

CODING  
NINJAS  
www.codingninjas.in



Featured Prediction Competition

## Santander Customer Transaction Prediction

Can you identify who will make a transaction?

\$65,000  
Prize Money

kaggle  
DAYS  
MEETUP



Banco Santander · 8,802 teams · 2 months ago

### WHAT WORKED FOR ME ?

- *Numerous experiments* – Number of experiments carried out in this competition just to let the data talk was huge
- Peaks in the distribution of few of the variables
- Value counts turned out to be the magic feature



Google

LOGICO

kaggle™

CODING  
NINJAS  
www.codingninjas.in

# CAN YOU LEVERAGE KAGGLE EXPERIENCE INTO YOUR DAY JOB

?



- ✓ YES, Definitely!
- ✓ You get to know about the best practices in your field. For e.g. *how to avoid overfitting*
- ✓ Habit of doing the *feature engineering* on every dataset is very healthy
- ✓ You're updated on latest developments in the field of Data Science
- ✓ You get to learn from the best data scientists across the globe, on how to approach
  - ✓ A classical ML problem
  - ✓ Computer vision problem
  - ✓ NLP problem

## WHAT CHANGES FROM KAGGLE TO DAY JOB WORK?



Google

LOGICO

kaggle™

CODING  
NINJAS  
www.codingninjas.in

- You **never get cleaned & prepared dataset** in day job, the way you get on Kaggle (If you do, tell me your company name :P)
  - One needs to spend major amount of time preparing or cleaning the dataset
- If you look at the bigger picture, you don't have to fight for **minor increase in performance** but you have to fight for it on Kaggle
- On kaggle, you don't always have to have a model which can be explained, whereas business mostly requires a model which has explainability factors attached to it. For e.g. using Shapley values to explain the contribution of each feature in the prediction at local level as well as global level

## WILL KAGGLE HELP YOU GET A JOB?



Google

LOGICO

kaggle™

CODING  
NINJAS  
www.codingninjas.in

- There is a job board on Kaggle too, though no jobs in India were posted on that portal yet, as far as i know
- Kaggle on its own won't help you get a job directly, but the experience, the peer group you get on kaggle is unmatched
- If you happen to win a medal, that's a plus for your resume
- You get to meet people from the community who might guide you to a good opportunity

## AWESOME RESOURCES (FREE)



Google

LOGICO

kaggle™

CODING  
NINJAS  
www.codingninjas.in

CLICK ON THE PICTURES



# Thank you !!

## Q/A Session



Google

LOCIN

kaggle™

CODING  
NINJAS  
www.codingninjas.in