# Kaggle 에서 얻을 수 있는 건?

이유한

카이스트
생명화학공학과
분자 시뮬레이션 실험실
(Molecular Simulation Laboratory)

# Kaggle 이란?

kaggle

## 2010년 설립된
## 빅데이터 솔루션 대회 플랫폼 회사

## 2017년 3월 구글에 인수

# Data Race for 데이터 과학자!

기업, 정부기관, 단체, 연구소, 개인

**Dataset
With Prize**

kaggle

**Dataset & Prize
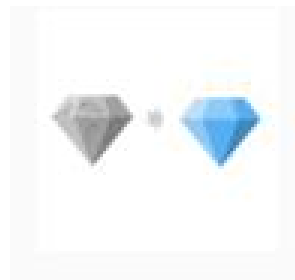개발 환경(kernel)
커뮤니티(follow, discussion)**

전 세계 데이터 사이언티스트

# 참가하려면?

By clicking on the "I understand and accept" button below, you are indicating that you agree to be bound to the competition rules.

I Do Not Accept    **I Understand and Accept**

# Kaggle 에서 competition 을 주최한 단체, 기업들

# Competition Example – Tensorflow competition



**Dataset: 65,000개의
word audio file**


**Prize :
1st - $8,000
2nd – $6,000
3rd – $3,000
+ spectial price $8,000**

Yes, no, up, down, left, right, on, off, stop, go, silence, others 로 이루어진 단어들을 구별하는 AI를 만들어달라!

# Competition Example – Tensorflow competition

## 약 두달간의 RACE
## 1315 팀 참가

# 또 다른 competition 들

**Mercedes-Benz Greener Manufacturing**

Can you cut the time a Mercedes-Benz spends on the test bench?

Featured · 10 months ago · automobiles, tabular data, regression

$25,000

**Quora Question Pairs**

Can you identify question pairs that have the same intent?

Featured · a year ago · linguistics, internet, tabular data, text data, duplicate detection

$25,000

**Passenger Screening Algorithm Challenge**

Improve the accuracy of the Department of Homeland Security's threat recognition algorithms

Featured · 5 months ago · terrorism, image data, object detection

$1,500,000

**Bosch Production Line Performance**

Reduce manufacturing failures

Featured · 2 years ago · manufacturing, tabular data, binary classification

$30,000

KAIST  Molecular Simulation Group
KAIST – Department of Chemical and Biomolecular Engineering

여지껏 다뤄본 것이
IRIS dataset, MNIST 뿐인데

저런 걸 어떻게
분석해야 하나?

# 고수의 발자취를 따라가자

# 모방은 창조의 시작

# 공부해서 함께 나누자! – 캐글 속 선순환

다른 이의 커널
(데이터 분석 보고서)을
공부한다

내 커널을
수정한다

내 커널(분석 보고서)을
만든다

Discussion
참고한다

피드백
받는다

KAIST    Molecular Simulation Group
KAIST – Department of Chemical and Biomolecular Engineering

# My first kaggle story



약 60만명의 정보를 가지고 머신러닝 알고리즘을 만들어, 40만명의 개인이 향후에 보험을 계속 사용할 것인지 예측하라

# My first kaggle story – learning the kernels

# My first kaggle story – making my own kernel

# My first kaggle story – ask anything to authors

# My first kaggle story – 친절한 올리비에 아저씨

# My first kaggle story – ask anything to author

# My first kaggle story – 멋진 스승님

# My first kaggle story

# **My first kaggle story – 범접할 수 없는 은하계 고수**

# My first kaggle story

# My first kaggle story – Get insight from discussion

# My first kaggle story – Submission

# My first kaggle story – After competition



91  2nd place solution NN model
6mo ago

614  1st place with representation learning
Michael Jahrer 5 months ago

106  3rd place solution
utility 6 months ago

94  Solution 1178 Public / 29 Private
CPMP 6 months ago

# Data Playground for 데이터 과학자!

기업, 정부기관, 단체, 연구소, 개인

**Dataset
With or w/o Prize**

kaggle

**Dataset & Prize
개발 환경(kernel)
커뮤니티(follow, discussion)**

전 세계 데이터 사이언티스트

KAIST    Molecular Simulation Group
KAIST - Department of Chemical and Biomolecular Engineering

**25**

# Competition Example – Tensorflow competition



1633  **Credit Card Fraud Detection**
Anonymized credit card transactions labeled as fraudulent or genuine
Machine Learning Group - ULB  updated 2 months ago
crime
finance

1185  **European Soccer Database**
25k+ matches, players & teams attributes for European Professional Football
Hugo Mathien  updated 2 years ago
association football
europe

1119  **TMDB 5000 Movie Dataset**
Metadata on ~5,000 movies from TMDb
The Movie Database (TMDb)  updated 8 months ago
film

870  **Global Terrorism Database**
More than 170,000 terrorist attacks worldwide, 1970-2016
START Consortium  updated 10 months ago
crime
terrorism
international relati...

735  **Bitcoin Historical Data**
Bitcoin data at 1-min intervals from select exchanges, Jan 2012 to March 2018
Zielak  updated a month ago
history
finance

# 2017 Kaggle survey – 캐글러 전공

# 2017 Kaggle survey – 캐글러 현재 직업

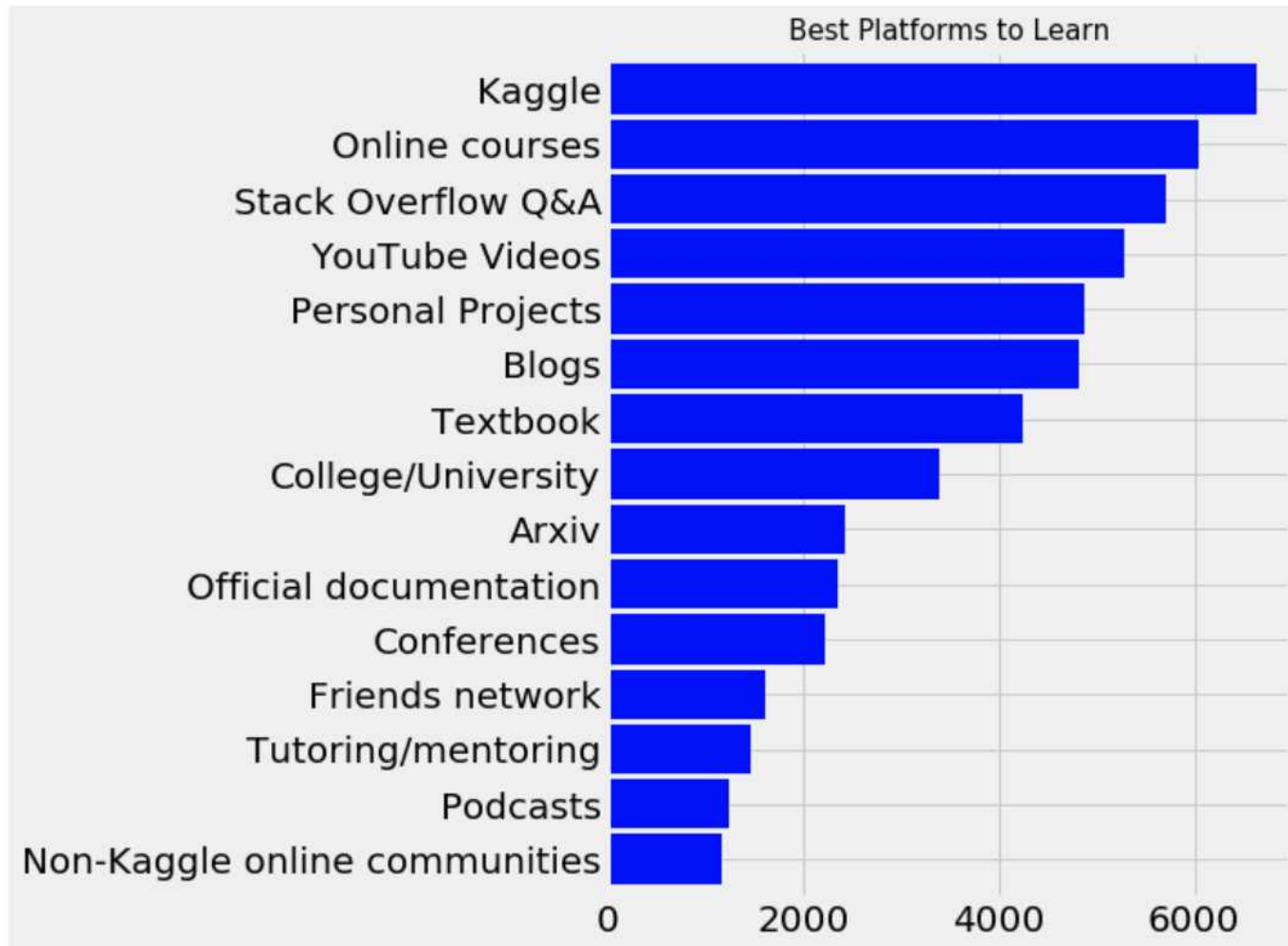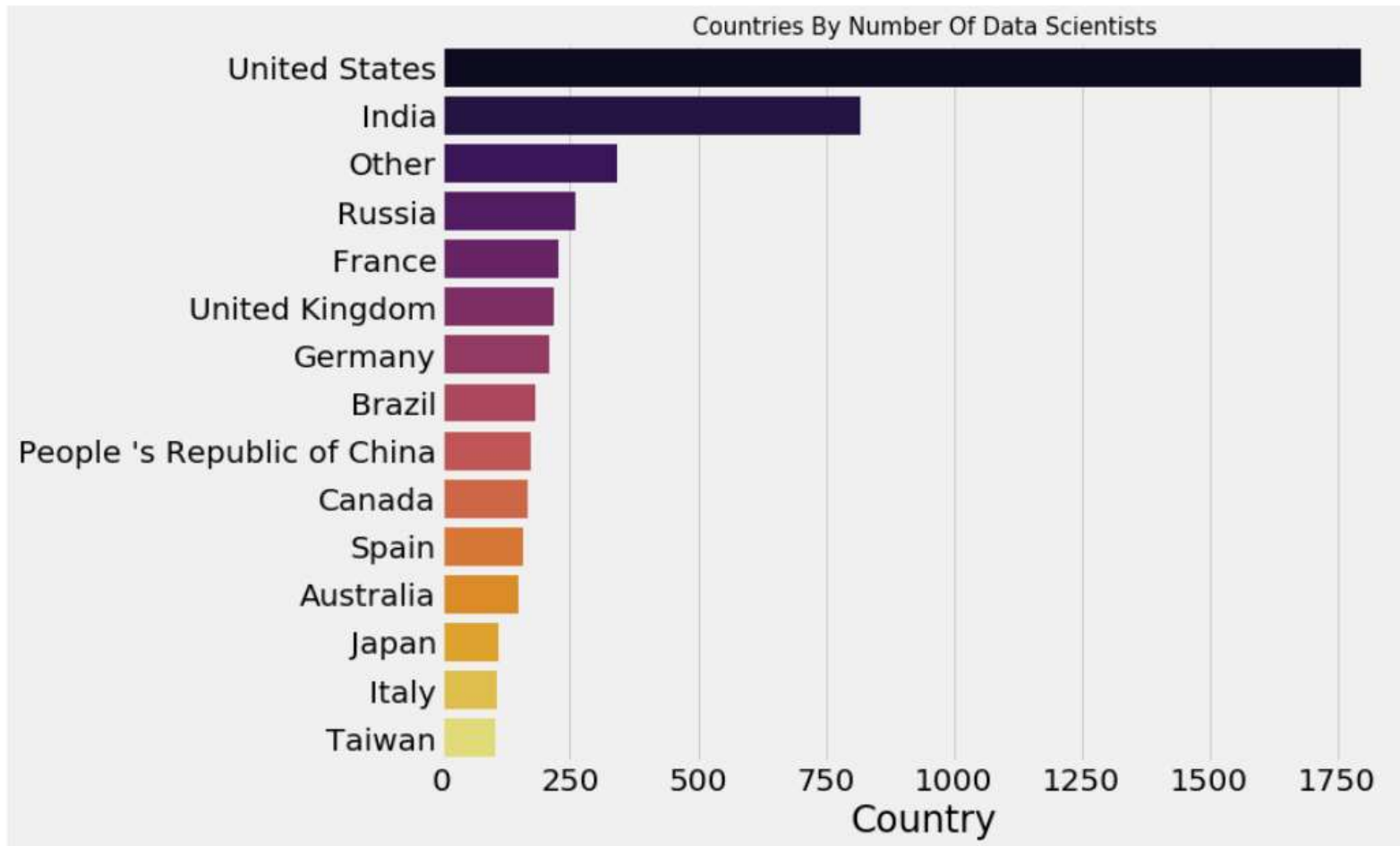# 2017 Kaggle survey – 캐글러는 어디서 배우나?



Best Platforms to Learn

# 2017 Kaggle survey – 캐글러들의 국적?

# 캐글에서 뭘 얻을 수 있나?

1. 현업에 바로 적용 가능한 데이터 분석 기술 습득

2. 현업에 바로 적용 가능한 머신러닝, 딥러닝 알고리즘 적용 방법 습득

3. 인적 네트워크 형성

# 캐글 합시다!